1   # Recent evolutionary origin and localized diversity hotspots of
2   # mammalian coronaviruses.

3

4   Renan Maestri[a,b,1], Benoît Perez-Lamarque[a,c,1], Anna Zhukova[d], Hélène Morlon[a]

5   [a]Institut de Biologie de l'École Normale Supérieure (IBENS), École Normale Supérieure, CNRS,
6   INSERM, Université PSL, Paris, France

7   [b]Departamento de Ecologia, Instituto de Biociências, Universidade Federal do Rio Grande do Sul,
8   Porto Alegre, RS, Brazil

9   [c]Institut de Systématique, Évolution, Biodiversité (ISYEB), Muséum national d'histoire naturelle,
10  CNRS, Sorbonne Université, EPHE, UA, Paris, France

11  [d]Institut Pasteur, Université Paris Cité, Bioinformatics and Biostatistics Hub, Paris, France

12  [1]R.M. and B.P.-L. contributed equally to this work.

13  *Corresponding author: Renan Maestri

14  **Email:** renanmaestri@gmail.com

15  **Author Contributions:** R.M. and H.M. conceived the study; R.M, B.P-L., A.Z. and H.M. designed
16  and performed the research; B.P-L. and R.M. analyzed the data; R.M, B.P-L., A.Z. and H.M.
17  wrote the paper.

18  **Competing Interest Statement:** The authors declare no competing interest.

19  **Classification:** Biological Sciences. Evolution.

20  **Keywords:** Paste the keywords here. There should be at least three and no more than five.

21  coronavirus evolution | diversity of coronaviruses | preferential host switches | parasite
22  diversification

23

24 **Abstract**

25 Several coronaviruses infect humans, with three, including the SARS-CoV2, causing diseases.
26 While coronaviruses are especially prone to induce pandemics, we know little about their
27 evolutionary history, host-to-host transmissions, and biogeography, which impedes the prediction
28 of future transmission scenarios. One of the difficulties lies in dating the origination of the family, a
29 particularly challenging task for RNA viruses in general. Previous cophylogenetic tests of virus-
30 host associations, including in the Coronaviridae family, have suggested a virus-host
31 codiversification history stretching many millions of years. Here, we establish a framework for
32 robustly testing scenarios of ancient origination and codiversification *versus* recent origination
33 and diversification by host switches. Applied to coronaviruses and their mammalian hosts, our
34 results support a scenario of recent origination of coronaviruses in bats and diversification by host
35 switches, with preferential host switches within mammalian orders. Hotspots of coronavirus
36 diversity, concentrated in East Asia and Europe, are consistent with this scenario of relatively
37 recent origination and localized host switches. Spillovers from bats to other species are rare, but
38 have the highest probability to be towards humans than to any other mammal species, implicating
39 humans as the evolutionary intermediate host. The high host-switching rates within orders, as
40 well as between humans, domesticated mammals, and non-flying wild mammals, indicates the
41 potential for rapid additional spreading of coronaviruses across the world. Our results suggest
42 that the evolutionary history of extant mammalian coronaviruses is recent, and that cases of long-
43 term virus–host codiversification have been largely over-estimated.

44
45
46 **Main Text**
47

48 **Introduction**
49
50 Coronaviruses are RNA-viruses of the family Coronaviridae, comprising positive-sense and
51 single-stranded viruses that have the largest genomes among nidoviruses (1, 2). As several other
52 RNA viruses, they may cause diseases in humans and other animals (3). Depending on the
53 taxonomic arrangement, seven (4–6) or eight (7) species of coronaviruses infect humans, three of
54 which being pathogenic: the SARS-CoV (8, 9), the MERS-CoV (10), and the SARS-CoV2 (11).
55 The latter is at the origin of the recent COVID-19 pandemic that infected more than 620 million
56 people and caused the death of more than six and a half million (12). Coronaviruses' high
57 frequency of recombination (13), broad host range, and high mutation rates (7) make them
58 especially prone to causing yet future diseases. Nevertheless, their evolutionary history and
59 biogeography are very poorly understood. Resolving the evolutionary origins of Coronaviridae,
60 understanding how they diversified, and characterizing their geographic diversity patterns would
61 facilitate attempts to predict future zoonoses (7, 14, 15).
62
63 Coronaviruses infect mammals, birds, and fish (2), although they predominate in
64 mammalian species (16–20). A consensus exists on the taxonomic segregation of four genera
65 within Coronaviridae: Orthocoronavirinae, namely Alpha-, Beta-, Gamma- and Deltacoronavirus
66 (2, 21). Alpha- and Betacoronaviruses are found exclusively in mammals, while Delta- and
67 Gammacoronaviruses infect mostly birds but also mammals to a lesser extent (20, 22, 23).
68 Coronaviruses are most numerous and genetically diversified in mammals (2, 23), in particular
69 bats, suggesting a mammalian origin in bats (2, 20, 23, 24), although this remains to be tested.
70
71 The timing of origination of the Coronaviridae family is debated, with results that vary by
72 several orders of magnitude. Woo et al (23) found a recent origin, around 10 thousand years ago.
73 This dating was obtained by sequencing the well-conserved RNA-dependent RNA-polymerase
74 (RdRp) genome region of representatives of all four coronavirus genera, and fitting to these

2

75  sequences a neutral nucleotide-based substitution model with an uncorrelated log-normal relaxed
76  clock (25) calibrated with serial samples. This calibration provided a mean substitution rate
77  estimate of 1.3 x 10-4 substitutions per site per year. Wertheim et al. (26) used this estimate and
78  the same genome region (RdRp), but with a codon-based substitution model accounting for the
79  effect of selection. Indeed, purifying selection can lead to an underestimation of viral origins when
80  not accounted for (26, 27). They found an ancient origin, around 293 (95% confidence interval,
81  190 to 489) million years ago (26). More recently, Hayman & Knox (28) obtained similar results,
82  but using the splitting times of hosts as constraints, therefore assuming a priori that coronaviruses
83  codiversified with their hosts.
84
85      More generally, dating the phylogenies of RNA virus families is a difficult task (29). While
86  for some of them dated calibration points can be used, based on orthologous copies of
87  endogenous virus elements (EVEs) present in the genomes of related mammalian species with
88  known times of divergence (30), in many others, including in the Coronaviridae, such elements
89  have not been found (31). Despite the difficulty in dating viral families, it has been proposed, from
90  cophylogenetic analyses investigating the congruence of the host and viral phylogenetic trees
91  (32), that vertebrate-associated RNA viruses have codiversified with their hosts over hundreds of
92  millions of years (31, 33). Indeed, RNA virus phylogenies tend to mirror that of their hosts; for
93  example, closely-related coronaviruses infect closely-related mammals (e.g. (28)). However, a
94  major caveat is that such cophylogenetic signals can emerge when viruses diversify by host
95  switches preferentially occurring among closely-related hosts, in the absence of any cospeciation
96  event (34). Event-based cophylogenetic methods can in principle identify cospeciation and host
97  switches events (32, 34), but their behavior in the presence of diversification by preferential host
98  switches is not well understood. Under a perfect codiversification scenario, host and symbiont
99  phylogenies would be identical. Events of host switches, duplications and losses induce
100  mismatches, and cophylogenetic methods aim to identify parsimonious sets of events that allow
101  "reconciling" the two phylogenies (34, 35). However, most of these methods rely entirely on tree
102  topology (and not branching times), such that time-inconsistent host switches between non-
103  contemporary host lineages are allowed during the reconciliation. In the presence of preferential
104  host switches, these methods may thus favor biologically unrealistic reconciliations that involve
105  cospeciation events and 'back-in-time' host switches to reconciliations that involve more frequent
106  contemporary host switches. This would have remained unnoticed, unless users of the methods
107  specifically looked at the time consistency of the inferred host switches, which is usually not done.
108
109      Here, we establish a framework for testing scenarios of ancient origination and
110  codiversification versus recent origination and diversification by host switches that combines
111  probabilistic cophylogenetic models and biogeographic analyses (Fig 1). We then apply this
112  framework to the Coronaviridae-mammals association. We assemble a dataset of all mammalian
113  hosts of coronaviruses and a complete association matrix between host species and species-like
114  Operational Taxonomic Units (sOTUs) of coronaviruses, as well as geographic repartition of
115  Coronaviridae and their mammalian hosts. We construct a new Coronaviridae tree based on a
116  recent proposition for the use of a well-conserved region of their RNA genome (36, 37). Under the
117  ancient origination scenario (Fig 1A), long-term vertical transmission of Coronaviridae within
118  mammalian lineages could lead to events of mammal-coronavirus cospeciations. Coronaviruses'
119  diversification would then be modulated by both cospeciations and horizontal host switches from
120  one mammalian lineage to another (26, 31). The most recent common ancestor of coronaviruses
121  could even have infected the most recent common ancestor of mammals and birds (26). Under
122  the recent origination scenario (Fig 1B), codiversification with hosts is virtually impossible, and
123  coronaviruses' diversification would then be largely dominated by recent host switches.
124  Expectations for the output of reconciliation and biogeographic analyses under these different
125  scenarios, as well as a scenario of random associations, are explicated in Fig 1. We identify the
126  likely origination of coronaviruses in the mammalian tree, quantify the frequency of cospeciation
127  and host-switching events, and locate these host switches, therefore identifying 'reservoirs' of
128  Coronaviridae and potential transmission routes across mammals.

129

## Results

131

132 By screening the 46 sOTUs of Coronaviridae identified by Edgar et al. (36) in public databases,
133 we found 35 that were associated with mammalian hosts. Our trees of these 35 sOTUs support a
134 well-defined split between Alphacoronaviruses and the other genera, regardless of the
135 phylogenetic method used (Fig 2; *SI Appendix*, Fig. S1). Overall, Alphacoronaviruses form a
136 monophyletic clade, Delta- and Gammacoronoviruses form sister clades, with the main
137 uncertainty being on the placement of their ancestor in relation to Beta (i.e. as a sister to a
138 monophyletic Beta-clade (Fig. 2) or within the Beta-clade (*SI Appendix*, Fig. S1).

139

140 We found that mammalian hosts of coronaviruses belong to 31 families and 10 orders of
141 mammals, and are widely distributed throughout the mammalian phylogeny (SI Appendix, Fig.
142 S2). Most mammalian hosts are bats (Chiroptera - 55 species), followed by rodents (Rodentia -
143 22 species), artiodactyls (Artiodactyla - 15 species), carnivores (Carnivora - 11 species), and
144 primates (Primates - 5 species). Five other orders have at least one representative species:
145 Eulipotyphla (4), Lagomorpha (1), Perissodactyla (1), Pholidota (1), and Sirenia (1). The number
146 of mammalian hosts per coronavirus' sOTU varies across the Coronaviridae tree, ranging from 1
147 to 22 species, with an average of 4.94 (Fig. 2). Of the 35 sOTUs, 23 are found in at least one bat
148 species and 17, mostly in alphacoronaviruses, are found exclusively in bats (Fig. 2). Eight sOTUs
149 are found in humans, six of which, including the three pathogenetic sOTUs, are
150 betacoronaviruses. Betacoronaviruses infect a larger average number of hosts and a larger
151 diversity of non-bat species than Alphacoronaviruses. Twenty-two coronaviruses occur in more
152 than one species; of those, 11 are found in multiple orders (Fig. 2; *SI Appendix*, Fig. S3) and 11
153 in multiple species of a single order (*SI Appendix*, Fig. S3).

154

155 We first tested whether closely-related coronaviruses tend to infect closely-related
156 mammals. A negative answer to this question would suggest that the diversification of
157 Coronaviridae is independent of mammalian history, excluding the scenarios of codiversification
158 or diversification per preferential host switches (Fig. 1). To the contrary, we found a significant
159 phylogenetic signal for the overall association between coronaviruses and mammals (Mantel test:
160 r= 0.38; *P*= 0.0001) and vice-versa (r= 0.29; *P*= 0.0001), after accounting for the confounding
161 phylogenetic signal in the number of partners (38). Mantel tests across sub-clades of both
162 phylogenies revealed that this overall phylogenetic signal is linked to phylogenetic signal in the
163 deep nodes of the Coronaviridae and mammal phylogenies rather than at shallow phylogenetic
164 scales (*SI Appendix*, Fig. S4). This pattern could arise from ancient codiversification followed by
165 un-preferential host switches, or from recent host switches preferentially occurring between hosts
166 from the same high-level taxonomic grouping (such as mammalian orders). We also found that
167 closely related coronaviruses tend to infect a similar number of hosts (r= 0.29; *P*=0.002), while
168 closely related mammals do not tend to host a similar number of distinct coronaviruses (r= 0.04,
169 *P*=0.1), suggesting that coronaviruses' specificity towards hosts is evolutionarily conserved while
170 hosts' susceptibility to coronaviruses is not.

171

172 To further investigate the hypotheses of ancient codiversification versus recent host
173 switches, we used a probabilistic cophylogenetic model, the amalgamated likelihood estimation
174 (ALE - (39)), that reconciliates the host and symbiont phylogenies using events of cospeciations,
175 host switches, duplications, or losses, while accounting for phylogenetic uncertainty in the
176 symbiont phylogenies and undersampling of the host species (35, 39, 40). The main version of
177 ALE we used is an "undated" version that accounts for topology but not branch lengths, as the
178 dated version did not perform well on our data (see Methods). Time-inconsistent host switches
179 are thus allowed during the reconciliation. If the scenario of ancient diversification holds, we
180 expect to find reconciliations requiring more cospeciations and fewer host switches than expected
181 under a scenario of independent evolution (hereafter referred to as 'significant reconciliation'),

182 and few time-inconsistent switches (Fig. 1A). Under the alternative scenario of recent origination
183 and diversification by preferential host switches, we also expect to infer a significant reconciliation,
184 but with many time-inconsistent switches, as the algorithm tends to explain the cophylogenetic
185 signal in the interactions by cospeciation events (Fig. 1B). We indeed found a significant
186 reconciliation between the Coronaviridae and the mammalian trees, confirming the non-
187 independence of their evolution, which we evaluated by randomly shuffling mammal species
188 across the full tree or within biogeographic regions (*SI Appendix*, Fig. S7). ALE reconciliations
189 inferred average numbers of 145 cospeciations, 65 losses, 0 duplication, and 92 host switches.
190 Without investigating the time-consistency of the host switches, we would conclude that there are
191 almost 1.5 more diversification events of Coronaviridae that are related to ancient
192 codiversification rather than host switches. However, on average 20% of the inferred host
193 switches are time-inconsistent, including "back-in-time" host switches of >50 Myr (*SI Appendix*,
194 Fig. S8), which suggests instead that extent Coronaviridae originated recently and diversified by
195 frequent preferential host switches. 64% of the reconciliations found an origination of
196 coronaviruses within bats, in particular within the Pteropodidae family (Fig. 2A-C). We no longer
197 found an origination in bats when randomly shuffling the dataset (*SI Appendix*, Fig. S5,S6),
198 suggesting that this result is not artifactual. We checked the interpretation of our results by
199 simulating the two scenarios of (i) ancient origination in the ancestors of bats followed by
200 codiversification and (ii) recent origination in an extant bat species and a subsequent
201 diversification by preferential host switches. On the first set of simulations, ALE correctly inferred
202 an origination in bats and a few time-inconsistent switches (*SI Appendix*, Fig. S12). On the
203 second set, ALE correctly inferred an origination in bats, although with lower confidence, and a
204 large fraction (~20%) of time-inconsistent host switches, similar to what we observed for
205 Coronaviridae. These results therefore indicate a scenario of recent origination of coronaviruses
206 in bats followed by diversification by preferential host switches.
207
208    To investigate this scenario in more detail, we gradually applied a tree transformation to
209 the mammalian phylogeny, which excludes the possibility of an ancient origination happening
210 earlier than a given time. We found that we had to impose a very recent time of origination
211 (younger than 5 Myr) to obtain few time-inconsistent switches (*SI Appendix*, Table S1). We thus
212 carried out our follow-up analyses with a mammals' tree transformation (star phylogeny) that
213 assumes an origination in an extant mammalian lineage, such that coronavirus diversification is
214 explained entirely by host switches between extant mammalian species. Simulations validated
215 this approach in terms of properly inferring originations and identifying preferential host switches
216 (*SI Appendix*, Fig. S13). Applied to the data, the approach inferred a high probability of origination
217 in bats (56%, Fig, 2B-C, *SI Appendix*, Fig. S6) and a scenario of diversification by preferential
218 host switches: 68% of the inferred host switches happened within mammal orders (Fig 2D, *SI
219 Appendix*, Fig. S10), whereas we would expect on average only 28% of within-order host
220 switches if happening at random. We also inferred more-than-expected host switches between
221 closely related mammal orders (e.g. between Artiodactyla and Perissodactyla) and between the
222 order containing humans (Primates) and those of their domesticated animals, such as
223 Artiodactyla and Carnivora (*SI Appendix*, Fig. S10, Table S2). In contrast, host switches were five
224 times less numerous than expected by chance between bats and other orders (10.7%, against
225 50.2% on average if host switches were randomly distributed, Fig. 2D), in particular Artiodactyla
226 and Rodentia (*SI Appendix*, Fig. S11, Table S2). When occurring, host switches from bats often
227 occurred toward humans (1.9 host switches per reconciliation on average) or toward urban-living
228 and/or domesticated animals, such as rats, camels, or pigs (>1 host switch on average; *SI
229 Appendix*, Table S3). Host switches to humans occurred mostly from domesticated mammals
230 (camels, pigs, dogs), the house shrew and the house mouse, then followed by Asian palm civets,
231 and lastly by bats and other rodents (*SI Appendix*, Table S4). Results were consistent when
232 subsampling the dataset to have an equal sampling effort per host species, suggesting that our
233 results are not artifactually explained by the enhanced monitoring of coronaviruses in humans or
234 domesticated animals (*SI Appendix*, Supplementary Information Text). Finally, we found that
235 some sOTUs, in particular from Betacoronaviruses (e.g, u24667 and u175, both with humans

236 among their hosts), have experienced frequent host switches, whereas others have not (e.g.
237 u165, which is restricted to pigs). In particular, u944 (SARS-Cov-2) has experienced an
238 intermediate number of host switches compared to other coronaviruses (*SI Appendix*, Fig. S9).
239

240      We found qualitatively similar results when applying ALE on different sub-parts of the
241 palmprint region, suggesting that the potential occurrence of recombination does not bias our
242 conclusions (*SI Appendix*, Table S5). The percentage of originations inferred to occur in bats
243 decreased in the analyses on the first sub-part, probably because using such a short fragment
244 (75 aa-long) does not allow robust reconciliations. We also obtained consistent results using a
245 reconciliation method based on maximum parsimony (eMPRess) instead of maximum likelihood
246 (ALE). Whatever the costs that we set for the different reconciliation events, eMPRess estimated
247 significant reconciliations (p-values<0.01). For instance, when favoring host switches, we inferred
248 a recent origination in bats in 54% of the reconciliations and observed on average 32
249 cospeciations (s.d.   3), 2 losses (s.d.   1), 0.1 duplication (s.d.   0.3), and 140 host switches
250 (s.d.   3) including several "back-in-time" host switches of >30 Myr. eMPRess therefore also
251 supports a scenario of recent origination in bats and diversification by preferential host switches
252 (Fig. 1B). Without investigating the time-consistency of the host switches, we would have wrongly
253 concluded that almost one fourth of the diversification of Coronaviridae is related to ancient
254 cospeciation events.
255

256      An additional piece of evidence for a recent origination scenario comes from the
257 geographical distribution of coronaviruses, with a hotspot of diversity in Eurasia that has not
258 colonized the whole world (Fig. 3A, Fig. 1B). The coronavirus' hotspot is more strongly influenced
259 by the diversity of alphacoronaviruses than of betacoronaviruses (*SI Appendix*, Fig. S14). The
260 higher host switches rates and broader host range of betacoronaviruses is reflected in a more
261 widespread geographic distribution, with less pronounced hotspots when compared to
262 alphacoronaviruses (*SI Appendix*, Fig. S14). Mammalian hosts of coronaviruses have a hotspot
263 of species diversity concentrated in East Asia (Fig. 3C). The richness of coronaviruses presents a
264 similar pattern, but with two comparable hotspots of species diversity in East Asia and Southern
265 Europe (Fig. 3A), suggesting that the European hotspot is composed by fewer host species,
266 together carrying as diverse a set of coronaviruses as the Asian hotspot. Other regions with a
267 relatively high richness of coronaviruses and their hosts include parts of the African continent.
268 The Americas and Australia have relatively low richness of coronaviruses and their hosts.
269 Phylogenetic diversity of both hosts and coronaviruses (Fig. 3C,D) depict a similar pattern but
270 with phylogenetic diversity more evenly distributed across most world regions, including the
271 Americas.
272


273 **Discussion**
274
275 Together, our results suggest that the common ancestor of extant mammalian coronaviruses
276 originated recently in a bat species, and that coronaviruses diversification occurred via
277 preferential host switches rather than through codiversification with mammals. Although we
278 cannot unequivocally reject that ancestors of present-day coronaviruses were not present several
279 million years ago, we demonstrate that Coronaviridae is a highly dynamic clade in which
280 diversification operates through host switches at a much faster pace than that of their hosts.
281 sOTUs are rapidly replaced by newly-generated ones, with little role for codiversification with the
282 hosts. The high diversity and endemicity of coronaviruses among bats has led others to anticipate
283 that bats might be implicated in the origin of coronaviruses (2, 20, 23, 24), although definitive
284 proof was lacking. We provided evidence for that hypothesis using a probabilistic cophylogenetic
285 model after sampling the entire diversity of coronaviruses across mammals. Independent
286 evidence for coronavirus recent host switches among different species exists in the literature (41,
287 42). The envisioned scenario suggests a timing of origination for extant Coronaviridae that is
288 much more recent than the hundreds of millions of years ago suggested by (26). This is not

6

289  surprising given the difficulties in estimating divergence times and inferring branch lengths for
290  viral phylogenies (26, 27, 29), and provided that the dating of Wertheim et al. (26) relied on a
291  substitution rate estimated from data with limited temporal signal (~50 serially sampled
292  contemporary sequences of a short gene fragment, (23)).
293
294        Our results contradict previous suggestions that codiversification with vertebrate hosts
295  played an important role in Coronaviridae diversification (26, 28, 31). They also suggest that
296  previously reported cases of long-term codiversification in vertebrate RNA viruses have been
297  largely over-estimated, as many of them may instead be cases of diversification by host switches
298  occurring preferentially among closely-related hosts. Indeed, these two scenarios both generate
299  cophylogenetic signal in host-symbionts associations, such that cophylogenetic signal alone is
300  not evidence for long-term codiversification (34). In addition, under a scenario of recent
301  origination and preferential host switches, event-based cophylogenetic methods tend to
302  artifactually favor biologically unrealistic scenarios with codiversification and back-in-time host
303  switches, as we have shown here. As the time-consistency of host switches is typically not
304  investigated, this has remained unnoticed, and evidence for codiversification has been taken for
305  real. Ideally, cophylogenetic reconciliation methods would not allow such time-inconsistent host
306  switches. However, imposing time constraints in methods based on parsimony is NP-hard (43),
307  and the 'dated' version of ALE is not well adapted when recent host switch events dominate
308  evolutionary history. We have found two ways to get around the problem, by interpreting time-
309  inconsistent host switches as evidence for recent preferential host switches, and by gradually
310  transforming the host tree to avoid large back-in-time switches, however future efforts should
311  focus on developing time-consistent cophylogenetic methods. This would allow more robust and
312  precise inferences of host-virus (and more generally host-symbiont) evolutionary history.
313
314        During their evolution, coronavirus' host switches occurred more frequently within than
315  between mammalian orders. This suggests that mammalian characteristics shared between
316  relatives (e.g., genetic, behavioral, ecological), and the frequency of encounters among hosts
317  play important roles in determining coronavirus' host switches. Additionally, between-order host
318  switches occurred more frequently among non-flying mammals and among orders containing
319  humans and urban and domesticated mammals, suggesting that contact frequency alone is likely
320  a key characteristic in host switches. Accordingly, amongst the most-likely host switches towards
321  humans were those coming from mammals suspected to be involved in the transfer of specific
322  coronavirus sOTUs likely through contact, for instance, camels in the case of MERS-CoV (41),
323  Asian palm civets with SARS-CoV (16, 17) and the house mouse with SARS-CoV2 (42).
324  Importantly, we found that host switches from bats to other mammalian species were rare during
325  the evolutionary history of Coronaviridae, even though coronaviruses originated and are more
326  diverse within bats. These pieces of evidences suggest that bats are a closed reservoir of the
327  Coronaviridae diversity.
328
329        Spillovers from bats to non-bat species, when they occurred, were found more likely to be
330  towards humans than to any other mammalian species, suggesting humans may have acted as
331  evolutionary intermediate hosts amongst mammals. From an ecological perspective, the large
332  abundance and widespread geographic distribution of humans, together with our habits of forcing
333  contact with other species, including bats, make it unsurprising that humans, among all mammal
334  species, have acted as intermediate hosts of ancestral forms of coronaviruses. Interestingly, for
335  some individual species of coronaviruses, such as the SARS-CoV2 and other SARS-like
336  coronaviruses, the dominant hypothesized scenario is that precursor forms spread from a bat to
337  another intermediate mammalian host before infecting humans (20, 44). Our molecular marker
338  lacks the intra-OTU resolution necessary to make species-level predictions, but our results
339  suggest that more ancient coronaviruses host switches may have occurred in the other direction:
340  from bats to humans to non-bat mammals. Many human activities lend credit to the human-as-
341  evolutionarily-intermediate-host-hypothesis, including human excursions to bat caves (45),
342  hunting (46), and habitat destruction and modification (47), all of which increase the contact

7

343 between bats and humans and their domesticated animals (47). Conservation of bats' natural
344 habitats, away from human contact, could help avoiding further spreads of coronaviruses among
345 humans.
346
347      Insights of past and future host switches are gained from coronavirus geographic
348 distribution. Coronaviruses are found worldwide and their hotspots of diversity are concentrated in
349 East Asia and Southern Europe, where they likely originated. Previous assessments of the
350 diversity of bat hosts of betacoronaviruses suggested similar hotspots but with a distribution of
351 coronaviruses more concentrated in the hotspots (7, 14, 15) than the more pervasive pattern we
352 found using all mammalian hosts. Moreover, the distribution of coronaviruses is less concentrated
353 in the hotspots when phylogenetic metrics of diversity are included, suggesting that species
354 richness alone is masking the global evolutionary potential of these viruses (48). Coronaviruses'
355 likely recent origination in bats, high within-order transmission rates, and their capacity to switch
356 between mammal orders in some cases suggest the potential for future fast spreading and
357 increase in the number of species across most world regions. Among alphacoronaviruses, the
358 spread is more likely to remain concentrated within bats, while betacoronaviruses have a higher
359 potential for among-orders spreading and infection of new mammalian hosts. The
360 betacoronaviruses lineages already detected in humans have especially high host generalism
361 and transmission rates, indicating that these lineages should be particularly monitored to avoid
362 future pandemics.
363
364      Finally, a few important limitations of our analyses deserve to be mentioned.
365 Recombination is an important mechanism of viral evolution (49), and approaches more
366 adequately designed to investigate the role of recombination are needed. The fact that different
367 subparts of the palmprint region lead to similar results indicates that recombination acting on the
368 palmprint region is unlikely to bias our conclusions. However, looking at other genomic regions
369 would allow gaining a more complete understanding of the role of recombination in coronavirus
370 evolution. Lastly, because the palmprint region is a conserved region, we could not reconstruct
371 the recent evolutionary history of coronaviruses (i.e. the within sOTU transmission dynamic):
372 combining the palmprint region with a fast-evolving region(s) would enable more precise
373 estimates of the recent routes of coronaviruses' transmission, including that of SARS-CoV-2.
374
375      Understanding the evolutionary origins and diversification of viruses is crucial to any
376 attempt of predicting new transmission routes, yet the relative frequencies of virus–host
377 cospeciation versus cross-species transmission in the evolution of vertebrate RNA viruses
378 remains uncertain (31). We found that coronaviruses originated in bats where they are more
379 diverse nowadays, and later diversified in other mammal orders through preferential host
380 switches. Spillovers from bats were rare but likely human-induced, suggesting humans are the
381 intermediate evolutionary bridge that facilitated the spread of coronaviruses across mammals.
382 Host switches between primates and artiodactyls, perissodactyls, and carnivorans happen at high
383 rates and we can thus expect a spread of coronaviruses amongst new mammalian hosts and
384 outside of their current diversity hotspots in East Asia and Europe, as well as future
385 coronaviruses-related pandemics. Our results suggest reducing human-bat contact, for example
386 by conserving bat habitats, as a mitigation strategy. They also suggest that cases of long-term
387 virus–host codiversification, reported on the basis of cophylogenetic tests, have been largely
388 over-estimated.
389

390 **Materials and Methods**
391
392 **Operational Taxonomic Units for Coronaviridae**

393 We used the 46 described species-like Operational Taxonomic Units (sOTUs) for Coronaviridae
394 delimited using 'palmprint' sequences by (36, 37). The palmprint is a conserved amino acid (aa)

395  sub-sequence (150 aa in Coronaviridae) of central importance in the viral RdRp (36), selected for
396  its homology across the large majority of sequences, allowing estimation of sequence divergence
397  and phylogenetic trees (37). sOTUs were identified by Edgar et al. (36) after clustering palmprint
398  sequences at 90% amino acid identity; and released through the Serratus project. Their approach
399  is equivalent to the species delimitation proposed by the International Committee on the
400  Taxonomy of Viruses for Coronaviridae ((2; *SI Appendix*, Supplementary Information Text),
401  which suggests 90% similarity of amino acid sequences for conserved domains (2, 37), and are,
402  therefore, ideal for species tree construction. Under the 'palmprint framework' a centroid definition
403  of species is applied to characterize a new OTU when a threshold of 90% amino acid identity is
404  surpassed, serving as a useful taxonomic barcode (37). We downloaded the palmprint amino acid
405  sequences of Coronaviridae sOTUs from the Serratus project (https://serratus.io/; (36)) on April
406  13 of 2022.
407

**Mammalian hosts of Coronaviridae**

409  All 46 sOTUs of Coronaviridae with a full palmprint and associated data in the NCBI database
410  were screened for the identification of its hosts. From those, 35 sOTUs were associated with
411  mammalian hosts and were kept for downstream analyses. Serratus' associated metadata was
412  used to identify GenBank accession codes linked to each sOTU. The complete set of 90,540
413  associated GenBank accession codes was screened to obtain the host information for each
414  sOTU (on NCBI, Features>source>/host=). All the host species with a full Linnean name were
415  kept as such. Accession codes with hosts leading to a generic level information were further
416  inspected to identify the associated publication and determine the complete species name.
417  Dubious cases or accession codes without publications had their hosts disregarded. Common
418  names or high-level host information (e.g., host="bats") were generally eliminated except in a few
419  cases where a domesticated species was found to be the host (i.e., host="dog","canine" were
420  Canis lupus; host="cat","feline" were Felis catus; host="pig","piglet","newborn piglet","sucking
421  piglet","porcine","swine" were Sus scrofa). A final dataset of 116 mammalian hosts associated
422  with the 35 sOTUs was assembled and used in downstream analyses. A matrix with the
423  association between Coronaviridae sOTUs and mammalian species is available in *SI Appendix*,
424  Dataset S1.
425

**Coronaviridae phylogenetic trees**

427  We constructed a Coronaviridae tree using the palmprint amino acid sequence information of the
428  35 sOTUs. We aligned the amino acid sequences with MAFFT (50) and trimmed them with trimAl
429  (51). The final alignment contained 150 amino acid positions. We used two main phylogenetic
430  software, BEAST2 (52) and PhyloBayes (53), both to assess the robustness of the tree to the
431  phylogenetic method, and because they have different advantages (e.g. rooting and time
432  calibration are performed in BEAST2 while PhyloBayes outputs are adapted to the
433  cophylogenetic algorithm we used). We visualized phylogenetic trees using R (54).
434      In order to run BEAST2, we generated an input file using BEAUti with the following
435  parameters: a WAG model with 4 classes of rates and invariant sites, a birth-death prior, and a
436  relaxed log-normal clock. BEAST2 sampled a posterior distribution of ultrametric trees using
437  Markov chain Monte Carlo (MCMC) with 4 independent chains each composed of 100,000,000
438  steps sampled every 10,000 generations. We checked the convergence of the 4 chains using
439  Tracer (55). We used LogCombiner to merge the results setting a 25% burn-in and TreeAnnotator
440  to obtain a Maximum Clade Credibility (MCC) tree with median branch lengths. PhyloBayes was
441  run using an LG model, 4 classes of rates, and a chain composed of 4,000 steps with a 25%
442  burn-in.
443      To further assess the robustness of the BEAST2 tree rooting, we estimated the root
444  position on a 46-sOTU maximum likelihood tree (from the Serratus project - (36)) assuming a

445 strict molecular clock and an ultrametric tree. We used an ultrametric setting as temporary
446 information from the tip dates (ranging between 1999 and 2022, a neglectable difference with
447 respect to the root age of dozens of thousands or even millions of years) was not sufficient to
448 infer the mutation rate (we assessed the temporal signal with TempEst, (56)). We performed
449 rooting and time-scaling with LSD2 (v2.3, (57)), assuming a tree of unknown scale (e.g. fixing all
450 the tips dates to 1 and the root date to 0) with outlier removal and root search on all branches.
451 LSD2 detected no outliers and positioned the root on the same branch as in the BEAST2 MCC
452 tree (between alpha and betacoronaviruses).
453

## Mammalian phylogenetic tree

455 We obtained a phylogenetic hypothesis for mammals from the consensus DNA-only tree of (58),
456 one of the most complete and updated phylogenies for mammals. We downloaded the node-
457 dated tree for 4,098 mammals, constructed based on a 31-gene supermatrix, from the VertLife
458 website (http://vertlife.org/data/mammals/). We used a pruned version of the tree with the 116
459 mammalian hosts of Coronaviridae in all analyses in this paper.
460

## Phylogenetic signal in the association between coronaviruses and mammals

462 To assess whether closely related coronaviruses interact with similar mammals, and vice-versa,
463 i.e. presence of phylogenetic signal in the association, we used Mantel tests following (38).
464 Mantel tests were constructed by taking the Pearson correlation between phylogenetic distances
465 and ecological distances. Phylogenetic distances of coronaviruses were computed on the
466 BEAST2 MCC phylogeny. Ecological distances were calculated based on the interaction network
467 matrix containing the association between coronavirus' sOTUs and mammals, accounting for the
468 evolutionary relationships among interaction partners using UniFrac distances (59). Firstly, we
469 conducted Mantel tests permuting the identity of species but keeping the number of partners per
470 species constant; this allows for assessing the effect of species identity while controlling for the
471 confounding effect of the number of partners. Then, we evaluated the phylogenetic signal in the
472 number of partners alone. Lastly, we calculated clade-specific Mantel tests for sub-networks
473 containing at least 10 species (38) to evaluate whether phylogenetic signal was stronger for
474 specific subclades of mammals or coronaviruses. Ten thousand permutations were used in each
475 analysis to assess significance. Analyses were conducted using the phylosignal_network and
476 phylosignal_sub_network functions in the R package RPANDA (60).
477

## Coronaviridae origination and host switches

479 We used the amalgamated likelihood estimation (ALE - (39)) to reconciliate the mammal and
480 coronaviruses evolutionary history using events of cospeciations, host switches, duplications, and
481 losses. Originally designed in the context of gene tree – species tree reconciliations (39), ALE
482 has also been particularly useful in the context of host-symbiont cophylogenetic analyses as it
483 considers both phylogenetic uncertainty of the symbiont evolutionary history and undersampling
484 of host species (35, 61, 62). ALE indeed assumes that host switches may imply an unsampled or
485 extinct intermediate host lineage (40). We ran ALE with the posterior distribution of phylogenetic
486 trees of coronaviruses generated with PhyloBayes to estimate the maximum likelihood rates of
487 host switches, duplications, and losses of the coronaviruses. We first tried running the "dated"
488 version of ALE, which accounts for the order of branching events in the host phylogeny, therefore
489 only allowing for time-consistent host switches (i.e. host switches that happen between two
490 contemporary host lineages). However, this led to unrealistic parameter estimates (such as very
491 high loss rates) and ALE was not able to output possible reconciliations, suggesting that the
492 mammalian and Coronaviridae trees are too incongruent to be reconciled with only time-

10

493   consistent host switches. We therefore used the "undated" version of ALE that only exploits the
494   topology of both the host and the symbiont tree and thus does not constrain the host switches to
495   be time-consistent. ALE generated a total of 5,000 reconciliations, from which we extracted the
496   mean number of cospeciations, host switches, duplications, and losses. We also reported the
497   likely origination of coronaviruses in mammals (i.e. the branch in the mammal phylogeny that was
498   first infected by coronaviruses) by computing, for each branch of the mammalian tree, the
499   frequency of reconciliations (among the 5,000) that supported an origination in that branch. If a
500   reconciliation requires more cospeciation events and fewer host switch events, than expected
501   under a null scenario of independent evolution, this indicates that the evolution of the symbiont
502   was not independent of that of the host, and in this case, we talk about a "significant
503   reconciliation" (63). We evaluated the significance of the reconciliation by comparing the
504   estimated number of cospeciation and host switch events to null expectations obtained with ALE
505   by shuffling the mammal host species across the mammal tree, both randomly or within major
506   biogeographic regions according to the proposal of regions by (64) for mammals (six
507   biogeographic regions: North American, South American, African, Eurasian, Oriental, and
508   Australian). We considered a reconciliation to be significant if the observed number of
509   cospeciations was higher than 95% of the null expectations and if the number of host switches
510   was lower than 95% of the null expectations (35). The likeliness of a host switch between two
511   mammal lineages is measured as the frequency of the reconciliations in which it occurs. Finally,
512   we reported the ratio of time-inconsistent host switches by focusing on "back-in-time" switches,
513   from a donor mammal lineage to an older receiver mammal lineage that never coexisted.

515   Because ALE estimated a large proportion of time-inconsistent host switches (see
516   Results), we first tested the scenario of a more recent origination by collapsing all mammalian
517   nodes anterior to X Myr into a polytomy at the root of the phylogeny (with X varying from 55 Myr
518   to 5 Myr), such that the coronavirus origination and host switches inferred by ALE could not
519   involve mammal lineages older than X Myr. Second, we investigated the scenario of
520   diversification by pure preferential host switches of the coronaviruses among extant mammals. To
521   do so, we ran ALE on a star mammalian phylogenetic tree. In this context, ALE could no longer
522   infer cospeciations, and only fit events of host switches, duplications, or losses. When inferring a
523   likely host switch between two specific mammalian lineages on a star phylogeny, there are often
524   as many reconciliations suggesting one directionality of the host switch (i.e. from one of the
525   lineages to the other) as the other. We then only kept host switches present in at least 10% of the
526   reconciliations and looked at the ratio between the number of host switches that were estimated
527   within versus between mammal orders. We compared this ratio to a null expectation obtained by
528   randomly shuffling the host mammal species.

530   Recombination is frequent in viruses and the palmprint region may be recombined, such
531   that different fragments of the palmprint region may have different evolutionary histories,
532   potentially biasing our inference. To test whether the results we obtained on the whole 150-amino
533   acid palmprint region were not impacted by recombination, we replicated the ALE analyses on
534   two sub-regions: the first part (positions 1-75) and the last part (positions 76-150).

536   Finally, we repeated our cophylogenetic analyses using eMPRess (43), another event-
537   based cophylogenetic approach that reconciliates host-symbiont evolutionary histories using
538   maximum parsimony. eMPRess is a recent improved version of the popular Jane approach (32);
539   it differs from Jane especially by not only relying on a heuristic and therefore guarantying that the
540   solution truly corresponds to the maximum parsimony reconciliation(s) (43). However, contrary to
541   Jane, eMPRess does not allow the same symbiont species to be present in different host species,
542   and does not offer the possibility to constrain host switches to occur only among lineages from
543   pre-specified time periods. eMPRess requires specifying cost values for the events of host
544   switches (t), duplications (d), and losses (l). We tested two sets of cost values: (1) cost values
545   that disadvantage host switches (d=6, t=6, l=1) and (2) uniform cost values that favor host
546   switches (d=1, t=1, l=1). As with ALE, we evaluated the significance of the reconciliations using

11

547 permutations. We ran eMPRess analyses on a set of 50 trees randomly sampled from the
548 posterior distribution of PhyloBayes.
549

550 **Simulation analyses**

551 By running the undated version of ALE either on the mammal phylogeny or a star phylogeny, we
552 proposed a framework to evaluate whether the cophylogenetic pattern is due to a history of
553 ancient codiversification (i.e. a mix of cospeciations, host switches, duplications, and losses; Fig.
554 1A) or to a scenario where the coronaviruses diversify more recently by preferential host switches
555 ((34); Fig. 1B). To validate the interpretation of our ALE results, we performed simulations under
556 the two alternative scenarios of codiversification and diversification by preferential host switches.
557 For the scenario of codiversification, we assumed that coronaviruses originated in the ancestors
558 of bats and that they subsequently codiversified with the mammals by experiencing events of
559 cospeciations, host switches, duplications, and losses. We used the function *sim_microbiota* in
560 the R-package HOME to obtain the corresponding coronavirus sequences and coronavirus-
561 mammals associations (65) (*SI Appendix*, Supplementary Information Text). For the scenario of
562 coronaviruses diversification by preferential host switches, we used a birth-death model (pbtree
563 function in the R-package phytools) to simulate a phylogenetic tree of the coronaviruses: in our
564 model, each coronavirus lineage is associated with a single host species, a birth event
565 corresponds to a host switch (at rate 50), while a death event corresponds to a loss of a
566 coronavirus in a host lineage (at rate 5). We started the diversification by assuming a single
567 coronavirus infection in *Eidolon helvum* (a bat host of external lineages within Betacoroviruses,
568 u25738 and u27845). Then, following de (66) and (67), we modeled preferential host switches by
569 assuming that for a host switch from a given donor mammal species, each potential receiver
570 species has a probability proportional to $\exp(-0.035*d)$ where $d$ is the phylogenetic distance
571 between the donor and receiver species. Finally, we simulated DNA sequences of the
572 coronavirus sequences using the function *simulate_alignment* in HOME. For each type of
573 simulation, we generated 50 simulated datasets of mammal-coronavirus associations. For each
574 dataset, we ran PhyloBayes and ALE on both the mammalian phylogeny and the star phylogeny.
575

576 **Geographic distribution of Coronaviridae**

577 We downloaded geographic range maps for each mammalian host species, with the exception of
578 *Homo sapiens*, from the Map of Life website (https://mol.org/species/); see (68). These maps
579 follow the taxonomy of the Mammal Diversity Database (69) supplemented with the Handbook of
580 the Mammals of the World (HWM) database and the Alien Checklist database for invasive
581 species (68).
582

583 We created a world map with hexagonal, equal-area grid cells of 220km on which we
584 mapped host and coronavirus species diversity, using the Mollweide world projection to
585 accurately represent areas. At large spatial scales, cells with ~220km resolution return more
586 reliable diversity estimates than smaller cells (70). We considered that a host species was
587 present in any given cell if its range covered at least 30% of the cell area to avoid overestimating
588 diversity. We calculated host species diversity as a simple sum of the species occurring in any
589 given cell, and host phylogenetic diversity as Faith's phylogenetic diversity index (PD - (71)) for
590 each cell. We mapped Coronaviridae diversity using the host-filling method (72): we constructed
591 a range map for each Coronaviridae sOTU by overlapping the range maps of all its hosts. We
592 consider the host filling method appropriate in this case because coronaviruses are obligatory
593 parasites that can only live inside hosts. Next, we calculated Coronaviridae sOTU diversity by
594 summing range maps overlapping on each cell, and Coronaviridae phylogenetic diversity as
595 Faith's PD (71). We created these maps in R (54) using the packages epm (73), sf (74), and ape
596 (75).
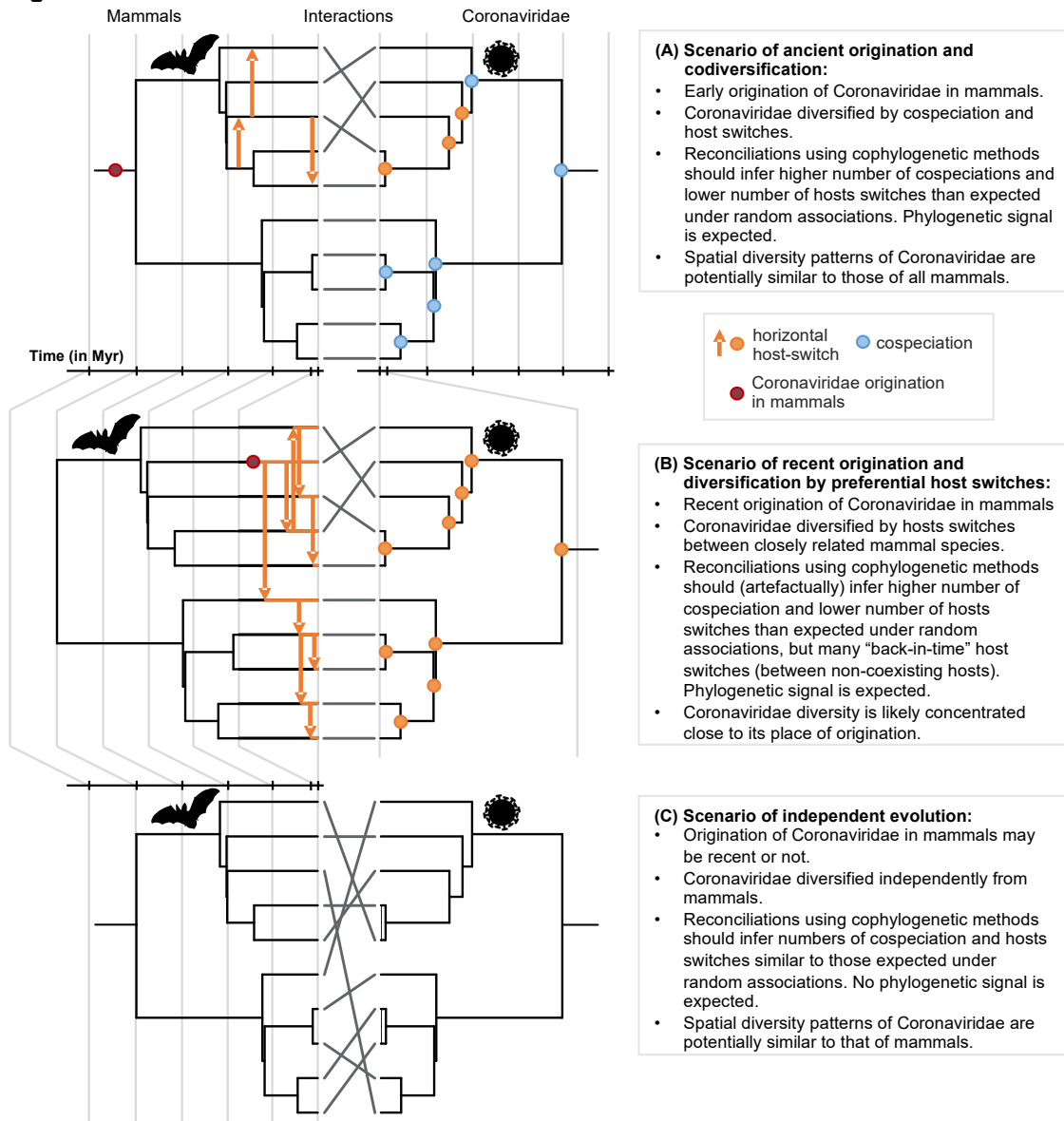
597

## Acknowledgments

603

## References

605    1.   K. P. Alekseev, et al., Bovine-Like Coronaviruses Isolated from Four Species of Captive
606         Wild Ruminants Are Homologous to Bovine Coronaviruses, Based on Complete Genomic
607         Sequences. J Virol 82, 12422–12431 (2008).
608    2.   R. J. de Groot, et al., Family: Coronaviridae. ICTV Ninth Report.
609         https://ictv.global/report_9th/RNApos/Nidovirales/Coronaviridae (2022).
610    3.   A. M. Q. King, Adams M.J., Carstens E.B., Lefkowitz E.J., Virus taxonomy: classification
611         and nomenclature of viruses: Ninth Report of the International Committee on Taxonomy of
612         Viruses (Elsevier, 2012).
613    4.   J. C. Leao, et al., Coronaviridae—Old friends, new enemy! Oral Dis 28, 858–866 (2022).
614    5.   M. Wardeh, M. Baylis, M. S. C. Blagrove, Predicting mammalian hosts in which novel
615         coronaviruses can be generated. Nat Commun 12 (2021).
616    6.   C. M. Zmasek, E. J. Lefkowitz, A. Niewiadomska, R. H. Scheuermann, Genomic evolution
617         of the Coronaviridae family. Virology 570, 123–133 (2022).
618    7.   D. J. Becker, et al., Optimising predictive models to prioritise viral discovery in zoonotic
619         reservoirs. Lancet Microbe 3, e625–e637 (2022).
620    8.   C. Drosten, et al., Identification of a Novel Coronavirus in Patients with Severe Acute
621         Respiratory Syndrome. New England Journal of Medicine 348, 1967–1976 (2003).
622    9.   J. Peiris, et al., Coronavirus as a possible cause of severe acute respiratory syndrome.
623         The Lancet 361, 1319–1325 (2003).
624    10.   A. M. Zaki, S. van Boheemen, T. M. Bestebroer, A. D. M. E. Osterhaus, R. A. M.
625         Fouchier, Isolation of a Novel Coronavirus from a Man with Pneumonia in Saudi Arabia.
626         New England Journal of Medicine 367, 1814–1820 (2012).
627    11.   P. Zhou, et al., A pneumonia outbreak associated with a new coronavirus of probable bat
628         origin. Nature 579, 270–273 (2020).
629    12.   H. Ritchie, et al., Coronavirus Pandemic (COVID-19). Published online at
630         OurWorldInData.org (2022).
631    13.   P. C. Y. Woo, Y. Huang, S. K. P. Lau, K. Y. Yuen, Coronavirus genomics and
632         bioinformatics analysis. Viruses 2, 1805–1820 (2010).
633    14.   S. J. Anthony, et al., Global patterns in coronavirus diversity. Virus Evol 3 (2017).
634    15.   N. F. R. Munoz, et al., The coevolutionary mosaic of bat betacoronavirus emer-gence risk.
635         EcoEvoRXiv (2022).
636    16.   Y. Guan, et al., Isolation and characterization of viruses related to the SARS coronavirus
637         from animals in Southern China. Science (1979) 302, 276–278 (2003).
638    17.   C. SMEC, Molecular Evolution of the SARS Coronavirus During the Course of the SARS
639         Epidemic in China. Science (1979) 303, 1666–1669 (2004).
640    18.   R. L. Graham, R. S. Baric, Recombination, Reservoirs, and the Modular Spike:
641         Mechanisms of Coronavirus Cross-Species Transmission. J Virol 84, 3134–3146 (2010).
642    19.   R. J. de Groot, et al., Commentary: Middle East Respiratory Syndrome Coronavirus
643         (MERS-CoV): Announcement of the Coronavirus Study Group. J Virol 87, 7790–7792
644         (2013).

20. V. M. Corman, D. Muth, D. Niemeyer, C. Drosten, "Hosts and Sources of Endemic Human Coronaviruses" in Advances in Virus Research, (Academic Press Inc., 2018), pp. 163–188.

21. E. Mavrodiev, M. L. Tursky, S. Vincent, D. M. Williams, L. Schroder, On Classication and Taxonomy of Coronaviruses (Riboviria, Nidovirales, Coronaviridae) with Special Focus on Severe Acute Respiratory Syndrome-Related Coronavirus 2 (SARS-CoV-2) (2021) https:/doi.org/10.21203/rs.3.rs-282371/v1.

22. P. C. Y. Woo, S. K. P. Lau, Y. Huang, K. Y. Yuen, Coronavirus diversity, phylogeny and interspecies jumping. Exp Biol Med 234, 1117–1127 (2009).

23. P. C. Y. Woo, et al., Discovery of Seven Novel Mammalian and Avian Coronaviruses in the Genus Deltacoronavirus Supports Bat Coronaviruses as the Gene Source of Alphacoronavirus and Betacoronavirus and Avian Coronaviruses as the Gene Source of Gammacoronavirus and Deltacoronavirus. J Virol 86, 3995–4008 (2012).

24. D. Vijaykrishna, et al., Evolutionary Insights into the Ecology of Coronaviruses. J Virol 81, 4012–4020 (2007).

25. A. J. Drummond, S. Y. W. Ho, M. J. Phillips, A. Rambaut, Relaxed Phylogenetics and Dating with Confidence. PLoS Biol 4, e88 (2006).

26. J. O. Wertheim, D. K. W. Chu, J. S. M. Peiris, S. L. Kosakovsky Pond, L. L. M. Poon, A Case for the Ancient Origin of Coronaviruses. J Virol 87, 7039–7045 (2013).

27. J. O. Wertheim, S. L. Kosakovsky Pond, Purifying selection can obscure the ancient age of viral lineages. Mol Biol Evol 28, 3355–3365 (2011).

28. D. T. S. Hayman, M. A. Knox, Estimating the age of the subfamily Orthocoronavirinae using host divergence times as calibration ages at two internal nodes. Virology 563, 20–27 (2021).

29. S. Duchêne, E. C. Holmes, S. Y. W. Ho, Analyses of evolutionary dynamics in viruses are hindered by a time-dependent bias in rate estimates. Proceedings of the Royal Society B: Biological Sciences 281, 20140732 (2014).

30. A. Katzourakis, R. J. Gifford, Endogenous Viral Elements in Animal Genomes. PLoS Genet 6, e1001191 (2010).

31. M. Shi, et al., The evolutionary history of vertebrate RNA viruses. Nature 556, 197–202 (2018).

32. C. Conow, D. Fielder, Y. Ovadia, R. Libeskind-Hadas, Jane: A new tool for the cophylogeny reconstruction problem. Algorithms for Molecular Biology 5 (2010).

33. Y.-Z. Zhang, W.-C. Wu, M. Shi, E. C. Holmes, The diversity, evolution and origins of vertebrate RNA viruses. Curr Opin Virol 31, 9–16 (2018).

34. D. M. de Vienne, et al., Cospeciation vs host-shift speciation: methods for testing, evidence from natural associations and relation to coevolution. New Phytologist 198, 347–385 (2013).

35. B. Perez-Lamarque, H. Morlon, Comparing different computational approaches for detecting long-term vertical transmission in host-associated microbiota in Molecular Ecology, (John Wiley and Sons Inc, 2022) https:/doi.org/10.1111/mec.16681.

36. R. C. Edgar, et al., Petabase-scale sequence alignment catalyses viral discovery. Nature 602, 142–147 (2022).

37. A. Babaian, R. Edgar, Ribovirus classification by a polymerase barcode sequence. PeerJ 10 (2022).

38. B. Perez-Lamarque, et al., Do closely related species interact with similar partners? Testing for phylogenetic signal in bipartite interaction networks. Peer Community Journal 2, XX (2022).

39. G. J. Szöllosi, W. Rosikiewicz, B. Boussau, E. Tannier, V. Daubin, Efficient exploration of the space of reconciled gene trees. Syst Biol 62, 901–912 (2013).

40. G. J. Szöllosi, E. Tannier, N. Lartillot, V. Daubin, Lateral gene transfer from the dead. Syst Biol 62, 386–397 (2013).

41. G. Dudas, L. M. Carvalho, A. Rambaut, T. Bedford, MERS-CoV spillover at the camel-human interface. Elife 7 (2018).

699  42.  C. Wei, et al., Evidence for a mouse origin of the SARS-CoV-2 Omicron variant. Journal
700       of Genetics and Genomics 48, 1111–1121 (2021).
701  43.  S. Santichaivekin, et al., eMPRess: a systematic cophylogeny reconciliation tool.
702       Bioinformatics 37, 2481–2482 (2021).
703  44.  C. C. S. Tan, et al., Transmission of SARS-CoV-2 from humans to animals and potential
704       host adaptation. Nat Commun 13 (2022).
705  45.  N. M. Furey, P. A. Racey, "Conservation Ecology of Cave Bats" in Bats in the
706       Anthropocene: Conservation of Bats in a Changing World, (Springer International
707       Publishing, 2016), pp. 463–500.
708  46.  T. Mildenstein, I. Tanshi, P. A. Racey, "Exploitation of Bats for Bushmeat and Medicine"
709       in Bats in the Anthropocene: Conservation of Bats in a Changing World, (Springer
710       International Publishing, 2016), pp. 325–375.
711  47.  I. Smith, L. F. Wang, Bats and their virome: An important source of emerging viruses
712       capable of infecting humans. Curr Opin Virol 3, 84–91 (2013).
713  48.  S. Leopardi, et al., Interplay between co-divergence and cross-species transmission in
714       the evolutionary history of bat coronaviruses. Infection, Genetics and Evolution 58, 279–
715       289 (2018).
716  49.  M. Pérez-Losada, M. Arenas, J. C. Galán, F. Palero, F. González-Candelas,
717       Recombination in viruses: Mechanisms, methods of study, and evolutionary
718       consequences. Infection, Genetics and Evolution 30, 296–307 (2015).
719  50.  K. Katoh, D. M. Standley, MAFFT Multiple Sequence Alignment Software Version 7:
720       Improvements in Performance and Usability. Mol Biol Evol 30, 772–780 (2013).
721  51.  S. Capella-Gutiérrez, J. M. Silla-Martínez, T. Gabaldón, trimAl: a tool for automated
722       alignment trimming in large-scale phylogenetic analyses. Bioinformatics 25, 1972–1973
723       (2009).
724  52.  R. Bouckaert, et al., BEAST 2: A Software Platform for Bayesian Evolutionary Analysis.
725       PLoS Comput Biol 10, e1003537 (2014).
726  53.  N. Lartillot, H. Philippe, A Bayesian Mixture Model for Across-Site Heterogeneities in the
727       Amino-Acid Replacement Process. Mol Biol Evol 21, 1095–1109 (2004).
728  54.  R Core Team, R: a language and environment for statistical computing (2018).
729  55.  A. Rambaut, A. J. Drummond, D. Xie, G. Baele, M. A. Suchard, Posterior Summarization
730       in Bayesian Phylogenetics Using Tracer 1.7. Syst Biol 67, 901–904 (2018).
731  56.  A. Rambaut, T. T. Lam, L. Max Carvalho, O. G. Pybus, Exploring the temporal structure
732       of heterochronous sequences using TempEst (formerly Path-O-Gen). Virus Evol 2,
733       vew007 (2016).
734  57.  T.-H. To, M. Jung, S. Lycett, O. Gascuel, Fast Dating Using Least-Squares Criteria and
735       Algorithms. Syst Biol 65, 82–97 (2016).
736  58.  N. S. Upham, J. A. Esselstyn, W. Jetz, Inferring the mammal tree: Species-level sets of
737       phylogenies for questions in ecology, evolution, and conservation. PLoS Biol 17 (2019).
738  59.  J. Chen, et al., Associating microbiome composition with environmental covariates using
739       generalized UniFrac distances. Bioinformatics 28, 2106–2113 (2012).
740  60.  H. Morlon, et al., RPANDA: an R package for macroevolutionary analyses on
741       phylogenetic trees. Methods Ecol Evol 7, 589–597 (2016).
742  61.  M. Groussin, et al., Unraveling the processes shaping mammalian gut microbiomes over
743       evolutionary time. Nat Commun 8 (2017).
744  62.  M. Bailly-Bechet, et al., How Long Does Wolbachia Remain on Board? Mol Biol Evol 34,
745       1183–1193 (2017).
746  63.  R. G. Dorrell, et al., Phylogenomic fingerprinting of tempo and functions of horizontal
747       gene transfer within ochrophytes https:/doi.org/10.1073/pnas.2009974118/-
748       /DCSupplemental.
749  64.  C. Barry Cox, The biogeographic regions reconsidered. J Biogeogr 28, 511–523 (2001).
750  65.  B. Perez-Lamarque, H. Morlon, Characterizing symbiont inheritance during host–
751       microbiota evolution: Application to the great apes gut microbiota. Mol Ecol Resour 19,
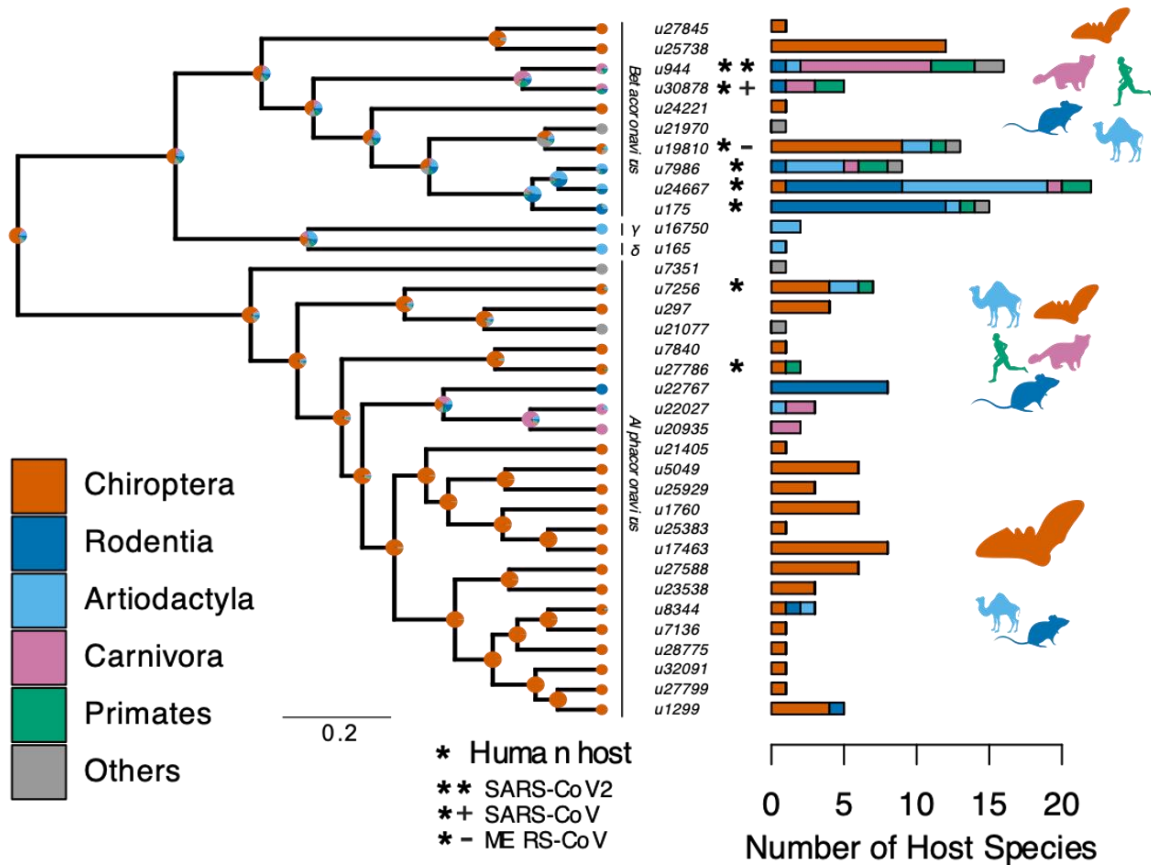752       1659–1671 (2019).

753    66.    D. M., de Vienne, T. Giraud, J. A. Shykoff, When can host shifts produce congruent host
754            and parasite phylogenies? A simulation approach. J Evol Biol 20, 1428–1438 (2007).
755    67.    B. Perez-Lamarque, H. Krehenwinkel, R. G. Gillespie, H. Morlon, Limited evidence for
756            microbial transmission in the phylosymbiosis between Hawaiian spiders and their
757            microbiota. mSystems 7, e01104-21 (2022).
758    68.    C. J. Marsh, et al., Expert range maps of global mammal distributions harmonised to
759            three taxonomic authorities. J Biogeogr 49, 979–992 (2022).
760    69.    C. J. Burgin, J. P. Colella, P. L. Kahn, N. S. Upham, How many species of mammals are
761            there? J Mammal 99, 1-14 (2018).
762    70.    A. H. Hurlbert, W. Jetz, Species richness, hotspots, and the scale dependence of range
763            maps in ecology and conservation. PNAS 104, 13384-13389 (2007).
764    71.    D. P. Faith, Conservation evaluation and phylogenetic diversity. Biol Cons 61, 1-10
765            (1992).
766    72.    P. Pappalardo, et al., Comparing methods for mapping global parasite diversity. Global
767            Ecol Biogeogr 29, 182-193 (2019).
768    73.    P. O. Title, D. L. Swiderski, M. L. Zelditch, EcoPhyloMapper: an R package for integrating
769            geographical ranges, phylogeny and morphology. Methods Ecol Evol 13, 1912-1922
770            (2022).
771    74.    E. Pebesma, Simple features for R: standardized support for spatial vector data. The R
772            Journal 10, 439-446 (2018).
773    75.    E. Paradis, K. Schliep, ape 5.0: an environment for modern phylogenetics and
774            evolutionary analyses in R. Bioinformatics 35, 526-528 (2019).
775

776     **Figures and Tables**



777
778
779     **Figure 1. A framework for testing scenarios of virus-host evolution, illustrated with the**
780     **example of Coronaviridae and their mammalian hosts:** In (A), a scenario of ancient origination
781     and codiversification; in (B) a scenario of recent origination and diversification by preferential host
782     switches; and in (C) a scenario of independent evolution. For each scenario, we indicate the
783     associated predictions in the grey boxes. Contrary to scenario C, both scenarios A and B are
784     expected to generate a cophylogenetic signal, *i.e.* closely-related coronaviruses tend to infect
785     closely-related mammals, resulting in significant reconciliations when using topology-based
786     probabilistic cophylogenetic methods, such as the undated version of ALE (Szollozi et al 2013),
787     Jane (Conow et al 2010), or eMPRess (Sanitchaivekin et al, 2021). However, we expect scenario
788     B to be distinguishable from scenario A in terms of the time consistency of host-switching events.
789     Under scenario B, cophylogenetic methods wrongly estimate a combination of cospeciations and
790     "back-in-time" host switches (see Methods & Results). We also expect different biogeographic
791     patterns under the different scenarios.
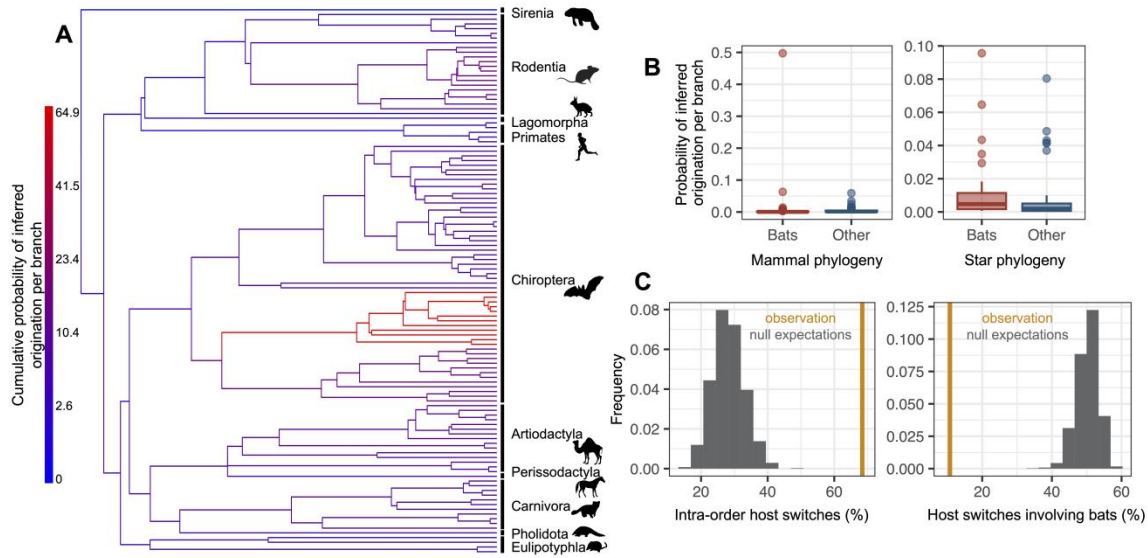792

793

**Figure 2. Species-level relationships among coronaviruses and their associated mammalian hosts.** A consensus Maximum Clade Credibility phylogenetic tree of coronaviruses is shown on the left. The tree was constructed with BEAST2 based on 150-aa palmprint amino acid sequences of the RdRp gene. sOTUs of Coronaviridae followed the definition of the Serratus project. The putative location of four genera of coronaviruses, Beta, Gamma, Delta, and Alphacoronaviruses, is shown. Bar scale is in units of aa substitution. On the right, a barplot gives the number of total mammalian host species and the number of host species by main mammalian order. Ancestral states on the left were obtained for illustrative purposes with the make.simmap function of the phytools R package (Revell 2012). Mammal silhouettes taken from open-to-use sources in phylopic.org, detailed credits given in SI Appendix Table S6.

804
805

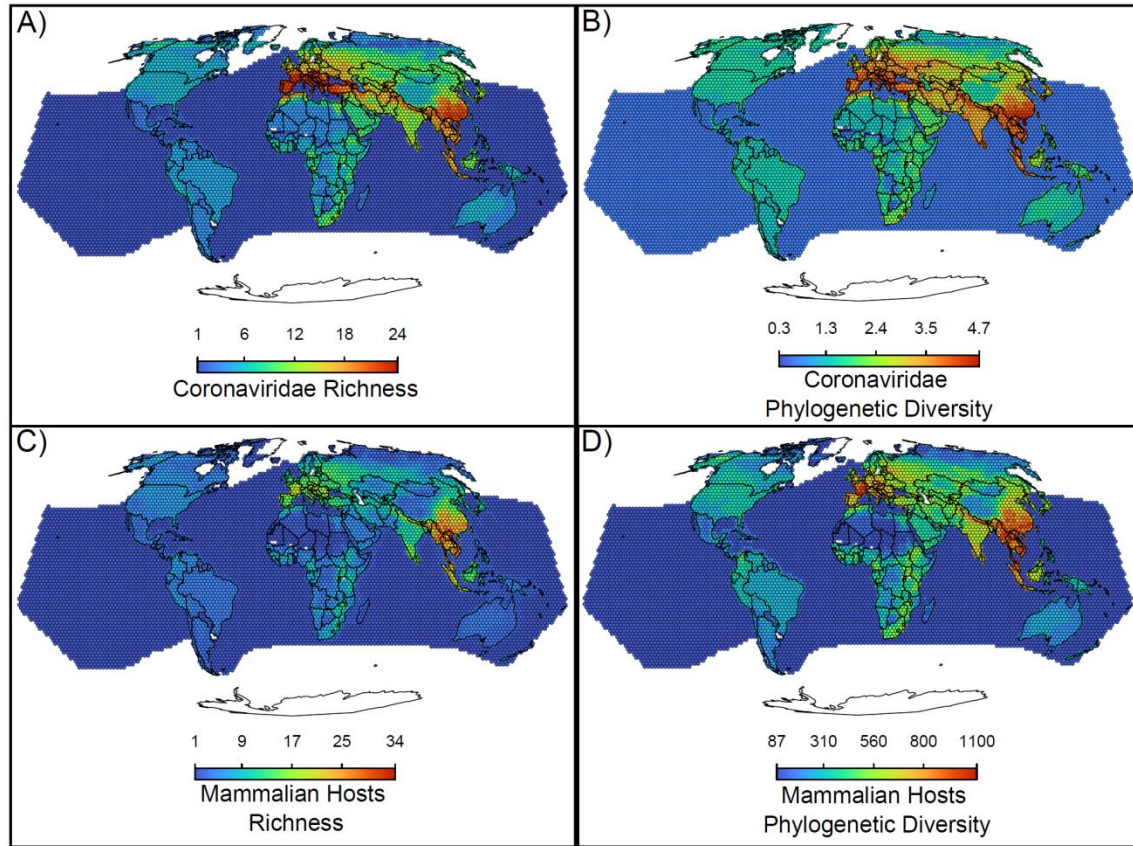**Figure 3. The origination of coronaviruses in mammals is estimated among bats, which tend to form a closed reservoir.** (A) Phylogenetic tree of the mammals with branches colored as the percentage of ALE reconciliations which inferred this branch or its ancestral lineages as the origination of coronaviruses in mammals. Red branches are likely originations, whereas blue branches are unlikely. (B) Boxplots recapitulating the probability of inferred origination per branch in bats *versus* other mammal orders, with ALE applied on the original mammal tree (left panel) or on the mammal tree transformed into a star phylogeny (right panel), therefore assuming an origination in extant species. (C) Distributions of the percentages of host switches occurring within mammalian orders (left panel) and between-orders involving bats (right panel). Observed values (in orange) are compared to null expectations if host switches were happening at random (in grey). Mammal silhouettes taken from open-to-use sources in phylopic.org, detailed credits given in SI Appendix Table S6.

820



**Figure 4. Maps of the diversity of coronaviruses and their mammal hosts.** In A) the richness of species of coronaviruses; geographic range maps of coronaviruses were constructed after applying the host-filling method on the geographic range maps of mammalian hosts of coronaviruses. In B) Faith's (1992) phylogenetic diversity of coronaviruses, calculated using the phylogenetic tree of coronaviruses (see main text). In C) and D), the richness and phylogenetic diversity of mammal hosts of coronaviruses, respectively. All maps are on the Mollweide projection.