

1 **Host shifts result in parallel genetic changes when viruses evolve in closely related species**

2

3 **Ben Longdon^{*1,2}, Jonathan P Day², Joel M Alves^{2,3}, Sophia CL Smith², Thomas M Houslay¹, John E**
4 **McGonigle², Lucia Tagliaferri² and Francis M Jiggins²**

5

6 ¹ Biosciences, College of Life & Environmental Sciences, University of Exeter, Penryn Campus, TR10
7 9FE, UK

8 ² Department of Genetics, University of Cambridge, Cambridge, CB2 3EH, UK

9 ³ CIBIO, Centro de Investigação em Biodiversidade e Recursos Genéticos, InBIO Laboratório

10 Associado, Universidade do Porto, 4485-661 Vairão, Portugal

11

12 *b.longdon2@exeter.ac.uk

13

14

15 **Abstract**

16

17 Host shifts, where a pathogen invades and establishes in a new host species, are a major source of
18 emerging infectious diseases. They frequently occur between related host species and often rely on
19 the pathogen evolving adaptations that increase their fitness in the novel host species. To
20 investigate genetic changes in novel hosts, we experimentally evolved replicate lineages of an RNA
21 virus (*Drosophila C Virus*) in 19 different species of *Drosophilidae* and deep sequenced the viral
22 genomes. We found a strong pattern of parallel evolution, where viral lineages from the same host
23 were genetically more similar to each other than to lineages from other host species. When we
24 compared viruses that had evolved in different host species, we found that parallel genetic changes
25 were more likely to occur if the two host species were closely related. This suggests that when a
26 virus adapts to one host it might also become better adapted to closely related host species. This
27 may explain in part why host shifts tend to occur between related species, and may mean that when
28 a new pathogen appears in a given species, closely related species may become vulnerable to the
29 new disease.

30 Introduction

31

32 Host shifts – where a pathogen jumps into and establishes in a new host species – are a major
33 source of emerging infectious diseases. RNA viruses seem particularly prone to host shift [1-4], with
34 HIV, Ebola virus and SARS coronavirus all having been acquired by humans from other host species
35 [5-7]. Whilst some pathogens may be pre-adapted to a novel host, there are increasing numbers of
36 examples demonstrating that adaptation to the new host occurs following a host shift [8, 9]. These
37 adaptations may allow a pathogen to enter host cells, increase replication rates, avoid or suppress
38 the host immune response, or optimise virulence or transmission [10, 11]. For example, in the 2013-
39 2016 Ebola virus epidemic in West Africa, a mutation in the viral glycoprotein gene that arose early
40 in the outbreak and rose to high frequency was found to increase infectivity in human cells and
41 decrease infectivity in bats, which are thought to be the source of Ebola virus [12, 13]. Likewise, a
42 switch of a parvovirus from cats to dogs resulted in mutations in the virus capsid that allowed the
43 virus to bind to cell receptors in dogs, but resulted in the virus losing its ability to infect cats [14, 15]

44

45 In some instances adaptation to a novel host relies on specific mutations that arise repeatedly
46 whenever a pathogen switches to a given host. For example, in the jump of HIV-1 from chimps to
47 humans, codon 30 of the *gag* gene has undergone a change that increases virus replication in
48 humans, and this has occurred independently in all three HIV-1 lineages [5, 16]. Similarly, five
49 parallel mutations have been observed in the two independent epidemics of SARS coronavirus
50 following its jump from palm civets into humans [17]. Similar patterns have been seen in
51 experimental evolution studies, where parallel genetic changes occur repeatedly when replicate viral
52 lineages adapt to a new host species in the lab. For example, when Vesicular Stomatitis Virus was
53 passaged in human or dog cells, the virus evolved parallel mutations when evolved on the same cell
54 type [18]. Likewise, a study passaging Tobacco Etch Potyvirus on four plant species found parallel
55 mutations occurred only when the virus infected the same host species [19]. These parallel
56 mutations provide compelling evidence that these genetic changes are adaptive, with the same
57 mutations evolving independently in response to natural selection [20].

58

59 The host phylogeny is important for determining a pathogen's ability to infect a novel host, with
60 pathogens tending to replicate most efficiently when they infect a novel host that is closely related
61 to their original host [2, 21-34]. Here, we asked whether viruses acquire the same genetic changes
62 when evolving in the same and closely related host species. We experimentally evolved replicate
63 lineages of an RNA virus called Drosophila C Virus (DCV; Discistroviridae) in 19 species of
64 Drosophilidae that vary in their relatedness and shared a common ancestor approximately 40 million
65 years ago [35, 36]. We then sequenced the genomes of the evolved viral lineages and tested
66 whether the same genetic changes arose when the virus was evolved in closely related host species.

67

68

69 Results

70

71 Parallel genetic changes occur in DCV lineages that have evolved in the same host species

72

73 To examine how viruses evolve in different host species we serially passaged DCV in 19 species of
74 Drosophilidae. In total we infected 22,095 adult flies and generated 173 independent replicate

75 lineages (6-10 per host species). We deep sequenced the evolved virus genomes to generate over
76 740,000 300bp sequence reads from each viral lineage. Out of 8989 sites, 584 contained a SNP with
77 a derived allele frequency >0.05 in at least one viral lineage, and 84 of these were tri-allelic. None of
78 these variants were found at an appreciable frequency in five sequencing libraries produced from
79 the ancestral virus, indicating that they had spread through populations during the experiment
80 (Figure 1). In multiple cases these variants had nearly reached fixation (Figure 1).

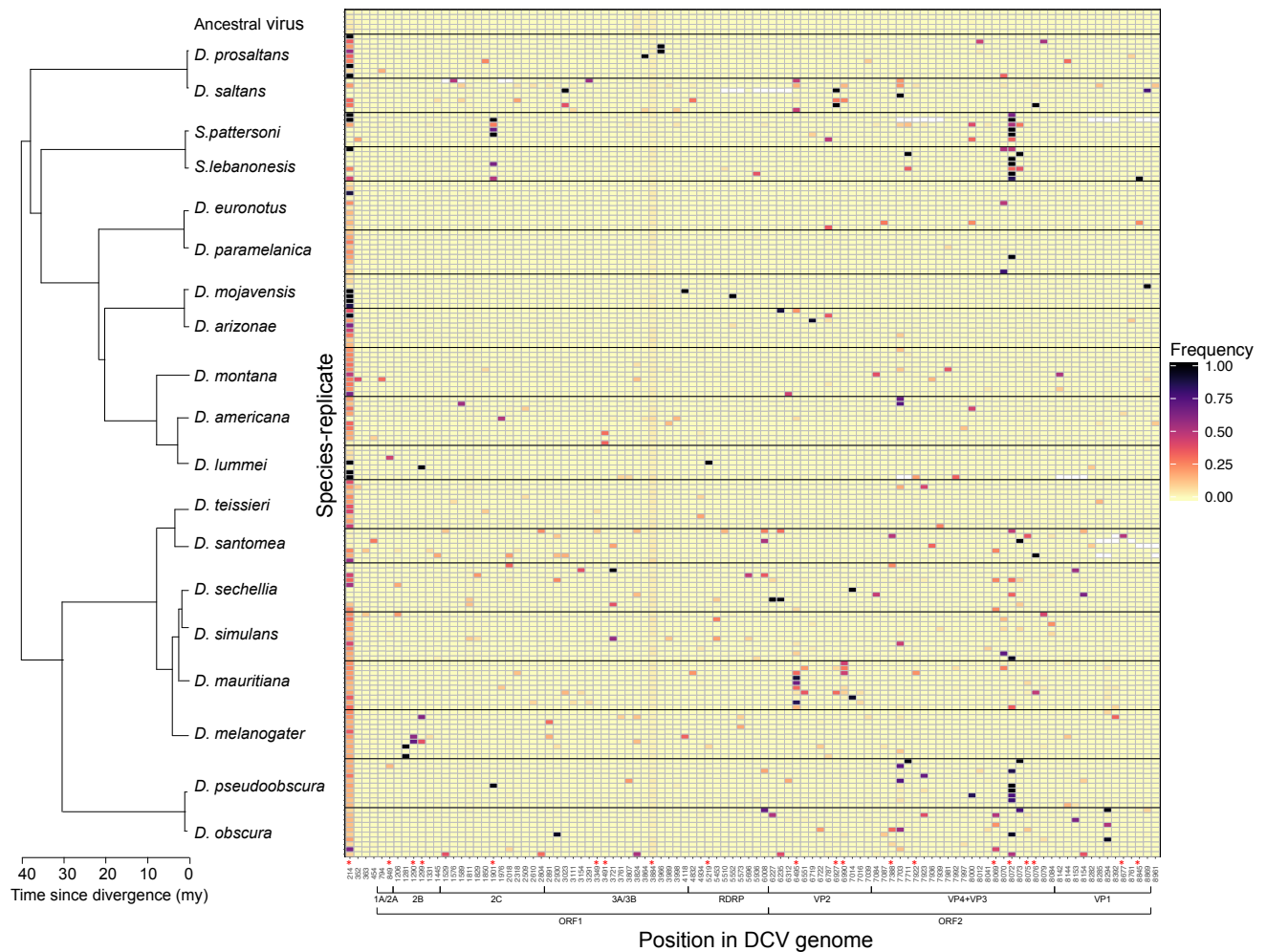
81

82 We next examined whether the same genetic changes occur in parallel when different populations
83 encounter the same host species. Of the 584 SNPs, 102 had derived allele frequencies >0.05 in at
84 least two viral lineages, and some had risen to high frequencies in multiple lineages (Figure 1). We
85 estimated the genetic differentiation between viral lineages by calculating F_{ST} . We found that viral
86 lineages that had evolved within the same host were genetically more similar to each other than to
87 lineages from other host species (Figure 2; $P<0.001$). Furthermore, we found no evidence of
88 differences in substitution biases in the different host species (Fisher Exact Test: $p=0.14$; see
89 methods), suggesting that this pattern is not driven by changes in the types of mutations in different
90 host species.

91

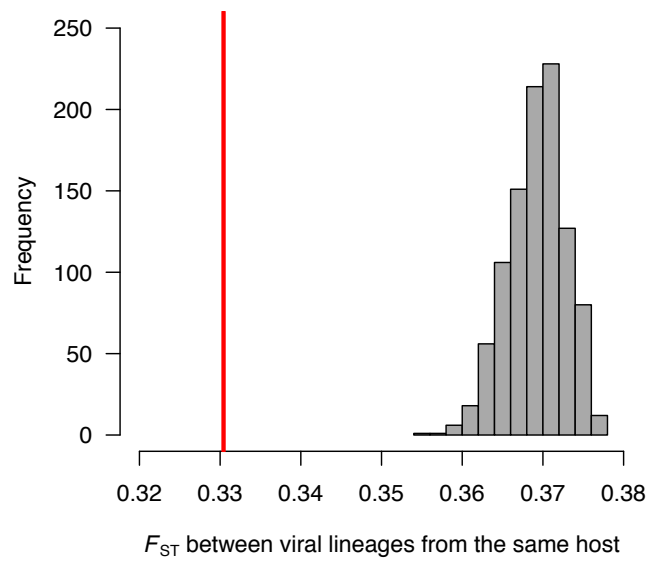
92 To examine the genetic basis of parallel evolution, we individually tested whether each SNP in the
93 DCV genome showed a signature of parallel evolution among viral lineages passaged in the same
94 host species (i.e. we repeated the analysis in Figures 2 for each SNP). We identified 56 polymorphic
95 sites with a significant signal of parallel evolution within the same host species ($P<0.05$; significantly
96 parallel sites are shown with a red asterisk in Figure 1; the false discovery rate is estimated to be
97 17% [37]).

98



99

100 **Figure 1. The frequency of SNPs in viral lineages that have evolved in different host species.** Each
 101 row represents an independent viral lineage. Viruses that evolved in different host species are
 102 separated by black horizontal lines. Each column represents a polymorphic site in the DCV genome,
 103 and only sites where the derived allele frequency >0.05 in at least two lineages are shown. The
 104 intensity of shading represents the derived allele frequency. Sites where there are three alleles have
 105 the two derived allele frequencies pooled for illustrative purposes. Sites with SNP frequencies that
 106 are significantly correlated among lineages from the same host species are shown by red stars at the
 107 bottom the column (permutation test; $p < 0.05$). Open reading frames (ORFs) and viral proteins based
 108 on predicted polyprotein cleavage sites [38-42] are below the x axis. Information on the distribution
 109 of mutations across the genome and whether they are synonymous or non-synonymous can be
 110 found in the supplementary results. Sites with missing data are shown in white. The phylogeny was
 111 inferred under a relaxed molecular clock [33, 43] and the scale axis represents the approximate age
 112 since divergence in millions of years (my) based on estimates from: [35, 36].



113

114 **Figure 2. Viral lineages from the same host species were genetically more similar to each other**
115 **than to lineages from different host species.** The mean pairwise F_{ST} between all possible pairs of
116 viral lineages from the same host species was calculated. The red line shows the observed value. The
117 grey bars are the null distribution of this statistic obtained by permuting the viral lineages across
118 host species 1000 times.

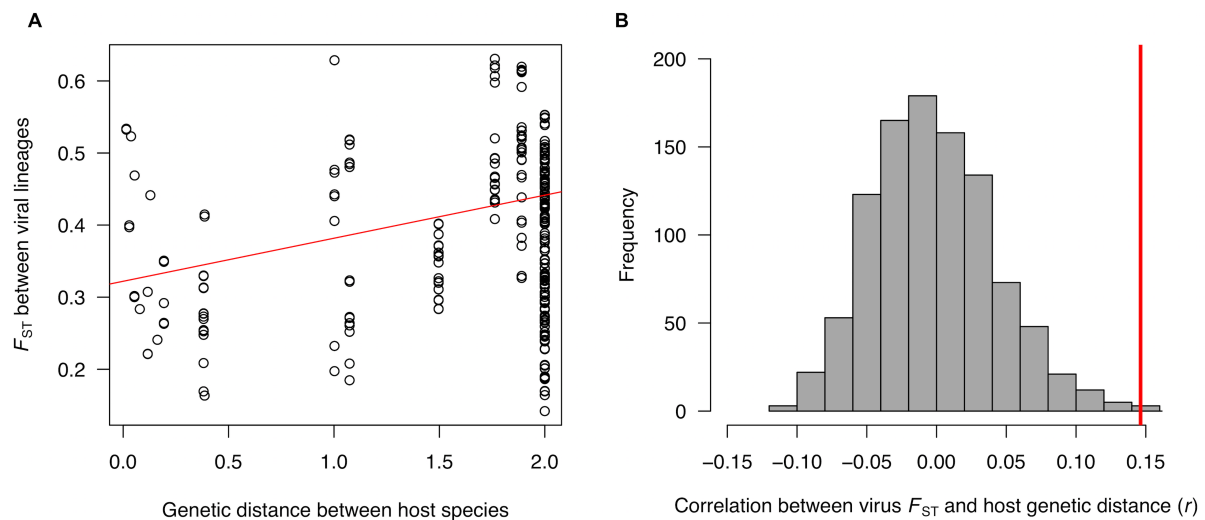
119

120 **Viruses in closely related hosts are genetically more similar**

121

122 We investigated if viruses passaged through closely related hosts showed evidence of parallel
123 genetic changes. We calculated F_{ST} between all possible pairs of viral lineages that had evolved in
124 different host species. We found that viral lineages from closely related hosts were more similar to
125 each other than viral lineages from more distantly related hosts (Figure 3A). This is reflected in a
126 significant positive relationship between virus F_{ST} and host genetic distance (Figure 3B, Permutation
127 test: $r=0.15$, $P=0.002$). We lacked the statistical power to identify the specific SNPs that are causing
128 the signature of parallel evolution in Figure 3 (false discovery rate >0.49 for all SNPs).

129



130

131

132 **Figure 3. Viral lineages from more closely related host species are genetically more similar.** (A) The
133 correlation between the genetic differentiation of viral lineages and the genetic distance between
134 the species they have evolved in. Linear regression line is shown in red. Genetic distances were
135 scaled so that the distance from the root to the tip of the tree was one. (B) Pearson's correlation
136 coefficient (r) of F_{ST} between pairs of viral lineage and the genetic distance between the host species
137 they evolved in. The observed value is in red and the grey bars are the null distribution obtained by
138 permutation.

138

139

140 Discussion

141

142 When a pathogen infects a novel host species, it finds itself in a new environment to which it must
143 adapt [4, 8, 10, 44]. When DCV was passaged through different species of Drosophilidae, we found
144 the same genetic changes arose repeatedly in replicate viral lineages in the same host species. Such
145 repeatable parallel genetic changes to the same host environment are compelling evidence that
146 these changes are adaptive [20]. We then examined whether these same genetic changes might
147 occur in closely related host species, as these are likely to present a similar environment for the
148 virus. We found that viruses evolved in closely related hosts were more similar to each other than
149 viruses that evolved in more distantly related species. Therefore, mutations that evolve in one host
150 species frequently arise when the virus infects closely related hosts. This finding of parallel genetic
151 changes in closely related host species suggests that when a virus adapts to one host it might also
152 become better adapted to closely related host species.

153

154 Phylogenetic patterns of host adaptation may in part explain why pathogens tend to be more likely
155 to jump between closely related host species. This pattern is seen in nature, where host shifts tend
156 to occur most frequently between closely related hosts, and in laboratory cross-infection studies,
157 where viruses tend to replicate more rapidly when the new host is related to the pathogens natural
158 host [2, 21-34]. For example, in a large cross-infection experiment involving *Drosophila sigma* viruses
159 (Rhabdoviridae) isolated from different species of *Drosophila*, the viruses tended to replicate most
160 efficiently in species closely related to their natural hosts [34]. This suggests that these viruses had

161 acquired adaptations to their host species that benefitted them when they infected closely related
162 species. Our results demonstrate that this pattern is apparent at the level of specific nucleotides,
163 and can arise very shortly after a host shift. The function of these mutations is unknown, but in other
164 systems adaptations after host shifts have been found to enhance the ability of the virus to bind to
165 host receptors [11], increase replication rates [16] or avoid the host immune response [8, 10, 45].
166

167 While the susceptibility of a novel host is correlated to its relatedness to the pathogens' original
168 host, it is also common to find exceptions to this pattern. This is seen both in nature when
169 pathogens shift between very distant hosts [46, 47], and in laboratory cross-infection experiments
170 [33, 34]. This pattern is also seen in our data where we also observe parallel genetic changes
171 occurring between more distantly related hosts. For example, a mutation at position 8072 was not
172 only near fixation in most of the lineages infecting two closely related species, but also occurred at a
173 high frequency in replicate lineages in a phylogenetically distant host (Figure 1).
174

175 In conclusion, we have found that host relatedness can be important in determining how viruses
176 evolve when they find themselves in a new host. This study suggests that while some genetic
177 changes will be found only in specific hosts, we frequently see the same changes occurring in closely
178 related host species. These phylogenetic patterns suggest that mutations that adapt a virus to one
179 host may also adapt it to closely related host species. Therefore, there may be a knock-on effect,
180 where a host shift leaves closely related species vulnerable to the new disease.
181

182

183

183 **Methods**

184

185 *Virus production*

186

187 DCV is a positive sense RNA virus in the family Discistroviridae that was isolated from *D.*
188 *melanogaster*, which it naturally infects in the wild [48, 49]. To minimise the amount of genetic
189 variation in the DCV isolate we used to initiate the experimental evolution study, we aimed to isolate
190 single infectious clones of DCV using a serial dilution procedure. DCV was produced in Schneider's
191 Drosophila line 2 (DL2) cells [50] as described in [51]. Cells were cultured at 25°C in Schneider's
192 Drosophila Medium with 10% Fetal Bovine Serum, 100 U/ml penicillin and 100 µg/ml streptomycin
193 (all Invitrogen, UK). The DCV strain used was isolated from *D. melanogaster* collected in Charolles,
194 France [52]. DL2 cells were seeded into two 96-well tissue culture plates at approximately 10⁴ cells in
195 100 µl of media per well. Cells were allowed to adhere to the plates by incubating at 25°C for five
196 hours or over-night. Serial 1:1 dilutions of DCV were made in complete Schneider's media, giving a
197 range of final dilutions from 1:10⁸ – 1:4x10¹⁴. 100 µl of these dilutions were then added to the cells
198 and incubated for 7 days, 8 replicates were made for each DCV dilution. Each well was then
199 examined for DCV infection of the DL2 cells, and a well was scored as positive for DCV infection if
200 clear cytopathic effects were present in the majority of the cells. The media was taken from the
201 wells with the greatest dilution factor that were scored as infected with DCV and stored at -80°C.
202 This processes was then repeated using the DCV samples from the first dilution series. One clone,
203 B6A, was selected for amplification and grown in cell culture as described above. Media containing
204 DCV was removed and centrifuged at 3000 x g for 5 minutes at 4°C to pellet any remaining cell

205 debris, before being aliquoted and stored at -80°C . The Tissue Culture Infective Dose 50 (TCID₅₀) of
206 the DCV was 6.32×10^9 infectious particles per ml using the Reed-Muench end-point method [53].

207

208 *Inoculating fly species*

209

210 We passaged the virus through 19 species of Drosophilidae, with 6-10 independent replicate
211 passages for each species. We selected species from across the phylogeny (that shared a common
212 ancestor approximately 40 million years ago [35, 36]), but included clades of closely related species
213 that recently shared common ancestors less than 5 million years ago (Figure 1). All fly stocks were
214 reared at 22°C . Stocks of each fly species were kept in 250ml bottles at staggered ages. Flies were
215 collected and sexed, and males were placed on cornmeal medium for 4 days before inoculation.
216 Details of the fly stocks used can be found in the supplementary materials.

217

218 4-11 day old males were infected with DCV using a 0.0125 mm diameter stainless steel needle
219 (26002–10, Fine Science Tools, CA, USA) dipped in DCV solution. For the first passage this was the
220 cloned DCV isolate in cell culture supernatant (described above), and then subsequently was the
221 virus extracted from the previous passage (described below). The needle was pricked into the
222 pleural suture on the thorax of flies, towards the midcoxa. Each replicate was infected using a new
223 needle and strict general cleaning procedures were used to minimise any risk of cross-contamination
224 between replicates. Species were collected and inoculated in a randomised order each passage. Flies
225 were then placed into vials of cornmeal medium and kept at 22°C and 70% relative humidity. Flies
226 were snap frozen in liquid nitrogen 3 days post-infection, homogenised in Ringer's solution (2.5 μl
227 per fly) and then centrifuged at 12,000g for 10 mins at 4°C . The resulting supernatant was removed
228 and frozen at -80°C to be used for infecting flies in the subsequent passage. The remaining
229 homogenate was preserved in Trizol reagent (Invitrogen) and stored at -80°C for RNA extraction.
230 The 3 day viral incubation period was chosen based on time course and pilot data showing that viral
231 load reaches a maximum at approximately 3 days post-infection. This process was repeated for 10
232 passages for all species, except *D. montana* where only 8 passages were carried out due to the fly
233 stocks failing to reproduce. Each lineage was injected into a mean of 11 flies at each passage (range
234 4-18). Experimental evolution studies in different tissue types have seen clear signals of adaptation
235 in 100 virus generations [18]. Based on \log_2 change in RNA viral load we estimate that we have
236 passaged DCV for approximately 100-200 generations.

237

238 *Sequencing*

239

240 After passaging the virus, we sequenced evolved viral lineages from 19 host species, with a mean of
241 9 independent replicate lineages of the virus per species (range 6-10 replicates). cDNA was
242 synthesised using Invitrogen Superscript III reverse-transcriptase with random hexamer primers
243 (25°C 5mins, 50°C 50mins, 70°C 15mins). The genome of the evolved viruses, along with the initial
244 DCV ancestor (x5) were then amplified using Q5 high fidelity polymerase (NEB) in nine overlapping
245 PCR reactions (see supplementary Table S2 for PCR primers and cycle conditions). Primers covered
246 position 62-9050bp (8989bp) of the Genbank refseq (NC_001834.1) giving 97% coverage of the
247 genome. PCRs of individual genomes were pooled and purified with Ampure XP beads (Agencourt).
248 Individual Nextera XT libraries (Illumina) were prepared for each viral lineage. In total we sequenced

249 173 DCV pooled amplicon libraries on an Illumina MiSeq (Cambridge Genomic Service) v3 for 600
250 cycles to give 300bp paired-end reads.

251

252 *Bioinformatics and variant calling*

253 FastQC, version 0.11.2 [54] was used to assess read quality and primer contamination. Trimmomatic,
254 version 0.32 [55] was used to removed low quality bases and adaptor sequences, using the following
255 options: MINLEN=30 (Drop the read if it is below 30 base pairs), TRAILING=15 (cut bases of the end
256 of the read if below a threshold quality of 15), SLIDINGWINDOW=4:20 (perform a sliding window
257 trimming, cutting once the average quality within a 4bp window falls below a threshold of 20), and
258 ILLUMINACLIP=TruSeq3-PE.fa:2:20:10:1:true (remove adapter contamination; the values correspond
259 in order to: input fasta file with adapter sequences to be matched, seed mismatches, palindrome clip
260 threshold, simple clip threshold, minimum adapter length and logical value to keep both reads in
261 case of read-through being detected in paired reads by palindrome mode).

262 To generate a reference ancestral Drosophila C Virus sequence we amplified the ancestral starting
263 virus by PCR as above. PCR products were treated with exonuclease 1 and Antarctic phosphatase to
264 remove unused PCR primers and dNTPs and then sequenced directly using BigDye reagents (ABI) on
265 an ABI 3730 capillary sequencer in both directions (Source Bioscience, Cambridge, UK). Sequences
266 were edited in Sequencher (version 4.8; Gene Codes), and were manually checked for errors. Fastq
267 reads were independently aligned to this reference sequence (Genbank accession: MG570143) using
268 BWA-MEM, version 0.7.10 {Li, 2009 #1605} with default options with exception of the parameter –
269 M, which marks shorter split hits as secondary. 99.5% of reads had mapping phred quality scores of
270 >60. The generated SAM files were converted to their binary format (BAM) and sorted by their
271 leftmost coordinates with SAMtools, version 0.1.19 (website: <http://samtools.sourceforge.net/>) [56].
272 Read Group information (RG) was added to the BAM files using the module
273 AddOrReplaceReadGroups from Picard Tools, version 1.126 (<https://broadinstitute.github.io/picard>).

274 The variant calling was then performed for each individual BAM using UnifiedGenotyper tool from
275 GATK, version 3.3.0. As we were interested in calling low frequency variants in our viruses, we
276 assumed a ploidy level of 100 (-sample_ploidy:100). The other parameters were set to their defaults
277 except --stand_call_conf:30 (minimum phred-scaled confidence threshold at which variants should
278 be called) and --downsample_to_coverage:1000 (down-sample each sample to 1000X coverage)

279 *Host phylogeny*

280 We used a trimmed version of a phylogeny produced previously [33]. This time-based tree (where
281 the distance from the root to the tip is equal for all taxa) was inferred using seven genes with a
282 relaxed molecular clock model in BEAST (v1.8.0) [43, 57]. The tree was pruned to the 19 species used
283 using the Ape package in R [58, 59].

284 *Statistical Analysis*

285

286 We examined the frequency of alternate alleles (single nucleotide polymorphisms: SNPs) in five
287 ancestral virus replicates (aliquots of the same virus stock that was used to found the evolved
288 lineages). SNPs in these ancestral viruses may represent pre-standing genetic variation, or may be

289 sequencing errors. We found the mean SNP frequency was 0.000923 and the highest frequency of
290 any SNP was 0.043 across the ancestral viruses. We therefore included a SNP in our analyses if its
291 frequency was >0.05 in any of the evolved viral lineages. For all analyses we included all three alleles
292 at triallelic sites.

293

294 *Parallel evolution within species*

295

296 As a measure of genetic differentiation we estimated F_{ST} between all the virus lineages based on the
297 heterozygosity (H) of the SNPs we called [60]:

$$298 \quad F_{ST} = \frac{H_b - H_w}{H_b} \quad (\text{Equation 1})$$

299 where H_b is the mean number of differences between pairs of sequence reads sampled from the two
300 different lineages. H_w is mean number of differences between sequence reads sampled from within
301 each lineage. H_b and H_w were calculated separately for each polymorphic site, and the mean across
302 sites used in equation (1). H_w was calculated separately for the two lineages being compared, and
303 the unweighted mean used in equation (1).

304

305 To examine whether there had been parallel evolution among viral lineages that had evolved within
306 the same fly species, we calculated the mean F_{ST} between lineages that had evolved in the same fly
307 species, and compared this to the mean F_{ST} between lineages that had evolved in different fly
308 species. We tested whether this difference was statistically significant using a permutation test. The
309 fly species labels were randomly reassigned to the viral lineages, and we calculated the mean F_{ST}
310 between lineages that had evolved in the same fly species. This was repeated 1000 times to
311 generate a null distribution of the test statistic, and this was then compared to the observed value.

312

313 To identify individual SNPs with a signature of parallel evolution within species, we repeated this
314 procedure separately for each SNP.

315

316 *Parallel evolution between species*

317

318 We next examined whether viral lineages that had evolved in different fly species tended to be more
319 similar if the fly species were more closely related. Considering all pairs of viral lineages from
320 different host species, we correlated pairwise F_{ST} with the genetic distance between the fly species.

321 To test the significance of this correlation, we permuted the fly species over the *Drosophila*
322 phylogeny and recalculated the Pearson correlation coefficient. This was repeated 1000 times to
323 generate a null distribution of the test statistic, and this was then compared to the observed value.

324 To identify individual SNPs whose frequencies were correlated with the genetic distance between
325 hosts we repeated this procedure separately for each SNP.

326

327 We confirmed there was no relationship between rates of molecular evolution (SNP frequency) and
328 either genetic distance from the host DCV was isolated from (*D. melanogaster*) or estimated viral
329 population size (see supplementary Figures S1 and S2) using generalised linear mixed models that
330 include the phylogeny as a random effect in the MCMCglmm package in R [61] as described

331 previously [34]. We also examined the distribution of SNPs and whether they were synonymous or
332 non-synonymous (see supplementary results).

333

334 To test whether there were systematic differences in the types of mutations occurring in the
335 different host species, we classified all the SNPs into the six possible types (A/G, A/T, A/C, G/T, G/C
336 and C/T). We then counted the number of times each type of SNP arose in each host species at a
337 frequency above 5% and in at least one biological replicate (SNPs in multiple biological replicates
338 were only counted once). This resulted in a contingency table with 6 columns and 19 rows. We
339 tested for differences between the species in the relative frequency of the 6 SNP types by simulation
340 [62].

341

342 Sequence data (fastq files) are available in the NCBI SRA (Accession: SRP119720). BAM files, data and
343 R scripts for analysis in the main text are available from the NERC data repository (funding
344 requirement - awaiting doi, temporary link to data and scripts
345 <https://figshare.com/s/b119ba86def8bca58782>).

346

347 **Author contributions**

348 Designing experiment: BL, FMJ. Lab work: JPD, SCLS, TMH, LT, BL. Bioinformatics analysis: JMA, JEM,
349 FMJ, BL. Statistical analysis: BL, FMJ. Manuscript written by BL and FMJ with input from all authors.

350

351 **Acknowledgements**

352 Thanks to the Drosophila species stock centre for providing fly stocks and four anonymous reviewers
353 for constructive comments.

354

355 **Funding**

356

357 BL and FMJ are supported by a Natural Environment Research Council (NE/L004232/1
358 <http://www.nerc.ac.uk/>) and by an European Research Council grant (281668, DrosophilaInfection,
359 <http://erc.europa.eu/>). JMA was supported by a grant from the Portuguese Ministério da Ciência,
360 Tecnologia e Ensino Superior (SFRH/BD/72381/2010). BL is supported by a Sir Henry Dale Fellowship
361 jointly funded by the Wellcome Trust and the Royal Society (Grant Number 109356/Z/15/Z).

362

363 **References**

364

- 365 1. Cleaveland S, Laurenson MK, Taylor LH. Diseases of humans and their domestic mammals:
366 pathogen characteristics, host range and the risk of emergence. *Philosophical Transactions of the*
367 *Royal Society of London Series B-Biological Sciences*. 2001;356(1411):991-9. PubMed PMID:
368 WOS:000170315900003.
- 369 2. Davies TJ, Pedersen AB. Phylogeny and geography predict pathogen community similarity in
370 wild primates and humans. *Proceedings of the Royal Society B-Biological Sciences*.
371 2008;275(1643):1695-701. doi: 10.1098/rspb.2008.0284. PubMed PMID: ISI:000256387500014.
- 372 3. Taylor LH, Latham SM, Woolhouse ME. Risk factors for human disease emergence. *Philos*
373 *Trans R Soc Lond B Biol Sci*. 2001;356(1411):983-9. Epub 2001/08/23. doi: 10.1098/rstb.2001.0888.
374 PubMed PMID: 11516376; PubMed Central PMCID: PMC1088493.
- 375 4. Woolhouse ME, Haydon DT, Antia R. Emerging pathogens: the epidemiology and evolution
376 of species jumps. *Trends Ecol Evol*. 2005;20(5):238-44. Epub 2006/05/17. doi: S0169-5347(05)00038-
377 8 [pii]

- 378 10.1016/j.tree.2005.02.009. PubMed PMID: 16701375.
- 379 5. Sharp PM, Hahn BH. The evolution of HIV-1 and the origin of AIDS. *Philosophical*
380 *Transactions of the Royal Society B-Biological Sciences*. 2010;365(1552):2487-94. doi:
381 10.1098/rstb.2010.0031. PubMed PMID: WOS:000280097000008.
- 382 6. Leroy EM, Kumulungui B, Pourrut X, Rouquet P, Hassanin A, Yaba P, et al. Fruit bats as
383 reservoirs of Ebola virus. *Nature*. 2005;438(7068):575-6. doi: 10.1038/438575a. PubMed PMID:
384 WOS:000233593100030.
- 385 7. Li WD, Shi ZL, Yu M, Ren WZ, Smith C, Epstein JH, et al. Bats are natural reservoirs of SARS-
386 like coronaviruses. *Science*. 2005;310(5748):676-9. doi: 10.1126/science.1118391. PubMed PMID:
387 WOS:000232997700042.
- 388 8. Parrish CR, Holmes EC, Morens DM, Park EC, Burke DS, Calisher CH, et al. Cross-species virus
389 transmission and the emergence of new epidemic diseases. *Microbiology and Molecular Biology*
390 *Reviews*. 2008;72(3):457-70. doi: 10.1128/mmbr.00004-08. PubMed PMID: WOS:000258951200004.
- 391 9. Russell CA, Fonville JM, Brown AE, Burke DF, Smith DL, James SL, et al. The potential for
392 respiratory droplet-transmissible A/H5N1 influenza virus to evolve in a mammalian host. *Science*.
393 2012;336(6088):1541-7. Epub 2012/06/23. doi: 10.1126/science.1222526. PubMed PMID:
394 22723414; PubMed Central PMCID: PMC3426314.
- 395 10. Longdon B, Brockhurst MA, Russell CA, Welch JJ, Jiggins FM. The Evolution and Genetics of
396 Virus Host Shifts. *PLoS Pathog*. 2014;10(11):e1004395. Epub 2014/11/07. doi:
397 10.1371/journal.ppat.1004395. PubMed PMID: 25375777; PubMed Central PMCID: PMC4223060.
- 398 11. Parrish CR, Kawaoka Y. The origins of new pandemic viruses: the acquisition of new host
399 ranges by canine parvovirus and influenza A viruses. *Annual review of microbiology*. 2005;59:553-86.
400 Epub 2005/09/13. doi: 10.1146/annurev.micro.59.030804.121059. PubMed PMID: 16153179.
- 401 12. Diehl WE, Lin AE, Grubaugh ND, Carvalho LM, Kim K, Kyaw PP, et al. Ebola Virus
402 Glycoprotein with Increased Infectivity Dominated the 2013-2016 Epidemic. *Cell*. 2016;167(4):1088-
403 98 e6. doi: 10.1016/j.cell.2016.10.014. PubMed PMID: 27814506; PubMed Central PMCID:
404 PMCPMC5115602.
- 405 13. Urbanowicz RA, McClure CP, Sakuntabhai A, Sall AA, Kobinger G, Muller MA, et al. Human
406 Adaptation of Ebola Virus during the West African Outbreak. *Cell*. 2016;167(4):1079-87 e5. doi:
407 10.1016/j.cell.2016.10.013. PubMed PMID: 27814505; PubMed Central PMCID: PMCPMC5101188.
- 408 14. Shackelton LA, Parrish CR, Truyen U, Holmes EC. High rate of viral evolution associated with
409 the emergence of carnivore parvovirus. *Proc Natl Acad Sci U S A*. 2005;102(2):379-84. Epub
410 2005/01/01. doi: 10.1073/pnas.0406765102. PubMed PMID: 15626758; PubMed Central PMCID:
411 PMCPMC544290.
- 412 15. Truyen U, Evermann JF, Vieler E, Parrish CR. Evolution of canine parvovirus involved loss and
413 gain of feline host range. *Virology*. 1996;215(2):186-9. Epub 1996/01/15. doi:
414 10.1006/viro.1996.0021. PubMed PMID: 8560765.
- 415 16. Wain LV, Bailes E, Bibollet-Ruche F, Decker JM, Keele BF, Van Heuverswyn F, et al.
416 Adaptation of HIV-1 to its human host. *Mol Biol Evol*. 2007;24(8):1853-60. Epub 2007/06/05. doi:
417 10.1093/molbev/msm110. PubMed PMID: 17545188.
- 418 17. Liu W, Tang F, Fontanet A, Zhan L, Wang TB, Zhang PH, et al. Molecular epidemiology of
419 SARS-associated coronavirus, Beijing. *Emerg Infect Dis*. 2005;11(9):1420-4. Epub 2005/10/19. doi:
420 10.3201/eid1109.040773. PubMed PMID: 16229772; PubMed Central PMCID: PMC3310602.
- 421 18. Remold SK, Rambaut A, Turner PE. Evolutionary genomics of host adaptation in vesicular
422 stomatitis virus. *Mol Biol Evol*. 2008;25(6):1138-47. Epub 2008/03/21. doi: 10.1093/molbev/msn059.
423 PubMed PMID: 18353798.
- 424 19. Bedhomme S, Lafforgue G, Elena SF. Multihost experimental evolution of a plant RNA virus
425 reveals local adaptation and host-specific mutations. *Mol Biol Evol*. 2012;29(5):1481-92. Epub
426 2012/02/10. doi: 10.1093/molbev/msr314. PubMed PMID: 22319146.

- 427 20. Bollback JP, Huelsenbeck JP. Parallel Genetic Evolution Within and Between Bacteriophage
428 Species of Varying Degrees of Divergence. *Genetics*. 2009;181(1):225-34. doi:
429 10.1534/genetics.107.085225. PubMed PMID: WOS:000262595500021.
- 430 21. Streicker DG, Turmelle AS, Vonhof MJ, Kuzmin IV, McCracken GF, Rupprecht CE. Host
431 Phylogeny Constrains Cross-Species Emergence and Establishment of Rabies Virus in Bats. *Science*.
432 2010;329(5992):676-9. doi: 10.1126/science.1188836. PubMed PMID: WOS:000280602700037.
- 433 22. Faria NR, Suchard MA, Rambaut A, Streicker DG, Lemey P. Simultaneously reconstructing
434 viral cross-species transmission history and identifying the underlying constraints. *Philos Trans R Soc
435 Lond B Biol Sci*. 2013;368(1614):20120196. Epub 2013/02/06. doi: 10.1098/rstb.2012.0196. PubMed
436 PMID: 23382420; PubMed Central PMCID: PMC3678322.
- 437 23. Cooper N, Griffin R, Franz M, Omotayo M, Nunn CL, Fryxell J. Phylogenetic host specificity
438 and understanding parasite sharing in primates. *Ecol Lett*. 2012;15(12):1370-7. Epub 2012/08/24.
439 doi: 10.1111/j.1461-0248.2012.01858.x. PubMed PMID: 22913776.
- 440 24. Waxman D, Weinert LA, Welch JJ. Inferring host range dynamics from comparative data: the
441 protozoan parasites of new world monkeys. *Am Nat*. 2014;184(1):65-74. Epub 2014/06/13. doi:
442 10.1086/676589. PubMed PMID: 24921601.
- 443 25. Huang S, Bininda-Emonds ORP, Stephens PR, Gittleman JL, Altizer S. Phylogenetically related
444 and ecologically similar carnivores harbor similar parasite assemblages. *Journal of Animal Ecology*.
445 2013;n/a-n/a. doi: 10.1111/1365-2656.12160.
- 446 26. Hadfield JD, Krasnov BR, Poulin R, Nakagawa S. A Tale of Two Phylogenies: Comparative
447 Analyses of Ecological Interactions. *The American Naturalist*. 2014;0(0):000. doi: 10.1086/674445.
- 448 27. Ramsden C, Holmes EC, Charleston MA. Hantavirus evolution in relation to its rodent and
449 insectivore hosts: no evidence for codivergence. *Mol Biol Evol*. 2009;26(1):143-53. Epub 2008/10/17.
450 doi: msn234 [pii]
451 10.1093/molbev/msn234. PubMed PMID: 18922760.
- 452 28. de Vienne DM, Hood ME, Giraud T. Phylogenetic determinants of potential host shifts in
453 fungal pathogens. *Journal of Evolutionary Biology*. 2009;22(12):2532-41. doi: 10.1111/j.1420-
454 9101.2009.01878.x. PubMed PMID: WOS:000271785800019.
- 455 29. Gilbert GS, Webb CO. Phylogenetic signal in plant pathogen-host range. *Proceedings of the
456 National Academy of Sciences of the United States of America*. 2007;104(12):4979-83. doi:
457 10.1073/pnas.0607968104. PubMed PMID: WOS:000245256700040.
- 458 30. Tinsley MC, Majerus MEN. Small steps or giant leaps for male-killers? Phylogenetic
459 constraints to male-killer host shifts. *Bmc Evolutionary Biology*. 2007;7. doi: 10.1186/1471-2148-7-
460 238. PubMed PMID: WOS:000252786000001.
- 461 31. Russell JA, Goldman-Huertas B, Moreau CS, Baldo L, Stahlhut JK, Werren JH, et al.
462 Specialization and geographic isolation among *Wolbachia* symbionts from ants and lycaenid
463 butterflies. *Evolution*. 2009;63(3):624-40. Epub 2008/12/05. doi: 10.1111/j.1558-5646.2008.00579.x.
464 PubMed PMID: 19054050.
- 465 32. Perlman SJ, Jaenike J. Infection success in novel hosts: An experimental and phylogenetic
466 study of *Drosophila*-parasitic nematodes. *Evolution*. 2003;57(3):544-57. PubMed PMID:
467 WOS:000182193800010.
- 468 33. Longdon B, Hadfield JD, Day JP, Smith SC, McGonigle JE, Cogni R, et al. The Causes and
469 Consequences of Changes in Virulence following Pathogen Host Shifts. *PLoS Pathog*.
470 2015;11(3):e1004728. Epub 2015/03/17. doi: 10.1371/journal.ppat.1004728. PubMed PMID:
471 25774803.
- 472 34. Longdon B, Hadfield JD, Webster CL, Obbard DJ, Jiggins FM. Host phylogeny determines viral
473 persistence and replication in novel hosts. *PLoS Pathogens*. 2011;7((9)):e1002260. doi:
474 10.1371/journal.ppat.1002260.
- 475 35. Obbard DJ, Maclennan J, Kim K-W, Rambaut A, O'Grady PM, Jiggins FM. Estimating
476 divergence dates and substitution rates in the *Drosophila* phylogeny. *Molecular Biology and
477 Evolution*. 2012;29(11):3459-73.

- 478 36. Tamura K, Subramanian S, Kumar S. Temporal patterns of fruit fly (*Drosophila*) evolution
479 revealed by mutation clocks. *Mol Biol Evol.* 2004;21(1):36-44. Epub 2003/09/02. doi:
480 10.1093/molbev/msg236
481 msg236 [pii]. PubMed PMID: 12949132.
- 482 37. Storey JD. A direct approach to false discovery rates. *J Roy Stat Soc B.* 2002;64:479-98. doi:
483 Unsp 1369-7412/02/64479
484 Doi 10.1111/1467-9868.00346. PubMed PMID: WOS:000177425500009.
- 485 38. Jan E. Divergent IRES elements in invertebrates. *Virus Res.* 2006;119(1):16-28. doi:
486 10.1016/j.virusres.2005.10.011. PubMed PMID: 16307820.
- 487 39. Johnson KN, Christian PD. The novel genome organization of the insect picorna-like virus
488 *Drosophila C virus* suggests this virus belongs to a previously undescribed virus family. *J Gen Virol.*
489 1998;79 (Pt 1):191-203. Epub 1998/02/14. PubMed PMID: 9460942.
- 490 40. Nakashima N, Nakamura Y. Cleavage sites of the "P3 region" in the nonstructural polyprotein
491 precursor of a dicistrovirus. *Arch Virol.* 2008;153(10):1955-60. doi: 10.1007/s00705-008-0208-5.
492 PubMed PMID: 18810573.
- 493 41. Nakashima N, Uchiumi T. Functional analysis of structural motifs in dicistroviruses. *Virus Res.*
494 2009;139(2):137-47. doi: 10.1016/j.virusres.2008.06.006. PubMed PMID: 18621089.
- 495 42. UniProtKB [cited 2017]. Available from: <http://www.uniprot.org/uniprot/O36966>.
- 496 43. Drummond AJ, Suchard MA, Xie D, Rambaut A. Bayesian phylogenetics with BEAUti and the
497 BEAST 1.7. *Mol Biol Evol.* 2012;29(8):1969-73. Epub 2012/03/01. doi: 10.1093/molbev/mss075.
498 PubMed PMID: 22367748; PubMed Central PMCID: PMC3408070.
- 499 44. Woolhouse ME, Gowtage-Sequeria S. Host range and emerging and reemerging pathogens.
500 *Emerg Infect Dis.* 2005;11(12):1842-7. Epub 2006/02/21. PubMed PMID: 16485468.
- 501 45. Sauter D, Schindler M, Specht A, Landford WN, Munch J, Kim KA, et al. Tetherin-Driven
502 Adaptation of Vpu and Nef Function and the Evolution of Pandemic and Nonpandemic HIV-1 Strains.
503 *Cell Host & Microbe.* 2009;6(5):409-21. doi: Doi 10.1016/J.Chom.2009.10.004. PubMed PMID:
504 ISI:000272539700006.
- 505 46. Webby RJ, Webster RG. Emergence of influenza A viruses. *Philos Trans R Soc Lond B Biol Sci.*
506 2001;356(1416):1817-28. Epub 2002/01/10. doi: 10.1098/rstb.2001.0997. PubMed PMID: 11779380;
507 PubMed Central PMCID: PMC1088557.
- 508 47. Weinert LA, Welch JJ, Suchard MA, Lemey P, Rambaut A, Fitzgerald JR. Molecular dating of
509 human-to-bovid host jumps by *Staphylococcus aureus* reveals an association with the spread of
510 domestication. *Biol Lett.* 2012;8(5):829-32. Epub 2012/05/26. doi: 10.1098/rsbl.2012.0290. PubMed
511 PMID: 22628096; PubMed Central PMCID: PMC3440972.
- 512 48. Christian PD. Studies of *Drosophila C* and *A* viruses in Australian populations of *Drosophila*
513 *melanogaster*: Australian National University; 1987.
- 514 49. Webster CL, Waldron FM, Robertson S, Crowson D, Ferrai G, Quintana JF, et al. The
515 discovery, distribution and evolution of viruses associated with *Drosophila melanogaster*. *PLOS*
516 *Biology.* 2015;13(7): e1002210.
- 517 50. Teixeira L, Ferreira A, Ashburner M. The Bacterial Symbiont *Wolbachia* Induces Resistance to
518 RNA Viral Infections in *Drosophila melanogaster*. *Plos Biology.* 2008;6(12):2753-63. doi: DOI
519 10.1371/journal.pbio.1000002. PubMed PMID: ISI:000261913700016.
- 520 51. Longdon B, Cao C, Martinez J, Jiggins FM. Previous Exposure to an RNA Virus Does Not
521 Protect against Subsequent Infection in *Drosophila melanogaster*. *Plos One.* 2013;8(9):e73833. doi:
522 10.1371/journal.pone.0073833.
- 523 52. Jousset FX, Plus N, Croizier G, Thomas M. [Existence in *Drosophila* of 2 groups of picornavirus
524 with different biological and serological properties]. *C R Acad Sci Hebd Seances Acad Sci D.*
525 1972;275(25):3043-6. Epub 1972/12/18. PubMed PMID: 4631976.
- 526 53. Reed LJ, Muench H. A simple method of estimating fifty per cent endpoints. *The American*
527 *Journal of Hygiene.* 1938;27:493-7.

- 528 54. Andrews S. FastQC: a quality control tool for high throughput sequence data. Available at:
529 <http://www.bioinformatics.bbsrc.ac.uk/projects/fastqc.2010>.
- 530 55. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data.
531 Bioinformatics. 2014;30(15):2114-20. doi: 10.1093/bioinformatics/btu170. PubMed PMID:
532 24695404; PubMed Central PMCID: PMC4103590.
- 533 56. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence
534 Alignment/Map format and SAMtools. Bioinformatics. 2009;25(16):2078-9. doi:
535 10.1093/bioinformatics/btp352. PubMed PMID: 19505943; PubMed Central PMCID:
536 PMC2723002.
- 537 57. Drummond AJ, Rambaut A. BEAST: Bayesian evolutionary analysis by sampling trees. BMC
538 Evolutionary Biology. 2007;7:214. doi: Artn 214
539 Doi 10.1186/1471-2148-7-214. PubMed PMID: ISI:000253468300001.
- 540 58. Team RDC. R: a language and environment for statistical computing. V 2.4. 2006.
- 541 59. Paradis E, Claude J, Strimmer K. APE: Analyses of Phylogenetics and Evolution in R language.
542 Bioinformatics. 2004;20(2):289-90. PubMed PMID: 14734327.
- 543 60. Hudson RR, Slatkin M, Maddison WP. Estimation of levels of gene flow from DNA sequence
544 data. Genetics. 1992;132(2):583-9. Epub 1992/10/01. PubMed PMID: 1427045; PubMed Central
545 PMCID: PMC1205159.
- 546 61. Hadfield JD. MCMC Methods for Multi-Response Generalized Linear Mixed Models: The
547 MCMCglmm R Package. Journal of Statistical Software. 2010;33(2):1-22. PubMed PMID:
548 WOS:000275203300001.
- 549 62. Patefield WM. Algorithm AS 159: An Efficient Method of Generating Random R × C
550 Tables with Given Row and Column Totals. Journal of the Royal Statistical Society Series C (Applied
551 Statistics). 1981;30(1):91-7. doi: 10.2307/2346669.
552