

1 **Viral outbreaks involve destabilized evolutionary**
2 **networks: evidence from Ebola, Influenza and Zika**

3 Stéphane Aris-Brosou,^{1,2,*} Neke Ibeh,¹ and Jessica Noël¹

4 ¹Department of Biology, University of Ottawa, Ottawa, ON K1N 6N5, Canada

5 ²Department of Mathematics and Statistics, University of Ottawa, Ottawa, ON K1N
6 6N5, Canada

7 *Correspondence and requests for materials should be addressed to S.A.B.

8 (sarisbro@uottawa.ca).

9 Abstract

10 Recent history has provided us with one pandemic (Influenza A/H1N1) and two severe
11 viral outbreaks (Ebola and Zika). In all three cases, post-hoc analyses have given us
12 deep insights into what triggered these outbreaks, their timing, evolutionary dynamics,
13 and phylogeography, but the genomic characteristics of outbreak viruses are still unclear.
14 To address this outstanding question, we searched for a common denominator between
15 these recent outbreaks, positing that the genome of outbreak viruses is in an unstable
16 evolutionary state, while that of non-outbreak viruses is stabilized by a network of cor-
17 related substitutions. Here, we show that during regular epidemics, viral genomes are
18 indeed stabilized by a dense network of weakly correlated sites, and that these networks
19 disappear during pandemics and outbreaks when rates of evolution increase transiently.
20 Post-pandemic, these evolutionary networks are progressively re-established. We finally
21 show that destabilization is not caused by substitutions targeting epitopes, but more likely
22 by changes in the environment *sensu lato*. Our results prompt for a new interpretation of
23 pandemics as being associated with evolutionary destabilized viruses.

24 *Keywords:* Ebola virus, Influenza virus, Zika virus, outbreak, pandemic, correlated evo-
25 lution

26 Introduction

27 Over the past few years, humanity has been affected by three major zoonotic events, with
28 an Influenza pandemic in 2009 [1], an Ebola virus outbreak in 2014-16 [2], and a Zika out-
29 break in 2015-16 [3]. In all these examples, the epidemiological and evolutionary dynamics
30 of the pathogens involved, *i.e.*, their phylodynamics [4], were meticulously reconstructed.
31 For instance, in the case of Ebola, an initial phylogenetic study showed evidence that
32 the outbreak originated from a single zoonotic event in an unknown animal reservoir [2],
33 and that the resulting epidemic then spread to the largest and closest neighboring cities
34 following the gravity model [5, 6], with some exceptions [7]. However, in this general
35 context of severe outbreaks, we still do not quite fully understand what characterizes the
36 evolutionary dynamics of the viruses during such events.

37 Recently, in an attempt to understand the genomic determinants of antigenic proper-
38 ties and drug resistance in influenza viruses, we described a novel algorithm to uncover
39 pair of amino acids in a protein that evolve in a correlated manner [8]. We found that
40 influenza A viruses show extensive evidence for correlated evolution, to such an extent
41 that some amino acids evolve correlatively with more than one other site, hereby forming
42 dense (undirected) networks (see also [9]). We furthermore uncovered that some of these
43 pairs of sites are known to be epistatically interacting – specifically, experimental studies
44 show that a mutation at one of these sites lowers viral fitness, which is then restored by a
45 compensatory mutation [10]. Moreover, we showed that similar networks of sites can be
46 found in the Ebola virus, with some of these sites also involved in episodes of adaptive
47 evolution [11]. In light of these results, we here hypothesized that during an outbreak
48 or a pandemic, these networks of tightly correlated sites might be transiently disrupted,
49 hereby leading to a virus that is, from an evolutionary point of view, destabilized.

50 Such a network destabilization would require that some of the intrinsic properties
51 describing these networks change in a similar manner across different viruses. One way of
52 studying these properties is by resorting to the theory used in social networks analysis, and
53 more generally developed in graph theory [12]. In our case, a network is made of nodes,
54 that are amino acid sites in viral proteins, and a link between two amino acids means
55 that these two sites show statistical evidence for evolving in a correlated manner. Both
56 the structure of this network, and the pattern of connections among its nodes, influence
57 its behavior: for instance, scale-free networks, where node connectivity follows a power
58 law, are extremely robust to disruptions [13], just like dense networks [14], while the most
59 connected nodes are also the most important ones in protein-protein interaction networks
60 [15]. Such properties can be derived by summarizing a network with different statistics,
61 such as the number of connections that a particular node has (its *degree*), or the shortest
62 distance between each pair of nodes (the *average path length*).

63 In order to contrast the evolutionary dynamics of pandemic *versus* non-pandemic
64 viruses, we here used these statistics to assess the stability of these networks of amino
65 acids that evolve in a correlated manner. We predicted that viral evolutionary dynamics
66 are weakened during a pandemic. As these dynamics often lead to complex networks
67 of interactions [9, 11], we more specifically tested how the structure of these correlation
68 networks is affected during an outbreak. We show that during a pandemic, the evolution-
69 ary dynamics of viral genes are severely disrupted, but also that they are progressively
70 restored after the pandemic.

71 Results

72 **Networks of correlated sites are destabilized during outbreaks.** In search for
73 evolutionary differences between regular epidemics and severe outbreaks, we first con-
74 trasted the glycoprotein precursor (GP) sequences of the Ebola virus that circulated
75 before and during the 2014/2016 outbreak. For this, we identified with a Bayesian graph-
76 ical model [16] the pairs of nucleotides that show evidence for correlated evolution in each
77 data set, before and during the outbreak. As in previous work [9, 11], we found that
78 these pairs of sites form a network. A first inspection of these networks of correlated sites
79 revealed a striking difference between pre-2014 and outbreak sequences: in particular at
80 weak correlations, the pre-2014 interaction networks are very dense and involve most sites
81 of GP, while only a small number of sites are interacting in outbreak viruses (Figure 1).
82 Furthermore, at increasing correlation strengths, outbreak networks become completely
83 disconnected faster: at posterior probability $Pr = 0.80$ some sites still interact in pre-
84 2014 proteins, while all interactions have disappeared from $Pr = 0.60$ in outbreak proteins
85 (Figure 1). Similar patterns for the Influenza (at two antigens, the hemagglutinin [HA]
86 and the neuraminidase [NA]; Figures S3-S4) and Zika viruses (polymerase NS5; Figures
87 S5) suggest that during a severe outbreak, an evolutionary destabilization of viral genes
88 occurs, especially among sites that entertain weak interactions.

89 **Destabilization affects weakly correlated sites.** To further investigate this desta-
90 bilization hypothesis, we analyzed the structure of these networks with the tools of social
91 network analysis and graph theory [12]. Again, we found a consistent pattern when con-
92 trasting regular and outbreak viruses: at weak to moderate interactions ($Pr \leq 0.50$),
93 outbreak viruses have networks of smaller diameter, shorter path length, and reduced
94 eccentricity (Figure 2a-c, columns 1-5). All these patterns point to fewer connected sites

95 in outbreak viruses. Betweenness is smaller for outbreak viruses (except Ebola), and
96 transitivity tends to be larger (except Zika). These last two measures also suggest that
97 interactions among sites are weakened in outbreak viruses. Other networks statistics failed
98 to show a clear pattern (Figure S6): in particular, there were no clear differences in terms
99 of degree, centrality or homophily – all properties that are not directly related to network
100 stability.

101 **Post-outbreak re-stabilization.** Should these weak interactions play a critical role in
102 the stabilization of viruses outside of pandemics, we would expect to observe the strength-
103 ening of all network statistics as years go by after the pandemic. To test this prediction
104 and estimate how long this re-stabilization process can take, we analyzed in a similar
105 way all influenza seasons in the Northern hemisphere following the 2009 pandemic (until
106 2015-16). Consistent with our prediction, both HA and NA genes show a gradual transi-
107 tion between a typical pandemic state to a regular state in two-to-three seasons (Figure
108 2, column 5-6, respectively).

109 **Non-genetic sources of destabilization.** To understand what the potential sources of
110 this destabilization are, we assessed the involvement of viral antigenic determinants / epi-
111 topes. Should mutations accumulating in such epitopes be responsible for destabilization,
112 we would expect (i) that weak interactions in non-pandemic viruses involve mostly epi-
113 topes, and (ii) that pandemics be associated with the disappearance of these interactions
114 at epitopes first. Figure 3 shows no evidence supporting this hypothesis ($X^2 = 0.0663$,
115 $df = 1$, $P = 0.7967$): non-pandemic viruses show a small number of predicted epitopes
116 in their interaction network, that do not act as central hubs of these networks, while
117 pandemic viruses may actually show an enrichment in interacting epitopes. This suggests
118 that non-genetic factors are likely responsible for the initial destabilization of the genome

119 of pandemic viruses. Changes in their ecology / environment (vector) cannot be ruled
120 out.

121 Discussion

122 To understand how evolutionary dynamics are affected during a viral outbreak, we com-
123 pared non-outbreak and outbreak viruses. Based on the hypothesis that non-outbreak
124 viruses are in a stable evolutionary equilibrium, and that such a stability is mediated by
125 correlated evolution among pairs of sites in viral genes, we reconstructed the coevolution
126 patterns in genes of non-outbreak and outbreak viruses. In line with our prediction, we
127 found that outbreak viruses exhibit fewer coevolving sites than their non-outbreak coun-
128 terparts, and that these interactions are gradually restored after the outbreak, at least in
129 the case of the Influenza (2009 H1N1) virus for both HA and NA.

130 Two independent lines of evidence are consistent with our destabilization hypothesis.
131 First, all three viruses showed temporary increases in their rate of molecular evolution
132 during each outbreak [2, 3, 1]; such increases can be expected to disrupt the coevolution-
133 ary structure, and hence, destabilize viral genomes. We showed that epitopes were not
134 particular targets of this mutational process. This can be expected, as mutations (i) most
135 likely affect sites randomly, and (ii) are quickly lost from the viral population. Second, a
136 probable cause of the epidemics can be identified in all cases studied here. For Influenza,
137 the 2009 pandemic was caused by a series of reassortment events that affected the two
138 genes studied here, HA (triple-reassortant swine) and NA (Eurasian avian-like swine) [1].
139 Such exchanges of segments can very well destabilize the evolutionary dynamics, at least of
140 the implicated segments. Similarly, a zoonotic event was implicated in the Ebola outbreak
141 [2], and a change of continent in the case of Zika [3, 17, 18]. These corresponding changes

142 of environment (*sensu lato*) might have triggered the destabilizations observed here. In
143 addition to such environmental changes, it is very likely that destabilization reflects a
144 complex interaction between the genetics of viruses, their demographic fluctuations and
145 environmental changes.

146 This argument is further supported by recent work in physics, where it was shown
147 that dense networks are more resilient, *i.e.* resistant to small perturbations, than sparser
148 ones [14]. Moreover, in their simplest example, these authors modeled abundances in
149 a community of mutualistic species, where the mutualistic term describes the pairs of
150 interacting species; perturbations were then applied to the system to assess resilience.
151 They showed that small perturbations did not affect average abundances, which remained
152 high – their ‘desirable’ state. However, above a particular perturbation threshold, a
153 bifurcation occurred and a new ‘undesirable’ state, at low abundances, was reached. Our
154 results are consistent with a similar system behavior, where the network of correlated
155 amino acids is resilient to perturbations up to a certain point, when a bifurcation to an
156 ‘undesirable’ state (the pandemic) occurs, and the system returns to its resilient state. One
157 major difference though is that we observed a progressive return to stability in the case of
158 influenza, while the resilience model suggests a second bifurcation, *i.e.* an instantaneous
159 change, to the ‘desirable’ state [14].

160 One outstanding question is about the importance of weak patterns of coevolution
161 within a gene: how can it be explained that it is essentially weak correlations (around
162 $Pr = 0.25$) that distinguish non-outbreak from outbreak viruses? In a recent study on
163 mice, four phenotypes were quantitatively analyzed following large intercrosses, and linear
164 regressions on pairs of quantitative trait loci were used to detect non-additive effects, *i.e.*,
165 epistasis; it was then shown that most epistatic interactions were weak and, critically,
166 tended to stabilize phenotypes towards the mean of the population [19]. Viruses are not

167 mice, and not all the correlations that we detected are involved in epistatic interactions,
168 but both this work in mice and the evidence presented here go in the same direction:
169 weak interactions have a stabilizing effect on viral genes and their phenotype (regular
170 epidemics). It is further possible that the intricate nature of these weak correlation net-
171 works has higher-order effects [19], that in turn increase canalization and hence may help
172 viruses weather modest environmental and genotypic fluctuations [20]. The elimination
173 of these many weak interactions has a destabilizing effect that may be caused by or lead
174 to outbreaks. Our findings call for a new interpretation of pandemics that, from an evo-
175 lutionary point of view, appeared to be associated with unhealthy or diseased viruses.
176 While the evidence shown here does not support the causal nature of this relationship,
177 monitoring correlation networks could help forecast imminent outbreaks.

178 **Methods**

179 **Sequence retrieval.** Nucleotide sequences were retrieved for three viruses: Influenza
180 A, Ebola, and Zika, for select protein-coding genes, chosen because they represent the
181 most sequenced / studied genes for each of these viruses [11, 21, 22, 23]. All sequences
182 were downloaded in May 2016 (Table S1).

183 Full-length Influenza A sequences were retrieved directly from the Influenza Virus
184 Resource [24]. Only H1N1 sequences circulating in humans for the hemagglutinin (HA)
185 and neuraminidase (NA) genes were downloaded. These two genes are also very commonly
186 studied and largely sampled in public databases [22, 23]. Two types of data sets were
187 constructed: one containing pandemic and non-pandemic sequences circulating in 2009,
188 the pandemic year, and one containing pandemic sequences circulating from August 1 to
189 July 31 of each season in the Northern temperate region between 2009/2010 and 2015/2016

190 (seven seasons in total). Only unique sequences were retrieved.

191 For Ebola, the virion spike glycoprotein precursor, GP, was retrieved because of its
192 key role in the emergence of the 2014 outbreak showing evidence for both correlated and
193 adaptive evolution [11] as follows. A GP sequence (KX121421) was drawn at random
194 from the 2014 strain used previously [11] and was employed as a query for a BLASTn
195 search [25] at the National Center for Biotechnology Information. A conservative E -value
196 threshold of 0 ($E < 10^{-500}$) was used, which led to 1,181 accession numbers. As most of
197 these accession numbers correspond to full genomes, while only GP is of interest, we (i)
198 retrieved all corresponding GenBank files, (ii) extracted coding sequences with ReadSeq
199 [26] of all genes, (iii) concatenated the corresponding FASTA files into a single file, (iv)
200 which was then used to format a sequence database for local BLASTn searches, and (v)
201 used GP from KX121421 in a second round of BLASTn searches ($E < 10^{-250}$, coverage
202 $> 75\%$).

203 In the case of Zika, sequences of 252 complete genomes were retrieved from the Virus
204 Pathogen Resource (www.viprbrc.org). The RNA-dependent RNA polymerase NS5 was
205 specifically extracted by performing local BLASTn searches as described above. It is one
206 of the most studied Zika genes [21, 27], as it is essential for the replication of the virus
207 [27].

208 **Phylogenetic analyses.** Sequences were all aligned with Muscle [28] with the fastest
209 options (-maxiters 1 -diags). Alignments were visually inspected with AliView [29] to
210 remove rogue sequences and sequencing errors. Phylogenetic trees were inferred by maxi-
211 mum likelihood under the General Time-Reversible model with among-site rate variation
212 [30] with FastTree [31]. As outbreak sequences (Ebola and Zika viruses) cluster away from
213 non-pandemic sequences, we used the `subtreepplot()` function in APE [32] to retrieve

214 accession numbers of pandemic sequences and hence separate them from non-pandemic
215 sequences with minimal manual input. FastTree was used a second time to estimate
216 phylogenetic trees of the subset alignments, with the same settings as above.

217 **Network analyses of correlated sites.** Amino acid positions (“sites”) that evolve
218 in a correlated manner were identified with the Bayesian graphical model (BGM) in
219 SpiderMonkey [16] as implemented in HyPhy [33]. Briefly, ancestral mutational paths
220 were first reconstructed under the MG94×HKY85 substitution model [34] along each
221 branch of the tree estimated above at non-synonymous sites. These reconstructions were
222 recoded as a binary matrix in which each row corresponds to a branch and each column
223 to a site of the alignment. A BGM was then employed to identify which pairs of sites
224 exhibit correlated patterns of substitutions. Each node of the BGM represents a site and
225 the presence of an edge indicates the conditional dependence between two sites. Such
226 dependence was estimated locally by a posterior probability. Based on the chain rule for
227 Bayesian networks, such local posterior distributions were finally used to estimate the full
228 joint posterior distribution [35]. A maximum of two parents per node was assumed to
229 limit the complexity of the BGM. Posterior distributions were estimated with a Markov
230 chain Monte Carlo sampler that was run for 10^5 steps, with a burn-in period of 10,000
231 steps sampling every 1,000 steps for inference. Analyses were run in duplicate to test for
232 convergence (Figures S1-S2).

233 The estimated BGM can be seen as a weighted network of coevolution among sites,
234 where each posterior probability measures the strength of coevolution. Each probability
235 threshold gives rise to a network whose topology can be analyzed based on a number
236 of measures [12] borrowed from social network analysis and graph theory. We focused in
237 particular on six statistics: average diameter, the length of the longest path between pairs

238 of nodes; average betweenness, measures the importance of each node in their ability to
239 connect to dense subnetworks; assortative degree, measures the extent to which nodes
240 of similar degree are connected to each other (homophily); eccentricity, is the shortest
241 path linking the most distant nodes in the network; average strength, rather than just
242 count the number of connections of each node (degree), strength sums up the weights of
243 all the adjacent nodes; average path length, measures the shortest distance between each
244 pair of nodes. All measures were computed using the igraph R package ver. 1.0.1 [36].
245 Thresholds of posterior probabilities for correlated evolution ranged from 0.01 (weak) to
246 0.99 (strong). LOESS regressions were then fitted to the results.

247 **Epitope analyses.** Epitopes were predicted using the NetCTL 1.2 Server [37]. Briefly,
248 Cytotoxic T lymphocyte (CTL) epitopes are predicted based on a neural network algo-
249 rithm trained on a database of human MHC class I ligands. Epitopes can be predicted
250 for 12 MHC supertypes (A1, A2, A3, A24, A26, B7, B8, B27, B39, B44, B58, B62),
251 that are broad families of very similar peptides for which independent neural network
252 models have been generated. As such, we ran the epitope prediction for each supertype
253 independently, on non-outbreak and outbreak viruses. Circos plots were generated with
254 the circlize R package ver. 0.3.10 [38]. Scripts and sequence alignments used are available
255 from github.com/sarisbro.

256 References

- 257 1. Smith, G. J. D. *et al.* Origins and evolutionary genomics of the 2009 swine-origin
258 H1N1 influenza A epidemic. *Nature* **459**, 1122–5 (2009).
- 259 2. Gire, S. K. *et al.* Genomic surveillance elucidates Ebola virus origin and transmis-
260 sion during the 2014 outbreak. *Science* **345**, 1369–72 (2014).
- 261 3. Faria, N. R. *et al.* Zika virus in the Americas: Early epidemiological and genetic
262 findings. *Science* **352**, 345–9 (2016).

- 263 4. Grenfell, B. T. *et al.* Unifying the epidemiological and evolutionary dynamics of
264 pathogens. *Science* **303**, 327–32 (2004).
- 265 5. Zipf, G. K. The P1 P2/D hypothesis: on the intercity movement of persons.
266 *American sociological review* **11**, 677–686 (1946).
- 267 6. Xia, Y., Bjørnstad, O. N. & Grenfell, B. T. Measles metapopulation dynamics: a
268 gravity model for epidemiological coupling and dynamics. *Am Nat* **164**, 267–81
269 (2004).
- 270 7. Dudas, G. *et al.* Virus genomes reveal factors that spread and sustained the Ebola
271 epidemic. *Nature* **544**, 309–315 (2017).
- 272 8. Nshogozabahizi, J. C., Dench, J. & Aris-Brosou, S. Widespread historical contin-
273 gency in influenza viruses. *Genetics* **205**, 409–420 (2017).
- 274 9. Poon, A. F. Y., Lewis, F. I., Pond, S. L. K. & Frost, S. D. W. An evolutionary-
275 network model reveals stratified interactions in the V3 loop of the HIV-1 envelope.
276 *PLoS Comput Biol* **3**, e231 (2007).
- 277 10. Gong, L. I., Suchard, M. A. & Bloom, J. D. Stability-mediated epistasis constrains
278 the evolution of an influenza protein. *Elife* **2**, e00631 (2013).
- 279 11. Ibeh, N., Nshogozabahizi, J. C. & Aris-Brosou, S. Both epistasis and diversifying
280 selection drive the structural evolution of the Ebola virus glycoprotein mucin-like
281 domain. *J Virol* **90**, 5475–84 (2016).
- 282 12. Newman, M. *Networks: an introduction* (OUP Oxford, 2010).
- 283 13. Albert, Jeong & Barabasi. Error and attack tolerance of complex networks. *Nature*
284 **406**, 378–82 (2000).
- 285 14. Gao, J., Barzel, B. & Barabási, A.-L. Universal resilience patterns in complex
286 networks. *Nature* **530**, 307–12 (2016).
- 287 15. Jeong, H., Mason, S. P., Barabási, A. L. & Oltvai, Z. N. Lethality and centrality
288 in protein networks. *Nature* **411**, 41–2 (2001).
- 289 16. Poon, A. F. Y., Lewis, F. I., Frost, S. D. W. & Kosakovsky Pond, S. L. Spi-
290 dermonkey: rapid detection of co-evolving sites using bayesian graphical models.
291 *Bioinformatics* **24**, 1949–50 (2008).
- 292 17. Zhang, Q. *et al.* Spread of Zika virus in the Americas. *Proc Natl Acad Sci U S A*
293 **114**, E4334–E4343 (2017).

- 294 18. Faria, N. R. *et al.* Establishment and cryptic transmission of zika virus in brazil
295 and the americas. *Nature* **546**, 406–410 (2017).
- 296 19. Tyler, A. L., Donahue, L. R., Churchill, G. A. & Carter, G. W. Weak epistasis
297 generally stabilizes phenotypes in a mouse intercross. *PLoS Genet* **12**, e1005805
298 (2016).
- 299 20. Waddington, C. H. Canalization of development and the inheritance of acquired
300 characters. *Nature* **150**, 563–565 (1942).
- 301 21. Faye, O. *et al.* Molecular evolution of Zika virus during its emergence in the 20(th)
302 century. *PLoS Negl Trop Dis* **8**, e2636 (2014).
- 303 22. Aris-Brosou, S. Inferring influenza global transmission networks without complete
304 phylogenetic information. *Evol Appl* **7**, 403–12 (2014).
- 305 23. Labonté, K. & Aris-Brosou, S. Automatic detection of rate change in large data
306 sets with an unsupervised approach: the case of influenza viruses. *Genome* **59**,
307 253–62 (2016).
- 308 24. Bao, Y. *et al.* The influenza virus resource at the National Center for Biotechnology
309 Information. *J Virol* **82**, 596–601 (2008).
- 310 25. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local
311 alignment search tool. *J Mol Biol* **215**, 403–10 (1990).
- 312 26. Gilbert, D. Sequence file format conversion with command-line readseq. *Curr*
313 *Protoc Bioinformatics* **Appendix 1**, Appendix 1E (2003).
- 314 27. Zhao, B. *et al.* Structure and function of the Zika virus full-length NS5 protein.
315 *Nat Commun* **8**, 14762 (2017).
- 316 28. Edgar, R. C. Muscle: multiple sequence alignment with high accuracy and high
317 throughput. *Nucleic Acids Res* **32**, 1792–7 (2004).
- 318 29. Larsson, A. AliView: a fast and lightweight alignment viewer and editor for large
319 datasets. *Bioinformatics* **30**, 3276–8 (2014).
- 320 30. Aris-Brosou, S. & Rodrigue, N. The essentials of computational molecular evolu-
321 tion. *Methods Mol Biol* **855**, 111–52 (2012).
- 322 31. Price, M. N., Dehal, P. S. & Arkin, A. P. Fasttree 2—approximately maximum-
323 likelihood trees for large alignments. *PLoS One* **5**, e9490 (2010).
- 324 32. Paradis, E., Claude, J. & Strimmer, K. APE: analyses of phylogenetics and evolu-
325 tion in R language. *Bioinformatics* **20**, 289–290 (2004).

- 326 33. Pond, S. L. K., Frost, S. D. W. & Muse, S. V. Hyphy: hypothesis testing using
327 phylogenies. *Bioinformatics* **21**, 676–9 (2005).
- 328 34. Kosakovsky Pond, S. L. & Frost, S. D. W. Not so different after all: a comparison of
329 methods for detecting amino acid sites under selection. *Mol Biol Evol* **22**, 1208–22
330 (2005).
- 331 35. Pearl, J. *Probabilistic reasoning in intelligent systems: networks of plausible infer-*
332 *ence* (Morgan Kaufmann, 1988).
- 333 36. Csardi, G. & Nepusz, T. The **igraph** software package for complex network re-
334 search. *InterJournal, Complex Systems* **1695**, 1–9 (2006).
- 335 37. Larsen, M. V. *et al.* Large-scale validation of methods for cytotoxic T-lymphocyte
336 epitope prediction. *BMC bioinformatics* **8**, 1 (2007).
- 337 38. Gu, Z., Gu, L., Eils, R., Schlesner, M. & Brors, B. **circIzize** implements and
338 enhances circular visualization in R. *Bioinformatics* btu393 (2014).

339 Acknowledgements

340 We thank Jonathan Dench and George S. Long for discussions and comments, as well as
341 two anonymous reviewers for additional comments. This work was supported by the Nat-
342 ural Sciences Research Council of Canada and by the Canada Foundation for Innovation
343 (S.A.B.) and by the University of Ottawa (N.I., J.N.).

344 Author information

345 Affiliations

346 **Department of Biology, University of Ottawa, Ottawa, ON K1N 6N5, Canada**

347 Stéphane Aris-Brosou, Neke Ibeh & Jessica Noël

348 **Department of Mathematics and Statistics, University of Ottawa, Ottawa, ON**
349 **K1N 6N5, Canada**

350 Stéphane Aris-Brosou

351 Contributions

352 S.A.B. designed the study, and wrote the paper. S.A.B., N.I. and J.N. performed re-
353 search and analyses, and edited the paper. All authors approved the final version of the
354 manuscript.

355 **Competing interests**

356 The authors declare that they have no competing interests.

357 **Corresponding author**

358 Stéphane Aris-Brosou (sarisbro@uottawa.ca).

359 **Supplementary information**

360 1. Supplemental Information

361 Figures

Figure 1. Correlation network of pre-outbreak and outbreak Ebola viruses. Networks of correlated sites in the GP protein are shown in each panel. The top row shows networks for the viruses circulating before the 2014 outbreak (blue); the bottom row shows networks for outbreak viruses (red). Each column shows networks for different strengths of correlation, from weak ($Pr = 0.05$) to strong ($Pr = 0.95$). Nodes represent amino acid sites, and edges correlations. Node sizes are proportional to diameter.

Figure 2. Network properties between pandemic and non-pandemic viruses. Results are shown for Ebola (column 1), Zika (2) and Influenza viruses: for HA and NA circulating in 2009 in (3) and (4), respectively, and for pandemic viruses circulating between the 2009-10 (deep red) and the 2015-16 (deep blue) season in (5) and (6). Pandemic viruses are shown in red, while non-pandemic ones are in blue. Shading: 95% confidence envelopes of the LOESS regressions. Five network measures are shown: (a) diameter, (b) average path length, (c) eccentricity, (d) betweenness, and (e) transitivity.

Figure 3. Interacting residues in pandemic and non-pandemic viruses. Results are shown for Ebola at weak correlations ($Pr = 0.20$). Coevolving positions in the alignment are identified with arabic numbers; for those that are predicted to be epitopes, supertypes (A1, A2, *etc.*) are shown.