

The coevolutionary mosaic of bat-betacoronaviruses spillover risk

Timothée Poisot^{1,2,‡} Peregrin Took^{3,4} Merriadoc Brandybuck^{5,4,‡}

¹ Université de Montréal ² Québec Centre for Biodiversity Sciences ³ Inn of the Prancing Pony

⁴ Fellowship of the Ring ⁵ Green Dragon Inn

‡ These authors contributed equally to the work

Correspondance to:

Timothée Poisot — timothee.poisot@umontreal.ca

This work is released by its authors under a CC-BY 4.0 license



Last revision: *April 6, 2022*

Purpose: This template provides a series of scripts to render a markdown document into an interactive website and a series of PDFs.

Motivation: It makes collaborating on text with GitHub easier, and means that we never need to think about the output.

Internals: GitHub actions and a series of python scripts. The markdown is handled with pandoc.

1 Spillover risk is not unidimensional. From the standpoint of an animal community, i.e. a pool of possible
2 hosts, there are a multiplicity of ecological factors that come into play (Plowright et al. 2017). The global
3 richness of hosts is one such component commonly mentioned/analysed (Anthony et al. 2017), but there
4 is an argument to be made that species who are not competent hosts of a specific virus genus may not
5 factor into this (Plowright et al. 2015), or that species who are assumed to share viruses at different rates
6 should be weighted accordingly (Albery et al. 2020). In mammals, key functional traits (for which
7 phylogeny is a reasonable proxy) are determinants of the spillover potential (Olival et al. 2017); these
8 include, notably, TK. Finally, especially when the pool of potential hosts spans the entire globe, there may
9 be local host pools that are highly unique; not having been observed in other locations, these can act on
10 the overall risk either by providing novel contact opportunities, reflecting unique host-environment
11 combinations (Engering et al. 2013), or facilitating rapid changes in specialism (Agosta et al. 2010). In the
12 specific case of generalist pathogens (as betacoronavirus clearly are), there is conceptual and empirical
13 support to the idea that these community- level mechanisms are even more important in driving the
14 overall risk (Power and Mitchell 2004).

15 Bats are important reservoir hosts for different classes of microorganisms (Chu 2008, Donaldson 2010, Li
16 2010), some of which can threaten human health. Especially concerning is the fact that bats are reservoirs
17 for a variety of emerging viruses (Calisher 2006), making balancing the needs for bat conservation and
18 disease prevention a potentially difficult act, especially in more densely populated areas (REF).

19 Chiropterans emerged around 64 million years ago and are one of the most diverse mammalian orders,
20 with an estimated richness of more than 12000 species, (Peixoto F et al, 2018) and 14325 known species
21 Simmons & Cirranello. They exhibit a broad variety of habitat use, behaviour, and feeding strategies,
22 resulting in their playing an essential role in the delivery of several ecosystem services (Kasso 2013),
23 including economic benefits. Over two-thirds of bats are either obligate or facultative insectivorous
24 mammals, therefore playing an important role in the regulation of insect pests that can affect crops
25 (Williams-Guillen 2011), and vectors of diseases that put a risk on human health (Gonsalves 2013).

26 Because bats are globally distributed and have a long evolutionary history, phylogeographic and
27 biogeographic approaches are required to shed light on the extant distribution of coevolutionary processes
28 between bats and the pathogens they carry. As a consequence, not all areas are facing a risk of human
29 spillover, and those who do might not be facing risks of the same nature and magnitude.

30 Yet a comprehensive assessment of the risk of spillover of betacoronaviruses from bat hosts to humans is

limited by the fact that we do not know the full diversity of viruses associated with every bat species. Predictive models can help fill in some of these gaps, by recommending hosts based on known host-virus associations BECKER REF. In this paper, we examine the biogeographic structure of bats-betacoronaviruses associations, based on a curated dataset of known and predicted hosts. We turn these associations into a spatially explicit additive mapping of zoonotic risk components, which reveals extreme heterogeneity of risk at the global scale; furthermore, we identify the Amazon and South-Eastern Asia as hotspots of phylogenetic distinctiveness of betacoronaviruses; surprisingly, current data suggest that viral sharing between hosts is high in the Amazon and low in South-Eastern Asia, which has the potential to result in different evolutionary dynamics between these two regions.

TK summary of the results

Methods

Known betacoronavirus hosts

We downloaded the CoV reservoir database from <https://www.viralemergence.org/betacov> on Aug. 2021. This database was assembled by a combination of data mining, literature surveys, and application of an ensemble recommender system classifying hosts as either “Suspected” or “Unlikely” (REF BECKER). The hosts considered for this study were all hosts with a known record of a betacoronavirus, and all those with a “Suspected” status in the ensemble model. This resulted in a list of TK TP unique host species.

Bats occurrences

We downloaded the rangemap of every extant bat species that was either classified as an empirically documented or a suspected host of beta-coronaviruses (Becker et al. 2020), according to recent IUCN data (IUCN 2021). The range maps were subsequently rasterized at a resolution of approximately TK TP. For every pixel in the resulting raster where at least one bat host of betacoronavirus was present, we extract the species pool, which was used to calculate the following risk assessment components: phylogenetic diversity, bat compositional uniqueness, and predicted viral sharing risk.

55 **Bats phylogeography**

56 For every pixel, we measured Faith's Phylogenetic Diversity (Faith 1992) based on a recent synthetic tree
57 with robust time calibration, covering about 6000 mammalian species (Upham et al. 2019). Faith's PD
58 measures the sum of unique branches from an arbitrary root to a set of tips, and comparatively larger
59 values indicate a more phylogenetic diverse species pool. We measured phylogenetic diversity starting
60 from the root of the entire tree (and not from Chiroptera); this bears no consequences on the resulting
61 values, since all branches leading up to Chiroptera are only counted one per species pool, and (as we
62 explain when describing the assembly of the composite risk map), all individual risk components are
63 ranged in [0,1]. This measure incorporates a richness component, which we chose not to correct for; the
64 interpretation of the phylogenetic diversity is therefore a weighted species richness, that accounts for
65 phylogenetic over/under-dispersal in some places.

66 **Bats compositional uniqueness**

67 For every species pool, we measured its Local Contribution to Beta-Diversity (Legendre and De Cáceres
68 2013); LCBD works from a species-data matrix (traditionally noted as Y), where species are rows and sites
69 are columns, and a value of 1 indicates occurrence. We extracted the Y matrix assuming that every pixel
70 represents a unique location, and following best practices (Legendre and Condit 2019) transformed it
71 using Hellinger's distance to account for unequal bat richness at different pixels. The correction of raw
72 community data is particularly important for two reasons: first, it prevents the artifact of richer sites
73 having higher importance; second, it removes the effect of overall species richness, which is already
74 incorporated in the phylogenetic diversity component. High values of LCBD indicate that the pixel has a
75 community that is on average more dissimilar in species composition than what is expected knowing the
76 entire matrix, i.e. a more unique community.

77 **Viral sharing between hosts**

78 For all bat hosts of betacoronaviruses, we extracted their predicted viral sharing network (Albery et
79 al. 2020). This network stores pairwise values of viral community similarity. To project viral sharing values
80 into a single value for every pixel, we averaged the pairwise scores. High values of the average sharing
81 propensity means that this specific extant bat assemblage is likely to be proficient at exchanging viruses.

82 **Composite risk map**

83 To visualize the aggregated risk at the global scale, we combine the three individual risk components
84 (phylogenetic diversity, compositional uniqueness, and viral sharing) using an additive color model
85 (Seekell et al. 2018). In this approach, every risk component gets assigned a component in the RGB color
86 model (phylogenetic diversity is green, compositional uniqueness is red, and viral sharing is blue). In
87 order to achieve a valid RGB measure, all components are re-scaled to the [0,1] interval. This additive
88 model conveys both the intensity of the overall risk, but also the nature of the risk as colors diverge
89 towards combinations of values for three risk components.

90 **Viral phylogeography**

91 We used the following query to pull all betacoronavirus sequence data from the GenBank Nucleotide
92 database except SARS-CoV-2; (“Betacoronavirus”[Organism] OR betacoronavirus[All Fields]) NOT
93 (“Severe acute respiratory syndrome coronavirus 2”[Organism] OR sars-cov-2[All Fields]). We added a
94 single representative sequence for SARS-CoV-2 and manually curated to remove sequences without the
95 RNA-dependent RNA polymerase (RdRp) sequence or that contained words indicating recombinant or
96 laboratory strains including “patent,” “mutant,” “GFP,” and “recombinant.” We filtered over-represented
97 taxa including betacoronavirus 1, hCoV-OC43, Middle East respiratory syndrome coronavirus, Murine
98 hepatitis virus, and hCoV-HKU1. Curated betacoronavirus RdRp sequences were then aligned using
99 MAFFT v 1.4.0 (Kato and Standley 2013, Supplemental X) and a maximum likelihood tree reconstructed
100 in IQ-TREE v 1.6.12 (Nguyen et al. 2015) with ModelFinder (Kalyaanamoorthy et al. 2017) ultrafast
101 bootstrap approximation (Hoang et al. 2018) and the following parameters (STEPH WILL ADD,
102 Supplemental X).

103 **Viral evolutionary diversification**

104 We first tested the hypothesis that hotspots of viral diversification would track hotspots of bat
105 diversification. To do so, we plotted the number of known bat hosts (specifically only those included in the
106 phylogeny, so there was a 1:1 correspondence between data sources) against the “mean evolutionary
107 distinctiveness” of the associated viruses. To calculate this, we derived the fair proportions evolutionary
108 distinctiveness (Isaac et al., 2007) for each of the viruses in the tree, then averaged these at the bat species

level, projected these values onto their geographic distributions, and averaged across every bat found in a given pixel. As such, this can be thought of as a map of the mean evolutionary distinctiveness of the known viral community believed to be associated with a particular subset of bats present.

Co-distribution of hosts and viral hotspots

Subsequently, we tested the hypothesis that the biogeography of bat betacoronaviruses should track the biogeography of their hosts. To test this idea, we loosely adapted a method from Kreft & Jetz (2010), who proposed a phylogenetic method for the delineation of animal biogeographic regions. In their original method, a distance matrix - where each row or column represents a geographic raster's grid cell, and the dissimilarity values are the "beta diversity similarity" of their community assemble - undergoes non-metric multidimensional scaling (NMDS); the first two axes of the NMDS are projected geographically using a four-color bivariate map.

Here, we build on this idea with an entirely novel methodology. First, we measure the phylogenetic distance between the different viruses in the betacoronavirus tree by using the cophenetic function in 'ape'; subsequently, we take a principal components analysis of that distance matrix (readily interchangeable for NMDS in this case) to project the viral tree into an n-dimensional space. We then take the first two principal components and, as with the evolutionary distinctiveness analysis, aggregated these to a mean host value and projected them using a four-color bivariate map.

Outbreaks data geo-referencing

Finally, we provide a summary visualization of what available information describes the spillover of zoonotic betacoronaviruses of bat origin where data was available before and up through the COVID-19 pandemic. The SARS-CoV-2 outbreak was georeferenced to the initial case cluster in Wuhan, China; SARS-CoV was georeferenced based on the cave with the closest known viruses circulating in nature (Hu et al. 2017 PLoS Pathogens), and a nearby location where serological (antibody) evidence has indicated human exposure to SARS-like viruses (Wang et al. 2018 Virologica Sinica). For MERS-CoV, we presented the index cases available from a recently-published compendium of MERS-CoV cases (Ramshaw et al. 2019); these are largely if not all presumed to be camel-to-human transmission, and the precise origin point of MERS-CoV in bats is uncertain. Not shown is a recent case of a recombinant canine coronavirus

136 that showed the ability to infect humans, both because this study was published after the beginning of the
137 COVID-19 pandemic and because bats' involvement in this cycle of transmission has been marginal to
138 non-existent.

139 **Results**