

The evolutionary history of ACE2 usage within the coronavirus subgenus *Sarbecovirus*

H. L. Wells,^{1,*} M. Letko,^{2,3} G. Lasso,⁴ B. Ssebide,⁵ J. Nziza,⁵ D. K. Byarugaba,^{6,7} I. Navarrete-Macias,⁸ E. Liang,⁸ M. Cranfield,^{9,10} B.A. Han,¹¹ M. W. Tingley,¹² M. Diuk-Wasser,¹ T. Goldstein,⁹ C. K. Johnson,⁹ J. A. K. Mazet,⁹ K. Chandran,⁴ V. J. Munster,³ K. Gilardi,^{6,9} and S. J. Anthony^{13,*}

¹Department of Ecology, Evolution, and Environmental Biology, Columbia University, 1200 Amsterdam Ave, New York, NY 10027, USA, ²Laboratory of Virology, Division of Intramural Research, National Institute of Allergy and Infectious Diseases, National Institutes of Health, 903 S. 4th St, Hamilton, MT 59840, USA, ³Paul G. Allen School for Global Animal Health, Washington State University, 1155 College Ave, Pullman, WA 99164, USA, ⁴Department of Microbiology and Immunology, Albert Einstein College of Medicine, 1300 Morris Park Ave, Bronx, NY 10462, USA, ⁵Gorilla Doctors, c/o MGVP, Inc., 1089 Veterinary Medicine Drive, Davis, CA 95616, USA, ⁶Makerere University Walter Reed Project, Plot 42, Nakasero Road, Kampala, Uganda, ⁷Makerere University, College of Veterinary Medicine, Living Stone Road, Kampala, Uganda, ⁸Center for Infection and Immunity, Mailman School of Public Health, Columbia University, 722 W 168th St, New York, NY 10032, USA, ⁹One Health Institute and Karen C. Drayer Wildlife Health Center, School of Veterinary Medicine, University of California Davis, 1089 Veterinary Medicine Drive, Davis, CA 95616, USA, ¹⁰Department of Microbiology and Immunology, University of North Carolina School of Medicine, 125 Mason Farm Road, Chapel Hill, NC 27599, USA, ¹¹Cary Institute of Ecosystem Studies, 2801 Sharon Turnpike, Millbrook, NY 12545, USA, ¹²Department of Ecology and Evolutionary Biology, University of California Los Angeles, 612 Charles E. Young Drive South, Los Angeles, CA 90095, USA and ¹³Department of Pathology, Microbiology, and Immunology, School of Veterinary Medicine, University of California Davis, One Shields Avenue, Davis, CA 95616, USA

*Corresponding authors: E-mails: hlw2124@columbia.edu, sjanthony@ucdavis.edu

[†]<https://orcid.org/0000-0002-1724-5843>

Abstract

Severe acute respiratory syndrome coronavirus 1 (SARS-CoV-1) and SARS-CoV-2 are not phylogenetically closely related; however, both use the angiotensin-converting enzyme 2 (ACE2) receptor in humans for cell entry. This is not a universal sarbecovirus trait; for example, many known sarbecoviruses related to SARS-CoV-1 have two deletions in the receptor binding domain of the spike protein that render them incapable of using human ACE2. Here, we report three sequences of a novel sarbecovirus from Rwanda and Uganda that are phylogenetically intermediate to SARS-CoV-1 and SARS-CoV-2 and demonstrate via in vitro studies that they are also unable to utilize human ACE2. Furthermore, we show that the observed pattern of ACE2 usage among sarbecoviruses is best explained by recombination not of SARS-CoV-2, but of SARS-CoV-1 and its relatives. We show that the lineage that includes SARS-CoV-2 is most likely the ancestral ACE2-using lineage, and that recombination with at least one virus from this group conferred ACE2 usage to the lineage including SARS-CoV-1 at some

© The Author(s) 2021. Published by Oxford University Press.

This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact journals.permissions@oup.com

time in the past. We argue that alternative scenarios such as convergent evolution are much less parsimonious; we show that biogeography and patterns of host tropism support the plausibility of a recombination scenario, and we propose a competitive release hypothesis to explain how this recombination event could have occurred and why it is evolutionarily advantageous. The findings provide important insights into the natural history of ACE2 usage for both SARS-CoV-1 and SARS-CoV-2 and a greater understanding of the evolutionary mechanisms that shape zoonotic potential of coronaviruses. This study also underscores the need for increased surveillance for sarbecoviruses in southwestern China, where most ACE2-using viruses have been found to date, as well as other regions such as Africa, where these viruses have only recently been discovered.

Key words: virus evolution; viral ecology; recombination; coronavirus.

1. Introduction

The recent emergence of *severe acute respiratory syndrome coronavirus 2* (SARS-CoV-2) in China and its rapid spread around the world demonstrates that coronaviruses (CoVs) from wildlife remain an urgent threat to global public health and economic stability. In particular, coronaviruses from the subgenus *Sarbecovirus* (which includes SARS-CoV-2, SARS-CoV-1, numerous bat viruses, and a small number of pangolin viruses) (Lefkowitz et al. 2018) are considered to be a high-risk group for potential emergence. As both sarbecoviruses that have caused human disease (SARS-CoV-1 and -2) use angiotensin-converting enzyme 2 (ACE2) as their cellular receptor (Wenhui Li et al. 2003; Zhou et al. 2020), the evolution of this trait is of particular importance for understanding the emergence pathway for sarbecoviruses. Bat SARS-like coronavirus Rp3 is a phylogenetically close relative of SARS-CoV-1 but is unable to bind human ACE2 (hACE2) *in vitro* (Ren et al. 2008). In contrast, other close relatives of SARS-CoV-1, including bat SARS-like coronavirus WIV1 and WIV16, do have the capacity to bind hACE2 (Ge et al. 2013; Yang et al. 2016). A number of other SARS-CoV-1-like viruses have also been tested for their ability to utilize hACE2 (Hu et al. 2017; Letko, Marzi, and Munster 2020; Menachery et al. 2015) and comparison of their spike protein sequences shows that viruses that are unable to utilize hACE2 unanimously have one or two deletions in their receptor binding domains (RBDs) that make them structurally very different from those that do use hACE2 (Letko, Marzi, and Munster 2020). As SARS-CoV-1, Rp3, WIV1, and WIV16 viruses are closely phylogenetically related, the evolutionary mechanism explaining the variation in their ability to utilize hACE2 (and likely also bat ACE2) as a cellular receptor has thus far been unclear.

Chinese horseshoe bats (*Rhinolophidae*) are thought to be the primary natural reservoir of sarbecoviruses (Lau et al. 2005; Li et al. 2005; Ge et al. 2013; He et al. 2014; Hu et al. 2017). Bats within this family are also considered to be the source of the progenitor virus to SARS-CoV-1, as related viruses with high sequence identity to SARS-CoV-1 have been sequenced from *Rhinolophid* bats, although none have high sequence similarity to SARS-CoV-1 across the entire genome (Hu et al. 2017; Hon et al. 2008). It is hypothesized that SARS-CoV-1 obtained genomic regions from different strains of bat SARS-1-like CoVs in or near Yunnan Province by recombination before spilling over into humans (Hon et al. 2008; Hu et al. 2017; Luk et al. 2019). In particular, one region of SARS-CoV-1 that is known to have a recombinant origin is the spike gene, as a breakpoint has been detected at the junction of ORF1b and the spike (Hon et al. 2008; Lau et al. 2010). The SARS-1-CoV spike is genomically very different from other viruses in the same clade that have large deletions in the RBD and are unable to use hACE2. The exact minor parent that contributed the recombinant region is still

unknown, but it was previously hypothesized that the recombination occurred with a yet undiscovered lineage of sarbecoviruses and that this event contributed strongly to its potential for emergence (Hon et al. 2008; Yuan et al. 2010). Recombination has also been shown within the spike genes of other CoVs that have spilled over into humans and domestic animals and is potentially an important driver of emergence for all coronaviruses (Woo et al. 2009; Graham and Baric 2010; Lu, Wang, and Gao 2015; Su et al. 2016; Menachery, Graham, and Baric 2017; Anthony et al. 2017a).

In order for CoVs to recombine, they must first have the opportunity to do so by sharing overlapping geographic ranges, host species tropism, and cell and tissue tropism. Sarbecoviruses in bats tend to phylogenetically cluster according to the geographic region in which they were found (Hu et al. 2017; Yu et al. 2019). Yu et al. (2019) showed that there are three lineages of SARS-CoV-1-like viruses: Lineage 1 from southwestern China (Yunnan, Guizhou, and Guangxi, and including SARS-CoV-1), Lineage 2 from other southern regions (Guangdong, Hubei, Hong Kong, and Zhejiang), and Lineage 3 from central and northern regions (Hubei, Henan, Shanxi, Shaanxi, Hebei, and Jilin). Studies in Europe and Africa have shown that there are distinct sarbecovirus clades in each of these regions as well, herein named 'Lineage 4' (Ar Gouilh et al. 2018; Lecis et al. 2019; Drexler et al. 2010; Rihtarić et al. 2010; Lelli et al. 2013; Tao and Tong 2019). Sarbecoviruses appear to switch easily among co-occurring *Rhinolophus* species (Cui et al. 2007; Leopardi et al. 2018); however, they appear to rarely occupy more than one geographic area, despite the fact that some of these bat species have widespread distributions across China.

Shortly after the emergence of SARS-CoV-2, Zhou et al. (2020) showed a high degree of homology across the genome between a bat virus (RaTG13) sampled from Yunnan Province in 2013 and SARS-CoV-2. RaTG13 has also been shown to bind hACE2, although with decreased affinity compared to SARS-CoV-2 (Shang et al. 2020). Subsequently, seven full- or near full-length SARS-CoV-2-like viruses were published that had been sampled from Malayan pangolins (*Manis javanica*) in 2017 and 2019 (Liu, Chen, and Chen 2019; Lam et al. 2020), one of which has also been tested and found to bind hACE2 (Wrobel et al. 2021). Neither SARS-CoV-2, RaTG13, nor the pangolin CoVs have deletions in their RBDs. In contrast, the most recently described bat virus (RmYN02) is even more closely related to SARS-CoV-2 than RaTG13 in the polymerase gene and was also found in Yunnan Province; however, this sequence has deletions in the RBD and homology modeling suggests it likely does not use hACE2 (Zhou et al. 2020). Together, these viruses form a fifth phylogenetic lineage ('Lineage 5') that is distinct from all other lineages of sarbecoviruses despite having been detected in

Yunnan, where all viruses found until this point had belonged to Lineage 1.

This finding of overlapping Lineage 1 and Lineage 5 viruses in geographic space is inconsistent with the previously observed pattern of biogeography for sarbecoviruses. SARS-CoV-2 was isolated first from people in Hubei Province and one of the pangolin viruses was isolated from an animal sampled in Guangdong, neither of which are Lineage 1 provinces. However, the true geographic origins of these viruses are unknown as it is possible that they were anthropogenically transported to the regions in which they were detected. For example, the Malayan pangolin (*Manis javanica*) has a natural range that reaches southwestern China (Yunnan Province) at its northernmost edge and extends further south into Myanmar, Lao PDR, Thailand, and Vietnam (Challender et al. 2014). So, if they were naturally infected (as opposed to infection via wildlife trade), the infection was potentially not acquired from Guangdong Province. Similarly, SARS-CoV-2 cannot be guaranteed to have emerged from bats in Hubei Province, as humans are highly mobile and the exact spillover event was not observed. If the clade containing SARS-CoV-2 and its close relatives is indeed endemic in animals in Yunnan and the nearby Southeast Asian regions as suggested by the presence of RaTG13, RmYN02, and the natural range of the Malayan pangolin, whatever mechanism is facilitating the biogeographical concordance of Lineages 1, 2, and 3 within China appears to no longer apply for the biogeography of Lineage 5, since they all appear to overlap in and around Yunnan Province.

Here, we report a series of observations that together suggest that SARS-CoV-1 and its close relatives gained the ability to utilize ACE2 through a recombination event that happened between an ancestor of SARS-CoV-1 and a Lineage 5 virus phylogenetically related to SARS-CoV-2, which could only have occurred with the lineages occupying the same geographic and host space. We also report three full-length genomes of sarbecoviruses from Rwanda and Uganda and demonstrate that the RBDs of these viruses are genetically intermediate between viruses that use ACE2 and those that do not. Accordingly, we also investigate the potential for these viruses to utilize hACE2 *in vitro*. Together, our findings help illuminate the evolutionary history of ACE2 usage within sarbecoviruses and provide insight into identifying their risk of emergence in the future. We also propose a mechanism that could explain the pattern of phylogeography across Lineages 1, 2, and 3, and why Lineage 5 viruses (including SARS-CoV-2 and its relatives) represent an inconsistency to this pattern.

2. Methods

2.1 Consensus polymerase chain reaction and sequencing of sarbecoviruses from Africa

Oral swabs, rectal swabs, whole blood, and urine samples collected from bats sampled and released in Uganda and Rwanda were assayed for CoVs using consensus polymerase chain reaction (PCR) as previously described (Anthony et al. 2017a). All sampling was conducted under UC Davis IACUC Protocol No. 16048. Bands of the expected size were purified and confirmed positive by Sanger sequencing, and the PCR fragments were deposited to GenBank (accessions MT738926-MT738928, MT732776). Samples were subsequently deep sequenced using the Illumina HiSeq platform and reads were bioinformatically de novo assembled using MEGAHIT v1.2.8 (Li et al. 2016) after quality control steps and subtraction of host reads using

Bowtie2 v2.3.5. Contigs were aligned to a reference sequence and any overlaps or gaps were confirmed with iterative local alignment using Bowtie2. The full genome sequences are deposited in GenBank. Cytochrome b, cytochrome oxidase I, and ACE2 host sequences were also extracted bioinformatically where possible by mapping reads to *Rhinolophus ferrumequinum* reference genes using Bowtie2 and deposited in GenBank.

2.2 Phylogenetic reconstruction

All publicly available full genome sarbecovirus sequences were collected from GenBank and SARS-CoV-2, pangolin virus genomes, RaTG13, and RmYN01/RmYN02 were downloaded from GISAID (Table 1). All relevant metadata (geographic origin, host species, date of collection) was retrieved from GenBank or the corresponding publications. The RdRp gene (nucleotides 13,431 to 16,222 based on SARS-CoV-2 sequence EPI_ISL_402125 from GISAID) and RBD region (nucleotides 22,506–23,174 based on the same SARS-CoV-2 reference genome) were extracted and aligned using Muscle v10.2.6. We chose RdRp as a backbone to which to compare because of the strong evolutionary constraints imposed by its fundamental biological role in viral replication (Ulferts et al. 2010). Indeed, the RdRp is generally considered to be a primary genetic trait in viral taxonomy (ICTV 2019; Gorbalenya et al. 2020) and most viruses exhibit strong purifying selection in this gene (Tang et al. 2009). Further, the orf1ab region of coronaviruses (which contains the RdRp) also tends to be more recombination-free as compared to the recombination-frequent latter half of the genome (Fu and Baric 1994; Boni et al. 2020). Since many of our conclusions are based around phylogenetic topology, we confirmed the robustness of the topology of our nucleotide trees by also building identical trees with alignments of other relatively stable genes in orf1ab frequently used for taxonomic classification (Gorbalenya et al. 2020) (Supplementary Fig. S1). Phylogenetic reconstruction was performed using BEAST v2.6.3 (Bouckaert et al. 2019) with partitioned codon positions, a GTR + Γ substitution model for each of the three codon positions, a constant size coalescent process prior, and a strict molecular clock model. Log files were examined using Tracer v1.7.1 to confirm that the model converged and that the effective sample size for each parameter was at least 100. Chains were run until these convergence criteria were met (~2–10 million samples) and multiple chains were run independently to ensure convergence to the same estimates. Use of Beagle 2.1.2 was chosen to increase computational speed.

Maximum clade credibility trees were built using TreeAnnotator and visualized with FigTree with branches scaled by distance. Posterior probabilities are shown on the preceding branch for each node and probabilities for nodes near the tips of the tree were removed for visual clarity as the exact reconstruction of the most recent divergence events are not within the scope of this study and bear no impact on the interpretation of evolutionary events deeper within the tree.

Finally, for time-calibrated phylogenies, we minimized the effect of recombination on our estimates by using regions of the genome that were free of recombination for the 13 Lineage 1 sequences of interest (further detailed below). In place of RdRp we used Region A, and in place of RBD we used Region E. These regions were determined to be completely breakpoint free for all sequences using 3SEQ. We started by adding tip dates to Region A and used a strict molecular clock with a normally distributed prior informed from estimates derived in Boni et al. (2020) (mean $5.5e-4$, sd $5.5e-5$). The prior distribution for the coalescent population size was set to lognormal with mean 1

Table 1. Full list of sequences and accession numbers used in this study.

| Accession | Name | Date | Country | Host | ACE2 usage |
|----------------|---|------|------------------|----------------------------------|--|
| AY304486 | SARS coronavirus SZ3 | 2003 | Guangdong, China | <i>Paguma larvata</i> (civet) | Li et al. (2005) ^a |
| AY304488 | SARS coronavirus SZ16 ^b | 2003 | Hong Kong, China | <i>Paguma larvata</i> (civet) | |
| AY572034 | SARS coronavirus civet007 ^b | 2004 | Guangdong, China | <i>Paguma larvata</i> (civet) | |
| DQ022305 | Bat SARS coronavirus HKU3 1 | 2005 | Hong Kong, China | <i>Rhinolophus sinicus</i> | |
| DQ071615 | Bat SARS coronavirus Rp3 | 2004 | Guangxi, China | <i>Rhinolophus pearsonii</i> | [Ren et al. 2008] ^a [Hu et al. 2017] ^a [Letko, Marzi, and Munster 2020] ^a |
| DQ084199 | Bat SARS coronavirus HKU3 2 ^b | 2005 | Hong Kong, China | <i>Rhinolophus sinicus</i> | |
| DQ084200 | Bat SARS coronavirus HKU3 3 ^b | 2005 | Hong Kong, China | <i>Rhinolophus sinicus</i> | |
| DQ412042 | Bat SARS coronavirus Rf1 | 2004 | Hubei, China | <i>Rhinolophus ferrumequinum</i> | Letko, Marzi, and Munster (2020) ^a |
| DQ412043 | Bat SARS coronavirus Rm1 | 2004 | Hubei, China | <i>Rhinolophus macrotis</i> | |
| DQ648856 | Bat coronavirus BtCoV/273/2005 | 2004 | Hubei, China | <i>Rhinolophus ferrumequinum</i> | Letko, Marzi, and Munster (2020) ^a |
| DQ648857 | Bat coronavirus BtCoV/279/2005 | 2004 | Hubei, China | <i>Rhinolophus macrotis</i> | Letko, Marzi, and Munster (2020) ^a |
| EPI_ISL_402125 | BetaCoV/Wuhan Hu 1 | 2019 | Hubei, China | human | Zhou et al. (2020) |
| EPI_ISL_402131 | BetaCoV/RaTG13 | 2013 | Yunnan, China | <i>Rhinolophus affinis</i> | Shang et al. (2020) ^a |
| EPI_ISL_412976 | BetaCoV/RmYN01 | 2019 | Yunnan, China | <i>Rhinolophus malayanus</i> | |
| EPI_ISL_412977 | BetaCoV/RmYN02 | 2019 | Yunnan, China | <i>Rhinolophus malayanus</i> | |
| EPI_ISL_410538 | BetaCoV/P4L ^b | 2017 | Guangxi, China | <i>Manis javanica</i> (pangolin) | |
| EPI_ISL_410539 | BetaCoV/P1E ^b | 2017 | Guangxi, China | <i>Manis javanica</i> (pangolin) | |
| EPI_ISL_410540 | BetaCoV/P5L ^b | 2017 | Guangxi, China | <i>Manis javanica</i> (pangolin) | |
| EPI_ISL_410541 | BetaCoV/P5E ^b | 2017 | Guangxi, China | <i>Manis javanica</i> (pangolin) | |
| EPI_ISL_410542 | BetaCoV/P2V | 2017 | Guangxi, China | <i>Manis javanica</i> (pangolin) | |
| EPI_ISL_410543 | BetaCoV/P3B ^b | 2017 | Guangxi, China | <i>Manis javanica</i> (pangolin) | |
| EPI_ISL_410544 | BetaCoV/P2S | 2019 | Guangdong, China | <i>Manis javanica</i> (pangolin) | Wrobel et al. (2021) ^a |
| FJ588686 | Bat SARS coronavirus Rs672/2006 | 2006 | Guizhou, China | <i>Rhinolophus sinicus</i> | |
| GQ153539 | Bat SARS coronavirus HKU3 4 ^b | 2005 | Hong Kong, China | <i>Rhinolophus sinicus</i> | |
| GQ153540 | Bat SARS coronavirus HKU3 5 ^b | 2005 | Hong Kong, China | <i>Rhinolophus sinicus</i> | |
| GQ153541 | Bat SARS coronavirus HKU3 6 ^b | 2005 | Hong Kong, China | <i>Rhinolophus sinicus</i> | |
| GQ153542 | Bat SARS coronavirus HKU3 7 ^b | 2006 | Guangdong, China | <i>Rhinolophus sinicus</i> | |
| GQ153543 | Bat SARS coronavirus HKU3 8 | 2006 | Guangdong, China | <i>Rhinolophus sinicus</i> | Letko, Marzi, and Munster (2020) ^a |
| GQ153544 | Bat SARS coronavirus HKU3 9 ^b | 2006 | Hong Kong, China | <i>Rhinolophus sinicus</i> | |
| GQ153545 | Bat SARS coronavirus HKU3 10 ^b | 2006 | Hong Kong, China | <i>Rhinolophus sinicus</i> | |
| GQ153546 | Bat SARS coronavirus HKU3 11 ^b | 2007 | Hong Kong, China | <i>Rhinolophus sinicus</i> | |
| GQ153547 | Bat SARS coronavirus HKU3 12 | 2007 | Hong Kong, China | <i>Rhinolophus sinicus</i> | |
| GQ153548 | Bat SARS coronavirus HKU3 13 ^b | 2007 | Hong Kong, China | <i>Rhinolophus sinicus</i> | Letko, Marzi, and Munster (2020) ^a |
| GU190215 | Bat coronavirus BM48-31/BGR/2008 | 2008 | Bulgaria | <i>Rhinolophus blasii</i> | Letko, Marzi, and Munster (2020) ^a |
| JX993987 | Bat coronavirus Rp/Shaanxi2011 | 2011 | Shaanxi, China | <i>Rhinolophus pusillus</i> | Letko, Marzi, and Munster (2020) ^a |
| JX993988 | Bat coronavirus Cp/Yunnan2011 | 2011 | Yunnan, China | <i>Chaerephon plicatus</i> | Letko, Marzi, and Munster (2020) ^a |
| KC881005 | Bat SARS-like coronavirus RsSHC014 | 2012 | Yunnan, China | <i>Rhinolophus sinicus</i> | [Letko, Marzi, and Munster 2020] ^a [Menachery et al. 2015] ^a |
| KC881006 | | 2012 | Yunnan, China | <i>Rhinolophus sinicus</i> | |

(continued)

Table 1.. (continued)

| Accession | Name | Date | Country | Host | ACE2 usage |
|-----------|--|------|-----------------|----------------------------------|--|
| | Bat SARS-like coronavirus Rs3367 | | | | |
| KF294457 | SARS related bat coronavirus Longquan 140 | 2012 | Guizhou, China | <i>Rhinolophus monoceros</i> | Letko, Marzi, and Munster (2020)^a |
| KF367457 | Bat SARS-like coronavirus WIV1 | 2012 | Yunnan, China | <i>Rhinolophus sinicus</i> | |
| KF569996 | <i>Rhinolophus affinis</i> coronavirus LYRa11 | 2011 | Yunnan, China | <i>Rhinolophus affinis</i> | Letko, Marzi, and Munster (2020)^a |
| KF636752 | Bat Hp betacoronavirus/ Zhejiang2013 | 2013 | Zhejiang, China | <i>Hipposideros pratti</i> | |
| KJ473811 | Bat coronavirus BtRf BetaCoV/ JL2012 | 2012 | Jilin, China | <i>Rhinolophus ferrumequinum</i> | Letko, Marzi, and Munster (2020)^a |
| KJ473812 | Bat coronavirus BtRf BetaCoV/ HeB2013 | 2013 | Hebei, China | <i>Rhinolophus ferrumequinum</i> | Letko, Marzi, and Munster (2020)^a |
| KJ473813 | Bat coronavirus BtRf BetaCoV/ SX2013 | 2013 | Shanxi, China | <i>Rhinolophus ferrumequinum</i> | |
| KJ473814 | Bat coronavirus BtRs BetaCoV/ HuB2013 | 2013 | Hubei, China | <i>Rhinolophus sinicus</i> | Letko, Marzi, and Munster (2020)^a |
| KJ473815 | Bat coronavirus BtRs BetaCoV/ GX2013 | 2013 | Guangxi, China | <i>Rhinolophus sinicus</i> | Letko, Marzi, and Munster (2020)^a |
| KJ473816 | Bat coronavirus BtRs BetaCoV/ YN2013 | 2013 | Yunnan, China | <i>Rhinolophus sinicus</i> | Letko, Marzi, and Munster (2020)^a |
| KP886808 | Bat SARS-like coronavirus YNLF 31C | 2013 | Yunnan, China | <i>Rhinolophus sinicus</i> | |
| KP886809 | Bat SARS-like coronavirus YNLF 34C | 2013 | Yunnan, China | <i>Rhinolophus sinicus</i> | |
| KT444582 | SARS-like coronavirus WIV16 | 2013 | Yunnan, China | <i>Rhinolophus sinicus</i> | Yang et al. (2016) |
| KU182964 | Bat coronavirus JTMC15 | 2013 | Yunnan, China | <i>Rhinolophus sinicus</i> | |
| KU182963 | Bat coronavirus MLHJC35 | 2012 | Jilin, China | <i>Rhinolophus sinicus</i> | |
| KU973692 | SARS related coronavirus F46 | 2012 | Yunnan, China | <i>Rhinolophus pusillus</i> | |
| KY352407 | SARS related coronavirus BtKY72 | 2007 | Kenya | <i>Rhinolophus</i> sp. | |
| KY417142 | Bat SARS-like coronavirus As6526 | 2014 | Yunnan, China | <i>Aselliscus stoliczkanus</i> | [Hu et al. 2017]^a [Letko, Marzi, and Munster 2020]^a |
| KY417143 | Bat SARS-like coronavirus Rs4081 | 2012 | Yunnan, China | <i>Rhinolophus sinicus</i> | [Hu et al. 2017]^a [Letko, Marzi, and Munster 2020]^a |
| KY417144 | Bat SARS-like coronavirus Rs4084 | 2012 | Yunnan, China | <i>Rhinolophus sinicus</i> | Letko, Marzi, and Munster (2020)^a |
| KY417145 | Bat SARS-like coronavirus Rf4092 | 2012 | Yunnan, China | <i>Rhinolophus ferrumequinum</i> | Letko, Marzi, and Munster (2020)^a |
| KY417146 | Bat SARS-like coronavirus Rs4231 | 2013 | Yunnan, China | <i>Rhinolophus sinicus</i> | [Hu et al. 2017]^a [Letko, Marzi, and Munster 2020]^a |
| KY417147 | Bat SARS-like coronavirus Rs4237 | 2013 | Yunnan, China | <i>Rhinolophus sinicus</i> | Letko, Marzi, and Munster (2020)^a |
| KY417148 | Bat SARS-like coronavirus Rs4247 | 2013 | Yunnan, China | <i>Rhinolophus sinicus</i> | Letko, Marzi, and Munster (2020)^a |
| KY417149 | Bat SARS-like coronavirus Rs4255 | 2013 | Yunnan, China | <i>Rhinolophus sinicus</i> | |
| KY417150 | Bat SARS-like coronavirus Rs4874 | 2013 | Yunnan, China | <i>Rhinolophus sinicus</i> | Hu et al. (2017) |
| KY417151 | Bat SARS-like coronavirus Rs7327 | 2014 | Yunnan, China | <i>Rhinolophus sinicus</i> | [Hu et al. 2017]^a [Letko, Marzi, and Munster 2020]^a |
| KY417152 | Bat SARS-like coronavirus Rs9401 | 2015 | Yunnan, China | <i>Rhinolophus sinicus</i> | |
| KY770858 | Bat coronavirus Anlong 103 | 2013 | Guizhou, China | <i>Rhinolophus sinicus</i> | |
| KY770859 | Bat coronavirus Anlong 112 | 2013 | Guizhou, China | <i>Rhinolophus sinicus</i> | |
| KY770860 | Bat coronavirus Jiyuan 84 | 2012 | Henan, China | <i>Rhinolophus ferrumequinum</i> | |

(continued)

Table 1.. (continued)

| Accession | Name | Date | Country | Host | ACE2 usage |
|-----------|---------------------------------------|------|-----------------|----------------------------------|---|
| KY938558 | Bat coronavirus 16BO133 | 2016 | South Korea | <i>Rhinolophus ferrumequinum</i> | |
| MG772933 | Bat SARS-like coronavirus SL CoVZC45 | 2017 | Zhejiang, China | <i>Rhinolophus sinicus</i> | Letko, Marzi, and Munster (2020) ^a |
| MG772934 | Bat SARS-like coronavirus SL CoVZXC21 | 2015 | Zhejiang, China | <i>Rhinolophus sinicus</i> | Letko, Marzi, and Munster (2020) ^a |
| MK211374 | Bat coronavirus BtRl BetaCoV/ SC2018 | 2018 | Sichuan, China | <i>Rhinolophus</i> sp. | |
| MK211375 | Bat coronavirus BtRs BetaCoV/ YN2018A | 2018 | Yunnan, China | <i>Rhinolophus affinis</i> | |
| MK211376 | Bat coronavirus BtRs BetaCoV/ YN2018B | 2018 | Yunnan, China | <i>Rhinolophus affinis</i> | |
| MK211377 | Bat coronavirus BtRs BetaCoV/ YN2018C | 2018 | Yunnan, China | <i>Rhinolophus affinis</i> | |
| MK211378 | Bat coronavirus BtRs BetaCoV/ YN2018D | 2018 | Yunnan, China | <i>Rhinolophus affinis</i> | |
| NC_004718 | SARS coronavirus | 2003 | Canada | human | Wenhui Li et al. (2003) |
| MT726044 | PREDICT PDF-2370 | 2013 | Uganda | <i>Rhinolophus</i> sp. | |
| MT726043 | PREDICT PDF-2386 | 2013 | Uganda | <i>Rhinolophus</i> sp. | |
| MT726045 | PREDICT PRD-0038 | 2010 | Rwanda | <i>Rhinolophus</i> sp. | |

All accession numbers are from GenBank with the exception of those beginning with EPI_ISL, which are from GISAID. Metadata includes sequencing year, geographic origin, and host species. Citations used to determine hACE2 binding capability are also included.

^aViruses that were not cultured but their spike was shown to enable (or not) hACE2-mediated entry using pseudotyped or recombinant viruses.

^bThese sequences were not included in the final phylogenetic reconstruction due to high genetic identity with another sequence in the alignment.

and SD 10 to help with convergence, as the default of 1/X is an improper prior. Our phylogenetics and time estimates are in accordance with those proposed by Boni et al. (2020). As the substitution rate in the spike gene is undoubtedly higher than in RdRp, the same clock rate prior could not be used for the Region E time-calibrated phylogeny because the divergence dates would not be comparable. Instead, we assumed the age of the root of this tree should be approximately the same as the age of the Region A tree and fixed the tree height to match the posterior estimate of the tree height for Region A (770 years before present, 1250 AD). This was done by adding a monophyletic time to most recent common ancestor (tMRCA) prior to all taxa with a Laplace distribution with mu 1250 and scale 0.1. To account for lineage-specific substitution rates, we also tested a relaxed lognormal clock model.

2.3 Screening for recombination using detection algorithms

We restricted our search for recombination breakpoints to the region of sequence beginning 750 base pairs upstream from RdRp (SARS-CoV-2 nucleotide 12,681) through the end of S2 (through SARS-CoV-2 nucleotide 25,176). There are undoubtedly other breakpoints outside of this region, but since our analysis focuses primarily on RdRp and the spike, the recombination events elsewhere in the genome are outside the scope of this study. We used the program 3SEQ (Lam, Ratmann, and Boni 2018) to test the 13 putative recombinants within Lineage 1 (SARS-CoV-1, SARS-SZ3, LYRa11, Rs3367, WIV1, RsSHC014, Rs4084, YN2018B, Rs7327, Rs9401, Rs4231, WIV16, Rs4874) and RmYN02 individually. If breakpoints were found, each subregion on either side of the breakpoint was assessed separately to fine-tune our assessments until no further breakpoints were

identified. We did not test any of the remaining sequences for recombination. We were able to identify six regions across all 13 recombinants that appear to be free of recombination and chose these for further phylogenetic analysis (above). The topologies of regions A and E are not significantly different from the topologies of RdRp and the RBD, respectively, suggesting that our use of RdRp and RBD phylogenies in Figs. 1, 2, and 5 is a sufficient representation despite some minor evidence of recombination (e.g. LYRa11).

2.4 Cell culture and transfection

BHK and 293T cells were obtained from the American Type Culture Collection and maintained in Dulbecco's modified Eagle's medium (DMEM; Sigma-Aldrich) supplemented with 10% fetal bovine serum (FBS), penicillin/streptomycin and L-glutamine. BHK cells were seeded and transfected the next day with 100 ng of plasmid encoding hACE2 or an empty vector using polyethylenimine (Polysciences). VSV plasmids were generated and transfected onto 293T cells to produce seed particles as previously described (Letko, Marzi, and Munster 2020). CoV spike pseudotypes were generated as described by Letko et al. (2018) and transfected onto 293T cells. After 24 h, cells were infected with VSV particles as described by Takada et al. (1997), and after 1 h of incubating at 37°C, cells were washed three times and incubated in 2 ml DMEM supplemented with 2% FBS, penicillin/streptomycin and L-glutamine for 48 h. Supernatants were collected and centrifuged at 500g for 5 min, then aliquoted and stored at -80°C.

2.5 Western blots

293T cells transfected with CoV spike pseudotypes (producer cells) were lysed in 1% sodium dodecyl sulfate, 150 mM NaCl,

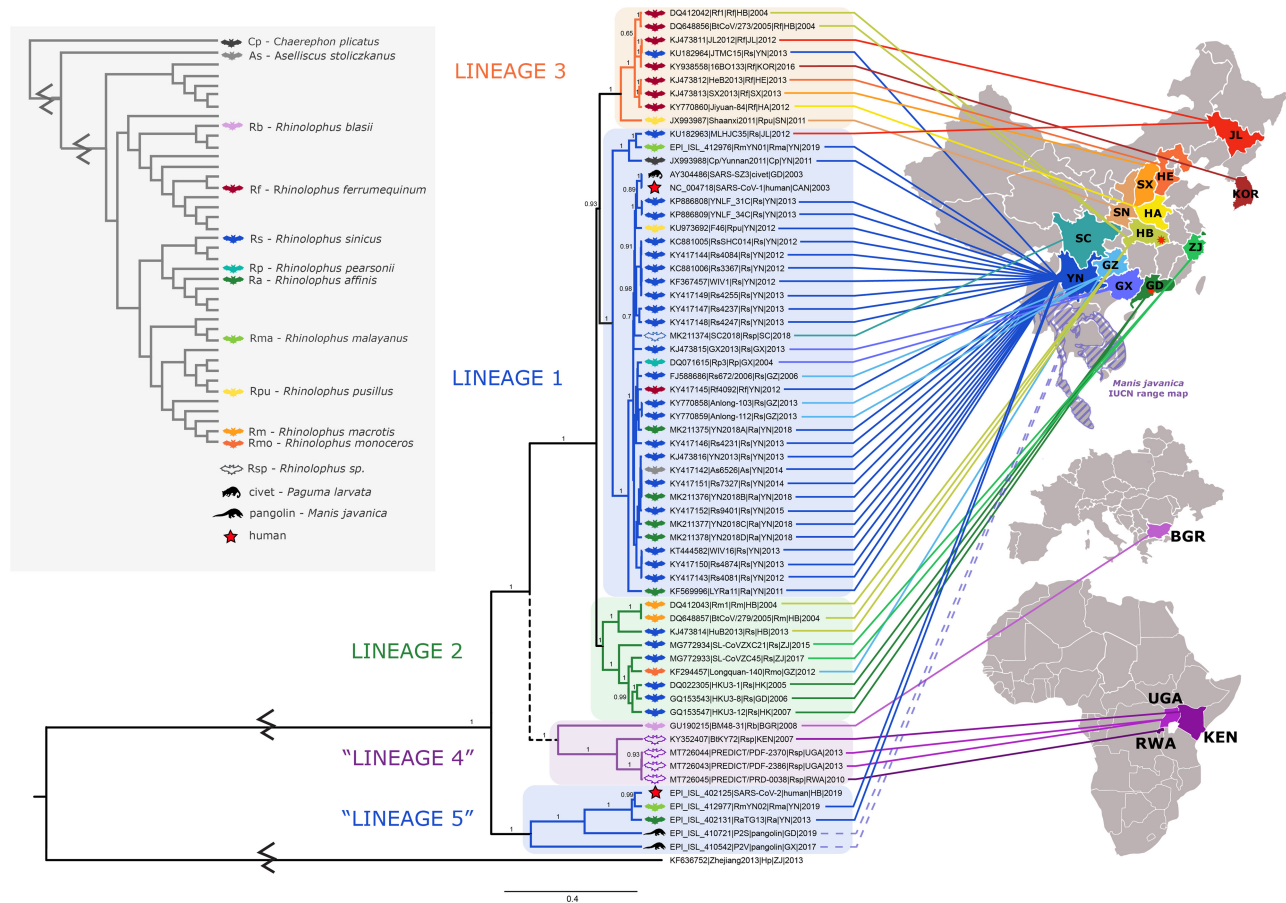


Figure 1. Phylogenetic tree of the RNA-dependent RNA polymerase (RdRp) gene (nsp12) and associated geographic origin and host species. Colors of clad bars represent the different geographic lineages. Lineage 1 is shown in blue, Lineage 2 in green, and Lineage 3 in orange. The clade of viruses from Africa and Europe is putatively named 'Lineage 4' and is shown in purple. The phylogeny shows strong posterior support for the branching order presented; however, different models or genes have produced trees with different branching orders placing Lineage 4 outside Lineage 5, so the branch to Lineage 4 is dashed to represent this uncertainty (Supplementary Fig. S1). The putative 'Lineage 5' containing SARS-CoV-2 is also shown in blue at the bottom of the tree to demonstrate that the sequences are from the same regions as Lineage 1 viruses. The geographic origin of each virus is indicated by the lines that terminate in the respective country or province with the same color code. The full province and country names for all two- and three-letter codes can be found in Table 1. As human, civet, and pangolin viruses cannot be certain to have naturally originated in the province in which they were first found, their locations are not illustrated, but the natural range of the pangolin (*Manis javanica*) is denoted with dashed shading and the origins of the SARS-CoV-1 and SARS-CoV-2 human outbreaks are designated with red stars in Guangdong and Hubei, respectively. Hosts are also shown with colored symbols according to the key on the left. The host phylogeny in the key was adapted from Agnarsson et al. (2011). The root of the tree was shortened for clarity.

50 mM Tris-HCl, and 5 mM EDTA and centrifuged at 14,000g for 20 minutes. Pseudotyped particles were concentrated from producer cell supernatants that were overlaid on a 10% OptiPrep cushion in PBS (Sigma-Aldrich) and centrifuged at 20,000g for 2 h at 4°C. Lysates and concentrated particles were analyzed for FLAG (Sigma-Aldrich; A8592; 1:10,000), GAPDH (Sigma-Aldrich; G8795; 1:10,000), and/or VSV-M (Kerafast; 23H12; 1:5,000) expression on 10% Bis-Tris PAGE gel (Thermo Fisher Scientific).

2.6 Cell entry assays

Luciferase-based cell entry assays were performed as described by Letko, Marzi, and Munster (2020). For each experiment, the relative light unit for spike pseudotypes was normalized to the plate relative light unit average for the no-spike control, and relative entry was calculated as the fold-entry over the negative control. Three replicates were performed for each CoV pseudotype.

2.7 Structural modeling

RBDs were modeled using Modweb (Pieper et al. 2011). Modeled RBDs were docked to hACE2 by structural superposition to the experimentally determined interaction complex between SARS-CoV-1 RBD and hACE2 (PDB 2ajf) (Li et al. 2005) using Chimera (Pettersen et al. 2004).

3. Results

To better understand the evolutionary history of sarbecoviruses, we first constructed a phylogenetic tree of the RNA-dependent RNA polymerase (RdRp) gene, also known as nsp12 (Fig. 1). The tree was constructed using sequences from GenBank as well as three sequences of a novel sarbecovirus detected in bats from Uganda and Rwanda as part of the USAID-PREDICT project. The three novel sequences share >99% nucleotide identity to each other and ~76% and ~74% nucleotide identity with SARS-CoV-1 and SARS-CoV-2, respectively. Phylogenetically, they lie within Lineage 4, clustering with

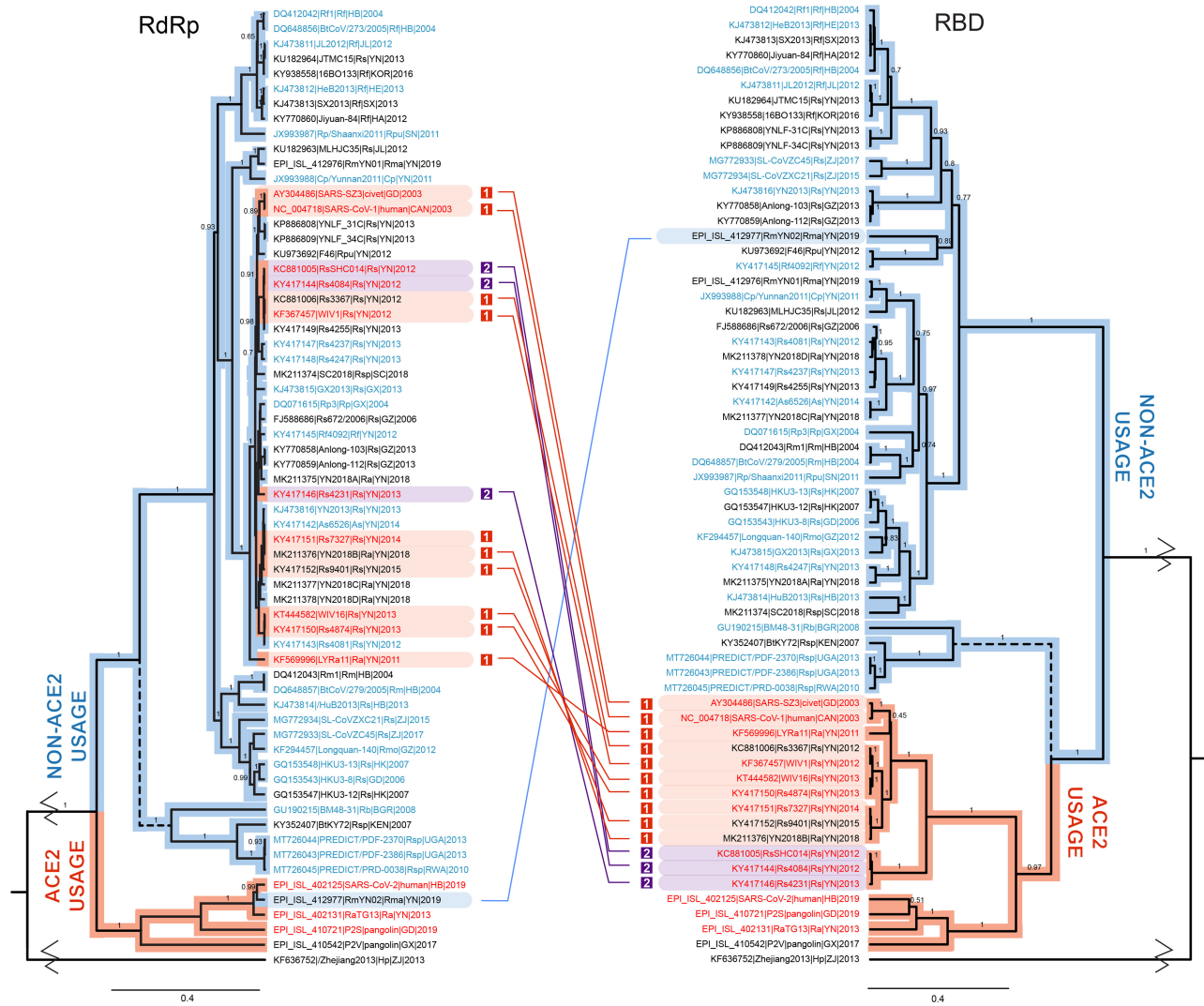


Figure 2. Phylogenetic trees of RdRp (left) and the RBD (right) demonstrating recombination events between ACE2-users and non-ACE2-users. Names of viruses that have been confirmed to use hACE2 are shown in red font, and those that have been shown to not use hACE2 are shown in blue font (citations can be found in Table 1). Viruses in black font have not yet been tested. The red and blue highlighted clade bars separate viruses with the structure associated with ACE2 usage (highly similar to viruses confirmed to use hACE2 specifically) and the structure with deletions that cannot use ACE2, respectively. Connecting lines indicate recombination events that resulted in a gain of ACE2 usage (red) or a loss of ACE2 usage (blue). The two different groups of RBD sequence within the Lineage 1 recombinants that gained ACE2 usage are distinguished in red (Type 1) and purple (Type 2) highlighting. The distances of the roots have been shortened for clarity. The branch leading to Lineage 4 is dashed to demonstrate uncertainty in its positioning.

previously reported SARS-related coronavirus BtKY72 found in bats in Kenya (Tao and Tong 2019) and bat coronavirus BM48-31 from Bulgaria (Drexler et al. 2010). The topology of the sarbecovirus phylogeny is uncertain with respect to the placement of the Lineage 4 viruses, with some models placing them between Lineage 5 and Lineages 1, 2, and 3, and others placing them at the base of the tree, depending on the methodology and alignment used (Zhou et al. 2020; Gorbalenya et al. 2020; Boni et al. 2020) (Supplementary Fig. S1). Our results place Lineage 4 in the former position with high posterior support for the RdRp gene, though the variability in this placement must be recognized. Figure 1 also demonstrates the same geographic pattern of concordance reported by Yu et al. (2019), where viruses in each lineage show a clear pattern of fidelity with particular geographic regions. However, SARS-CoV-2 does not lie within the clade of bat sarbecoviruses that have been detected in bats in China to date but rather forms a much deeper, separate lineage.

The discovery of the ‘Lineage 5’ clade containing SARS-CoV-2 and related viruses in pangolins and bats is a deviation from the geographic patterns observed for other sarbecoviruses. To investigate the evolutionary history of ACE2 usage, we built a second phylogenetic tree using only the RBD of the spike gene and compared it to the phylogeny of RdRp (Fig. 2). This region was selected because the spike protein mediates cell entry and because previous reports showed that SARS-CoV-1 and SARS-CoV-2 both use hACE2, despite being distantly related in the RdRp (Wenhui Li et al. 2003; Zhou et al. 2020). Within the RBD region of the genome, SARS-CoV-1 and all ACE2-using viruses are much more closely related to SARS-CoV-2 than to other Lineage 1 viruses (Fig. 2). Interestingly, bat virus RmYN02 is no longer associated with SARS-CoV-2 in the RBD and is instead within the clade of non-ACE2-using viruses. We also found that within the RBD, ACE2-using viruses and non-ACE2-using viruses are perfectly phylogenetically separated. The

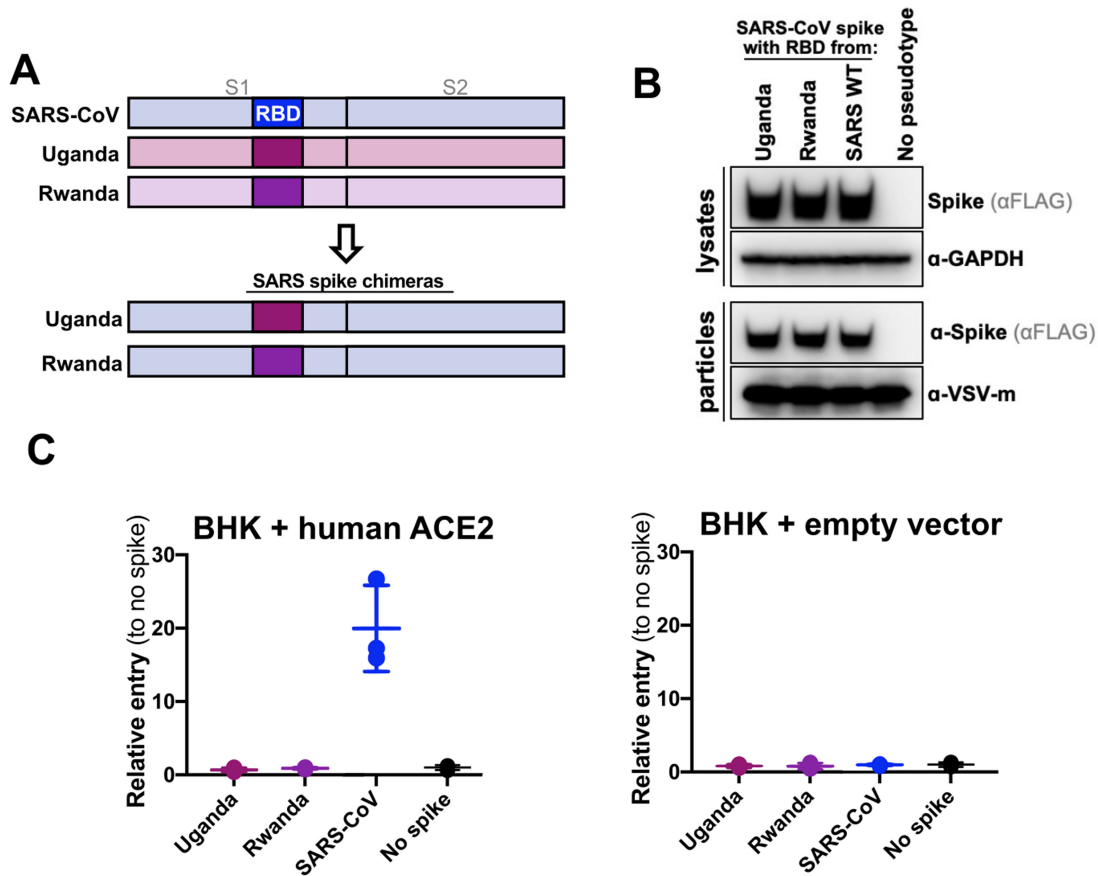


Figure 3. hACE2 usage of bat sarbecoviruses investigated using a surrogate VSV-pseudotyping system. (A) Schematic showing the structure of chimeric spike proteins. The SARS-CoV-1 spike backbone is used in conjunction with the RBD from the Uganda and Rwanda strains. (B) Incorporation of chimeric SARS-CoV-1 spike proteins into VSV. Western blots show successful expression of chimeric spikes (lysates) and their incorporation into VSV (particles). (C) hACE2 entry assays. Left, wildtype SARS-CoV spike protein is able to mediate entry into BHK cells expressing hACE2. In contrast, recombinant spike proteins containing either the Uganda or Rwanda RBD were unable to mediate entry. Entry is expressed relative to VSV particles with no spike protein. Right, control experiment for entry assay. BHK cells do not express hACE2 and therefore do not permit entry of hACE2-dependent VSV pseudotypes.

viruses from Africa and Europe form a distinct clade that is intermediate between the ACE2-using and non-ACE2-using groups but appears more closely related to the ACE2-using group.

While these viruses from Africa and Europe are slightly more similar to the ACE2-using group, they differ somewhat in amino acid sequence from the ACE2-users at the binding interface, including a small deletion in the middle of the sequence (Fig. 5, region 2). Thus, to determine the ability of these sarbecoviruses to use hACE2 and better delineate the boundaries of ACE2 usage, we performed *in vitro* experiments in which we replaced the RBD of SARS-CoV-1 with the RBD from the Uganda (PDF-2370, PDF-2386) and Rwanda viruses (PRD-0038) (Letko, Marzi, and Munster 2020). Single-cycle vesicular stomatitis virus (VSV) reporter particles containing the recombinant SARS-Uganda and SARS-Rwanda spike proteins were then used to infect BHK cells expressing hACE2. While VSV-SARS-CoV-1 showed efficient usage of hACE2, VSV-Uganda and VSV-Rwanda did not (Fig. 3).

To try and explain why the African sarbecoviruses are unable to use hACE2, we modeled the RBD domain of the sequences from Uganda (PDF-2370, PDF-2386) and Rwanda (PRD-0038). Unlike other non-ACE2 binders, homology modeling suggests that the RBDs of these viruses from Africa are structurally similar to SARS-CoV-1 and SARS-CoV-2 (Fig. 4A). However, modeling

the interaction with hACE2 reveals amino acid differences at key interfacial positions that can help explain the lack of interaction observed for the rVSV-Uganda and rVSV-Rwanda viruses (Fig. 4B and C). There are four regions of the RBD that lie within 10 Å of the interface with hACE2, one of which is the receptor binding ridge (SARS-CoV-1 residues 459–477) that is critical for hACE2 binding (Prabakaran, Xiao, and Dimitrov 2004; Shang et al. 2020). We have designated the remaining regions as regions 1 (residues 390–408), 2 (residues 426–443), and 3 (residues 478–491) (Fig. 5).

The sarbecoviruses from Africa evaluated here have a 2–3 amino acid deletion (SARS-CoV-1 residues 434–436) in region 2 (Fig. 5). As many of the residues in this region make close contact with hACE2 (<5 Å), it is possible that this contributes to the disruption of hACE2 binding. One of these residues, Y436, establishes hydrogen bonds with human ACE residues D38 and Q42 in both SARS-CoV-1 and SARS-CoV2 (Fig. 4C). Notably, all other non-ACE2 binders also have deletions in residues 432–436. While this deletion is thought to interfere or reduce binding, restoring a similar deletion (SARS-CoV-1 residues 432–437) in the S protein of an European CoV (BM48-31) with the corresponding consensus segment obtained from Lineage 1 ACE2-binding viruses did not restore hACE2-mediated entry; only replacing the receptor-binding motif increased hACE2-mediated entry (Letko, Marzi, and Munster 2020).

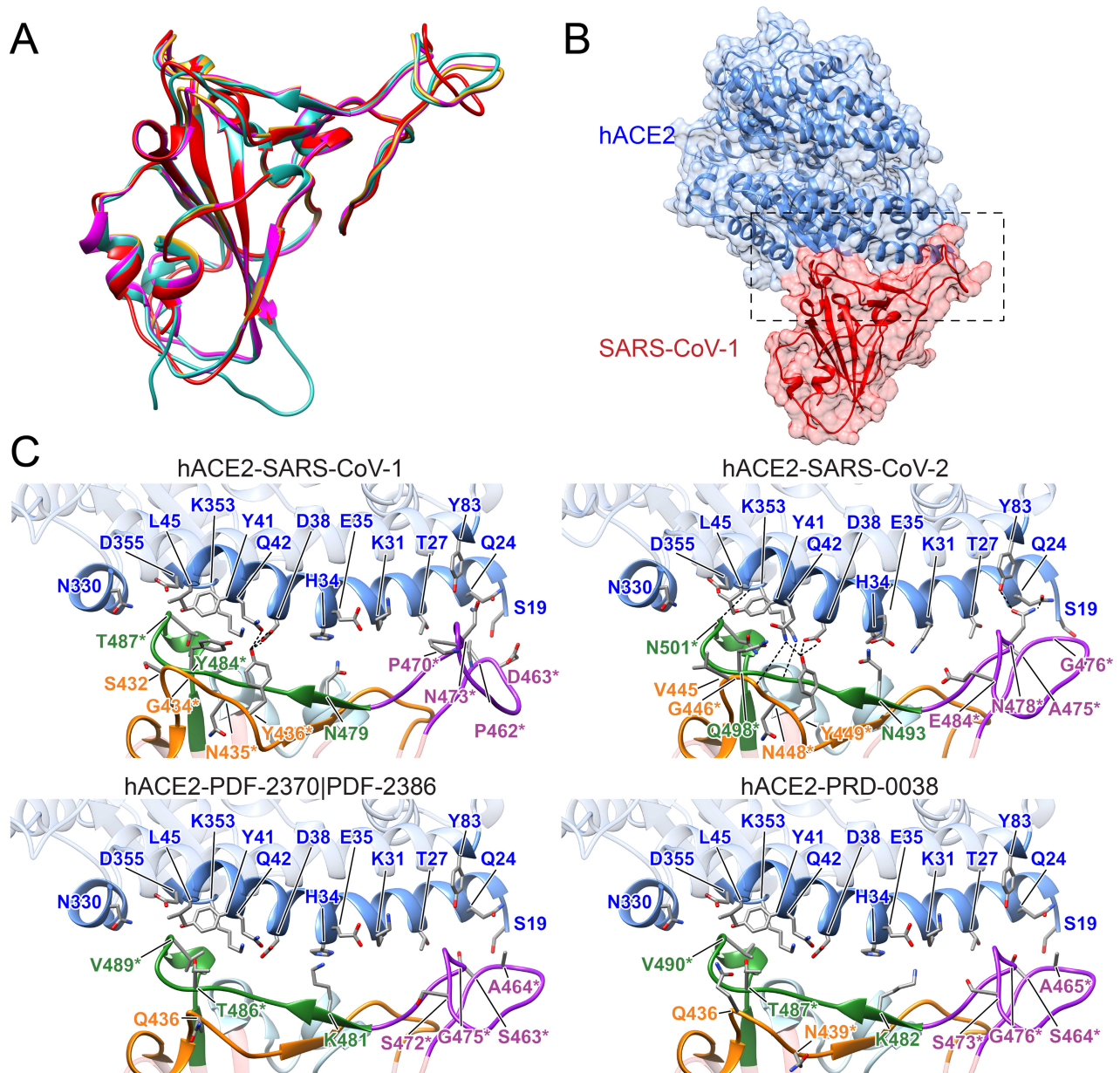


Figure 4. Structural modeling of sarbecovirus RBDs found in Uganda and Rwanda. (A) Structural superposition of the X-ray structures for the RBDs in SARS-CoV-1 (PDB 2ajf, red) (Li et al. 2005) and SARS-CoV-2 (PDB 6m0j, cyan) (Lan et al. 2020) and homology models for SARS-CoV found in Uganda (PDF2370 and PDF-2386, magenta) and Rwanda (PRD-0038, yellow). (B) Overview of the X-ray structure of SAR-CoV-1 RBD (red) bound to hACE2 (blue) (PDB 2ajf, red) (F. Li et al. 2005). (C) Close-up view of the interface between hACE2 (blue) and RBDs in SARS-CoV-1 (PDB 2ajf, top left) (Li et al. 2005) and SARS-CoV-2 (PDB 6m0j, top right) (Lan et al. 2020) and homology models for viruses found in Uganda (PDF-2370 and PDF-2386, bottom, left) and Rwanda (PRD-0038, bottom, right). The color of the RBD loops corresponds to the colors of the labeled sequence regions in Fig. 5: region 1 in cyan, region 2 in orange, the receptor binding ridge in purple, and region 3 in green. Labeled RBD residues correspond to interfacial residues whose identity differ in African sarbecoviruses and SARS-CoV-1 or SARS-CoV-2 (labels are included in all four panels to facilitate the identification of counterpart residues in each virus). Asterisks denote residues whose identity is not shared by any ACE-2 binding SARS-CoV as dictated by Fig. 5. Labeled hACE2 residues correspond to residues within 5 Å of RBD residues depicted.

Moreover, sarbecoviruses from Africa contain additional amino acid changes at the interface that can also contribute to hACE2 binding disruption (Fig. 4C). hACE2 contains two hotspots (K31 and K353) that are crucial targets for binding by SARS-RBDs and amino acid variations in the RBD sequence enclosing these ACE2 hotspots have been shown to shape viral infectivity, pathogenesis, and determine the host range of SARS-CoV-1 (Li et al. 2005; Li 2008; Wu et al. 2012). All sarbecoviruses from Africa contain a Lys (K) at SARS-CoV-1 position 479 within region 3 (positions 481 and 482 of Uganda and Rwanda,

respectively), which makes contact with these ACE2 hotspots (as compared to N479 or Q493 in SARS-CoV-1 and 2, respectively; Fig. 4C). K479 decreases binding affinity by more than 20-fold in SARS-CoV-1 (Li et al. 2005). The negative contribution of K479 in region 3 is likely due to unfavorable electrostatic contributions with ACE2 hotspot K31 (Fig. 4C) (Li 2008; Wan et al. 2020). On the other hand, SARS-CoV-1 residue T487 (N501 in SARS-CoV-2) interacts with ACE2 hotspot K353 and has a Val (V) in the viruses from Africa (residues 489 and 490) (Fig. 5). As with residue 479, the amino acid identity at position 487 contributes

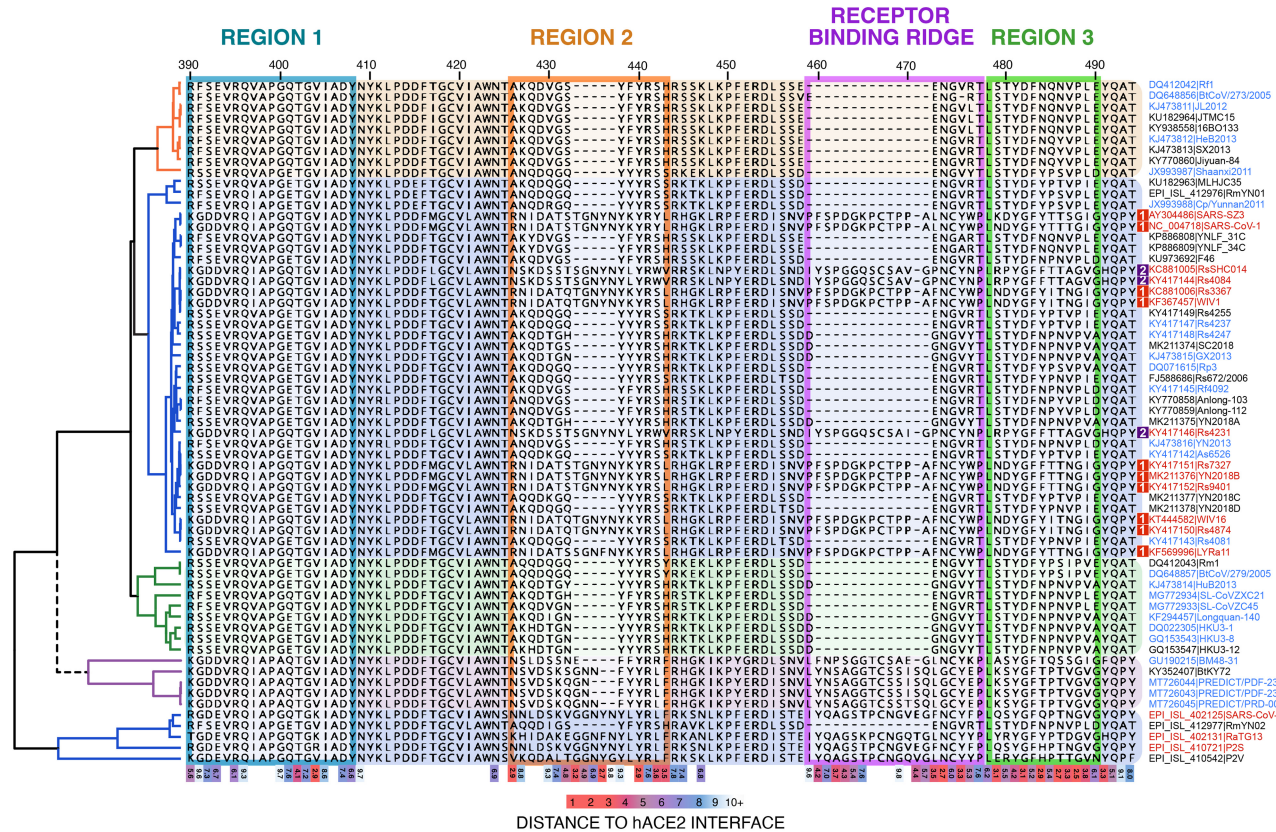


Figure 5. The phylogenetic backbone of the RdRp gene alongside the amino acid sequences of the RBM. Amino acid numbering is relative to SARS-CoV-1. Virus names in red font are known hACE2 users, those in blue are known non-users, and those in black have not been tested. Residues within 10 Å of the interface with hACE2 are considered interfacial, and exact distances between each interfacial residue and the closest hACE2 residue (based on structural modeling of SARS-CoV-1 bound with hACE2) are shown along the bottom. Residues that are closer to the interface (3 Å or less) and thus make strong interactions with hACE2 are shown in red, and as distance increases this color transitions to purple, blue, and finally to white. The receptor binding ridge sequences are highlighted in purple and the remaining interfacial segments have been numbered regions 1, 2, and 3 for clarity within the main text. The colors of these regions correspond with the colors in the structural models in Fig. 4. The branch leading to Lineage 4 is dashed to demonstrate uncertainty in its positioning.

to the enhanced hACE2 binding observed in SARS-CoV-2 (Li 2008; Wu et al. 2012; Wan et al. 2020). The presence of a hydrophobic residue at position 487, not previously observed in any ACE2 binding sarbecovirus, might lead to a local rearrangement at the K353 hotspot that hinders hACE2 binding. Indeed, most non-ACE2 binders have a Val (V) in SARS-CoV-1 position 487 (Fig. 5).

Finally, the receptor binding ridge, which is conspicuously absent from all non-ACE2 binders, is present in the sarbecoviruses from Africa but has amino acid variations that differ significantly from both SARS-CoV-1 and SARS-CoV-2 (Fig. 5). Changes in the structure of this ridge contribute to increased binding affinity of SARS-CoV-2, as a Pro-Pro-Ala (PPA) motif in SARS-CoV-1 (residues 469–471) replaced with Gly-Val-Glu-Gly (GVEG) in SARS-CoV-2 results in a more compact loop and better binding with hACE2 (Shang et al. 2020). Changes within this ridge may be negatively contributing to hACE2 binding of viruses from Africa, which have Ser-Thr-Ser-Gln (STSQ) or Ser-Iso-Ser-Gln (SISQ) in this position (Figs. 4C and 5).

While our studies suggest that these viruses from Africa do not utilize hACE2, it is not clear whether they are still ACE2-users but are adapted to divergent forms of bat ACE2 in their natural hosts. The specific bat host species for the Uganda and Rwanda viruses reported here could not be definitively identified in the field or in the lab but are all genetically identical. They may represent a cryptic species, as the mitochondrial

sequences are ~94% identical with *Rhinolophus ferrumequinum* in the cytochrome oxidase I gene (COI) and ~96% identical with *Rhinolophus clivosus* in the cytochrome b (cytb) gene, each of which have been deposited in GenBank (accessions MT738926–MT738928, MT732776). We were also able to extract ACE2 sequences from the deep sequencing reads of PDF-2370 (GenBank accession MW183243) to compare it to ACE2 sequences from species that are known to host ACE2 binders (human, civet, pangolin), non-ACE2 binders (*R. macrotis*, *pearsonii*, *pusillus*, *ferrumequinum*), and both (*R. sinicus*). Comparison of the ACE2 sequences shows that they are highly similar, with only a few amino acids that are changed in hosts of viruses that utilize ACE2 compared to the host of our African bat sample (Supplementary File S1). *R. sinicus* in particular is a known host of viruses that utilize ACE2 as well as the deletions that do not, suggesting that adaptation to divergent bat ACE2 is not a likely explanation for the deviation in sequence and structure of the RBD of viruses with deletions, including the novel sarbecoviruses from Uganda and Rwanda. These findings provide additional structural evidence that aids in distinguishing viruses which bind ACE2 from those that do not. They also demonstrate that ACE2 usage within sarbecoviruses is restricted to those viruses within the SARS-CoV-1 and SARS-CoV-2 clade in the RBD (Lineages 1 and 5, Fig. 2).

The finding of discordant evolutionary trees for RdRp and the RBD in Fig. 2 more strongly supports a recombination

scenario; however, to consider an alternate scenario where ACE2 usage arose in Lineages 1 and 5 independently through convergent evolution, we compared the RdRp phylogeny with the amino acid sequences of the interfacial residues in the RBD (Fig. 5). When mapped to the RdRp tree, the 'extra' RBD sequence present in the ACE2-using viruses is conspicuous within the Lineage 1 clade of otherwise non-ACE2-using viruses that have large deletions. We also note that there are two distinct groups of RBD sequences within ACE2-using Lineage 1 viruses: Type 1, containing SARS-CoV-1, SARS-SZ3 (civet), Rs3367, WIV1, Rs7327, YN2018B, Rs9401, WIV16, Rs4874, and LYRa11, and Type 2, containing Rs4231, Rs4084, and RsSHC014. Further, RmYN02 is within the Lineage 5 clade of ACE2-using viruses in RdRp but its RBD sequence contains both deletions (Fig. 5). Without recombination, the viruses with deletions in region 2 and in the receptor binding ridge would have had to be gained and lost in precisely the same positions for ACE2-using Lineage 1 viruses and RmYN02, respectively, which is not a parsimonious explanation. The phylogeny and sequence in Fig. 5 also illustrate that ACE2 usage appears to be an ancestral trait conserved in Lineage 5 (Boni et al. 2020) and a derived trait in each of the 13 Lineage 1 viruses with ACE2-using structure.

Finally, we further investigated support for the recombination scenario by examining the region of sequence between RdRp and the RBD for possible breakpoints. Only the 13 Lineage 1 viruses with ACE2-using structure were targets of this analysis as we were primarily interested in explaining the discordant phylogeny and variation in ACE2 usage (Fig. 2), not in fully describing the recombination history of every sarbecovirus. Using 3SEQ, we show that all of the ACE2-using Lineage 1 sequences show extensive evidence of recombination within S1 and the RBD specifically (Table 2, Fig. 6A). Further, the assignment of the parental sequence that donated the recombinant region (the minor parent) always resulted in the identification of one of the other recombinant sequences. This would not have been possible, as the recombinant region would have had to come from somewhere other than these 13 sequences, indicating that the true minor parent does not exist in our alignment. Using these breakpoints, we designated six subregions that were relatively free of recombination within these 13 sequences, mirroring the approach of Boni et al. 2020 (2020) and built phylogenetic trees for each region. We show that in orf1ab (region A) and S2 (region F), these 13 sequences fall within Lineage 1, but within S1 and particularly the RBD (B through E) they switch phylogenetic positions and cluster with Lineage 5 (Fig. 6B), supporting the recombination scenario.

Despite only investigating the Lineage 1 recombinants for the locations of sequence breakpoints, the phylogenetic trees provide evidence that recombination has occurred frequently in other sarbecoviruses in this genomic region as well (Fig. 6B). Of note, Rs4084 and RsSHC014 cluster with Type 1 RBDs in regions B, C, and D, but with swap to cluster with Rs4231 (Type 2) in Region E, even though Rs4084, RsSHC014, WIV1, and Rs3367 are all nearly identical in every other region. This suggests that a WIV1/Rs3367-like Type 1 virus which had already undergone recombination in regions B through E underwent a second recombination event with a Type 2 virus on top of the first in region E. A number of other viruses also appear to have recombinant history in regions B, C, and D (SL-CoVZC45 and SL-CoVZXC21, YN2013, Anlong-103, and Anlong 112), but these viruses do not show evidence of recombination that spans the RBD in region E, which contains the amino acid deletions in region 2 and the receptor binding ridge and appears to primarily determine ACE2-using potential. The frequency of recombination in this region

among Lineage 1 viruses strongly supports the hypothesis that after ACE2 usage was acquired in Lineage 1, it subsequently spread throughout the clade via additional recombination events with other Lineage 1 viruses.

As all of our evidence supports a recombination scenario over convergent evolution, we sought to construct a possible timeline of events that could explain our observations. Using tip dating in BEAST2, we constructed a time-calibrated phylogeny for RdRp using a substitution rate prior inferred from Boni et al. 2020 (2020). Using the RdRp tree as an evolutionary backbone, the deletions in region 2 and the receptor binding ridge of the RBD appear to have been lost in a stepwise fashion (Fig. 5). The small deletion in region 2 likely arose first, before the diversification of Lineage 4 in Africa and Europe (Fig. 5) and was dated using the tMRCA of Lineages 1, 2, 3 and 4 (Fig. 8). Alternatively, as the boundaries of the deletion in region 2 in Lineage 4 and Lineages 1, 2, and 3 do not align perfectly and there is uncertainty in the position of this branch in the phylogeny, it is equally possible that this deletion was lost independently in Lineage 4. The larger deletion in the receptor binding ridge, not present in known sequences from Lineage 4, likely arose second, but before the diversification of Lineages 1, 2, and 3 (Fig. 5) and was dated with the tMRCA of these three lineages (Fig. 8). Because no ACE2-using viruses have been discovered in Lineage 2 or 3 to date, we propose that the re-appearance of this trait arose after the MRCA of Lineage 1 on the tree (Fig. 8). As SARS-CoV-1 was the earliest Lineage 1 virus sequenced with ACE2-using structure, the emergence of ACE2 usage in Lineage 1 must have occurred in the time between the tMRCA of Lineage 1 (1852, 95% highest posterior density 1804–1901) and the emergence of SARS-CoV-1 in 2003.

Next, we constructed a time-calibrated phylogeny for RBD with a strict tMRCA age prior informed by the estimation of the tree height in RdRp (see Section 2), such that the timescale would be comparable even though the evolutionary rates between these two regions likely are not the same (Fig. 7). To account for variability in lineage-specific substitution rates, we also generated a time-calibrated model using a relaxed lognormal clock (Fig. 7). Comparing the time-calibrated RBD tree to the time-calibrated RdRp tree, the divergence dates for the two types of RBD sequence observed in the recombinant Lineage 1 sequences are incompatible, suggesting that more than one recombination event donating ACE2 usage from Lineage 5 to Lineage 1 must have occurred. The 13 Lineage 1 recombinants (both Type 1 and Type 2) coalesce between 119 and 216 years ago in RdRp and between 259 and 490 years ago in the RBD (Fig. 7). If these time estimates reflect true rates of diversification, a single introduction of the ACE2-using phenotype via recombination would not allow enough time for the sequence divergence between Type 1 and Type 2 RBDs to accumulate, even when accounting for the substitution rate in RBD being estimated as an order of magnitude higher than that of RdRp ($5.248e-4$ in RdRp, $2.181e-3$ in RBD). Further, the substitution rate that would be needed for the observed sequence divergence in the RBD of the 13 recombinants to have accumulated since their MRCA in RdRp (1852) is more than double the estimated rate of our time-calibrated tree ($5.899e-3$). Even with a relaxed clock assumption, the maximum value of the posterior distribution of the mean rate is only $4.733e-3$. From this, we conclude that two independent recombination events occurred between Lineage 5 and Lineage 1 resulting in two distinct RBD types.

We propose two main hypotheses for the acquisition and spread of the two distinct RBD types donating ACE2 usage from Lineage 5 to Lineage 1. The recombination hypothesis posits

Table 2. Recombination breakpoints detected in ACE2-using Lineage 1 viruses by the program 3SEQ.

| Major parent | Minor parent | Child | P | Length | Breakpoint estimates |
|----------------|----------------|------------|------------|--------|--|
| KU973692 | EPI_ISL_402131 | NC_004718 | 0 | 952 | 8836–8837 and 10510–10542 |
| F46 | RaTG13 | SARS-CoV-1 | | | 8836–8837 and 10726–10752 |
| MK211374 | EPI_ISL_412976 | NC_004718 | 0 | 1290 | 6497–6519 and 8363–8365 |
| SC2018 | RmYN01 | SARS-CoV-1 | | | 6401–6406 and 8363–8365 6440–6472 and 8363–8365 |
| KY417146 | KY417151 | NC_004718 | 0 | 573 | 9760–9772 and 10702–10704 |
| Rs4231 | Rs7327 | SARS-CoV-1 | | | |
| MG772933 | KY770860 | NC_004718 | 1.4775E-07 | 1072 | 11035–11037 and 12610–12624 |
| SL-CoVZC45 | Jiyuan-84 | SARS-CoV-1 | | | |
| KY770859 | KY352407 | AY304486 | 0 | 993 | 8620–8681 and 10732–10771 |
| Anlong-112 | BtKY72 | SARS-SZ3 | | | |
| MK211374 | KJ473814 | AY304486 | 1.1774E-07 | 1077 | 6755–6784 and 8397–8431 |
| SC2018 | HuB2013 | SARS-SZ3 | | | |
| KY417146 | MK211376 | AY304486 | 0 | 558 | 9760–9772 and 10702–10704 |
| Rs4231 | YN2018B | SARS-SZ3 | | | |
| MG772933 | KP886808 | AY304486 | 1.592E-07 | 791 | 11260–11273 and 12543–12558 |
| SL-CoVZC45 | YNLF_31C | SARS-SZ3 | | | |
| EPI_ISL_412976 | NC_004718 | KF569996 | 0 | 921 | 9107–9113 and 10700–10701 |
| RmYN01 | SARS-CoV-1 | LYRa11 | | | 9027–9043 and 10865–10869 9077–9095 and 10865–10869 9107–9113 and 10865–10869 9027–9043 and 10840–10842 9077–9095 and 10840–10842 9107–9113 and 10840–10842 9027–9043 and 10700–10701 9077–9095 and 10700–10701 |
| JX993988 | KY770859 | KF569996 | 0 | 1627 | 1658–1714 and 4151–4199 |
| Cp/Yunnan2011 | Anlong-112 | LYRa11 | | | 1368–1428 and 4229–4240 1487–1498 and 4229–4240 1658–1714 and 4229–4240 1368–1428 and 4151–4199 1487–1498 and 4151–4199 |
| NC_004718 | KY417142 | KC881006 | 0 | 2117 | 0–11 and 9245–9251 |
| SARS-CoV-1 | As6526 | Rs3367 | | | |
| KC881005 | KF569996 | KC881006 | 0 | 168 | 10201–10233 and 10549–10565 |
| RsSHC014 | LYRa11 | Rs3367 | | | |
| KY417151 | KY417142 | KC881006 | 0 | 3036 | 1853–3932 and 8288–8374 |
| Rs7327 | As6526 | Rs3367 | | | |
| NC_004718 | KY417142 | KF367457 | 0 | 2116 | 0–11 and 9245–9251 |
| SARS-CoV-1 | As6526 | WIV1 | | | |
| KC881005 | KF569996 | KF367457 | 0 | 168 | 10201–10233 and 10549–10565 |
| RsSHC014 | LYRa11 | WIV1 | | | |
| KY417151 | KY417142 | KF367457 | 0 | 3036 | 1853–3932 and 8288–8374 |
| Rs7327 | As6526 | WIV1 | | | |
| KF367457 | KY417146 | KC881005 | 0 | 378 | 9841–9915 and 10549–10572 |
| WIV1 | Rs4231 | RsSHC014 | | | |
| KY417151 | KY417142 | KC881005 | 0 | 3037 | 1853–3932 and 8288–8374 |
| Rs7327 | As6526 | RsSHC014 | | | |
| KF367457 | KY417146 | KY417144 | 0 | 378 | 9841–9915 and 10549–10572 |
| WIV1 | Rs4231 | Rs4084 | | | |
| KY417151 | KY417142 | KY417144 | 0 | 3034 | 1853–3932 and 8288–8374 |
| Rs7327 | As6526 | Rs4084 | | | |
| NC_004718 | MK211377 | MK211376 | 0 | 2417 | 411–551 and 9245–9251 |
| SARS-CoV-1 | YN2018C | YN2018B | | | |
| KC881005 | KF569996 | MK211376 | 0 | 122 | 10201–10233 and 10469–10497 |
| RsSHC014 | LYRa11 | YN2018B | | | |
| KY417151 | MK211378 | MK211376 | 0 | 2205 | 4541–5578 and 8766–8789 |
| Rs7327 | YN2018D | YN2018B | | | |
| NC_004718 | KY417142 | KY417151 | 0 | 2112 | 0–11 and 9245–9251 |
| SARS-CoV-1 | As6526 | Rs7327 | | | |
| KC881005 | KF569996 | KY417151 | 0 | 122 | 10201–10233 and 10469–10497 |
| RsSHC014 | LYRa11 | Rs7327 | | | |

(continued)

Table 2.. (continued)

| Major parent | Minor parent | Child | P | Length | Breakpoint estimates |
|------------------------------|-------------------------|--------------------------|-------|--------|--|
| KY417144 Rs4084 | MK211377 YN2018C | KY417151 Rs7327 | 0 | 3260 | 924–1939 and 8186–8374 |
| NC_004718 SARS-CoV-1 | KY417142 As6526 | KY417152 Rs9401 | 0 | 2112 | 0–11 and 9245–9251 |
| KC881005 RsSHC014 | KF569996 LYRa11 | KY417152 Rs9401 | 0 | 122 | 10201–10233 and 10469–10497 |
| KY417144 Rs4084 | MK211377 YN2018C | KY417152 Rs9401 | 0 | 3260 | 924–1939 and 8186–8374 |
| NC_004718 SARS-CoV-1 | KY417149 Rs4255 | KY417146 Rs4231 | 0 | 2296 | 0–11 and 8838–8840 |
| NC_004718 SARS-CoV-1 | KC881005 RsSHC014 | KY417146 Rs4231 | 0 | 1788 | 9769–9780 and 12448–12793 |
| NC_004718 SARS-CoV-1 | KY417143 Rs4081 | KT444582 WIV16 | 0 | 2293 | 0–32 and 8838–8840 |
| KF367457 WIV1 | KY417146 Rs4231 | KT444582 WIV16 | 0 | 541 | 0–8891 and 9973–10233 |
| KC881005 RsSHC014 | NC_004718 SARS-CoV-1 | KT444582 WIV16 | 0 | 403 | 0–8891 and 9769–9780 |
| KY417143 Rs4081 | KY417146 Rs4231 | KT444582 WIV16 | 4E-12 | 1781 | 5975–6133 and 8727–12793 3536–5782 and 8727–12793 |
| NC_004718 SARS-CoV-1 | KY417143 Rs4081 | KY417150 Rs4874 | 0 | 2294 | 0–32 and 8838–8840 |
| KF367457 WIV1 | KY417146 Rs4231 | KY417150 Rs4874 | 0 | 541 | 0–8891 and 9973–10233 |
| KC881005 RsSHC014 | NC_004718 SARS-CoV-1 | KY417150 Rs4874 | 0 | 403 | 0–8891 and 9769–9780 |
| KY417143 Rs4081 | KY417146 Rs4231 | KY417150 Rs4874 | 4E-12 | 1782 | 5975–6133 and 8727–12793 3536–5782 and 8727–12793 |
| EPI_ISL_402125 SARS-CoV-2 | KU182964 JTMC15 | EPI_ISL_412977 RmYN02 | 0 | 1111 | 8957–8957 and 10827–10828 8938–8941 and 10831–10845 8957–8957 and 10831–10845 8938–8941 and 10827–10828 |
| EPI_ISL_410542 P2V | KY770859 Anlong-112 | EPI_ISL_412977 RmYN02 | 0 | 3218 | 1904–1907 and 5126–5128 1862–1879 and 5126–5128 1883–1885 and 5126–5128 |

Each recombinant Lineage 1 virus was set as the child sequence, and the parental sequences between the breakpoints identified (minor parent) and on either side (major parent) are listed. The *p*-value indicates the level of significance indicated by 3SEQ. Breakpoint estimates are given as ranges, and the minimum length of the recombinant region between these breakpoints is given. Numbering is relative to the alignment, which begins at SARS-CoV-2 nucleotide 12,681. When 3SEQ identified more than one set of breakpoint estimates, all were included in the table. Each recombinant region was further analyzed separately for more breakpoints within, since 3SEQ identifies only one at a time.

that two recombination events donated Type 1 and Type 2 RBD sequence from Lineage 5 to Lineage 1; however, these two events are insufficient to explain the non-monophyletic pattern of ACE2 usage in Lineage 1. We further hypothesize that whichever Lineage 1 virus first gained Type 1 and Type 2 ACE2 usage in each group then donated the trait to other Lineage 1 viruses through subsequent recombination events (Fig. 8). It is difficult to approximate a date for such an event, but the tMRCA of the Type 1 recombinants in the RBD may be a close estimation (between 42 and 77 years ago) (Fig. 7). The events must have been recent enough that the observed diversity of Type 2 RBD sequences is quite low, yet not so recent such that there would not have been time for recombination to have occurred twice in region E for sequences Rs4084 and RsSHC014 (Fig. 6B).

The second hypothesis and only remaining possibility for ACE2 usage in Lineage 1 (besides convergence) is that perhaps the trait persisted in this Lineage from the ancestral state (Fig. 8). Because no viruses demonstrating ACE2 usage have been discovered in Lineages 2, 3, and 4, this would mean that

the ACE2 usage trait would have been lost via deletion in these lineages. Further, because of the non-monophyletic branching order of these lineages, this would require multiple independent and identical losses of the region 2 and receptor binding ridge deletions in all three of these lineages. If this did indeed occur, in order to then observe the pattern of ACE2 usage in Lineage 1 where some viruses, but not all, have the ACE2 usage trait, further independent losses would be required in individual viruses. In much the same manner as convergence would require multiple independent and identical events, persistence of ACE2 usage with multiple independent deletions for the entire clades of Lineages 2, 3, and 4 and only some of the viruses in Lineage 1 is also highly non-parsimonious. Persistence is also a poor explanation for the pattern of the two RBD types observed, particularly for Type 2, where the RBD sequences are highly similar but the RdRp sequences are quite divergent. If both genes were vertically inherited via persistence, we would expect these genes to have approximately equal MRCA ages. Instead, we observe that the MRCA age for Type 2 RBDs in region E are much younger than for RdRp.

| | RdRp | | S1 | | | | | RBD | | S2 | |
|----------|---------------|------------|---------------|--------|--------------|---------------|----------|----------|--|----|--|
| A | B | C | D | E | F | | | | | | |
| SARS | SC2018 | | 6401 | 8365 | 8836 | 9760 | 10752 | | | | |
| | | | RmYN01 | SC | Rs4231 | Rs7327 | Rs7327 | | | | |
| civet | SC2018 | | 6755 | 8365 | 8836 | 9760 | 10752 | | | | |
| | | | HuIP2013 | SC | Rs4231 | Rs7327 | Rs7327 | | | | |
| RsSHC014 | Rs7327 | As6526 | | | Rs7327 | Rs4231 | WIV1 | | | | |
| | 1853 | | | 8374 | | 9841 | 10752 | | | | |
| Rs4084 | Rs7327 | As6526 | | | Rs7327 | Rs4231 | WIV1 | | | | |
| | 1853 | | | 8374 | 9251 | 10201 | 10565 | | | | |
| Rs3367 | Rs7327 | As6526 | | | Rs7327 | Rs4231 | LY | RsSHC014 | | | |
| | 1853 | | | 8374 | 9251 | 10201 | 10565 | | | | |
| WIV1 | Rs7327 | As6526 | | | Rs7327 | Rs4231 | LY | RsSHC014 | | | |
| | | | | 8840 | 9769 | | | | | | |
| Rs4231 | Rs4255 | | SARS | | RsSHC014 | | | | | | |
| | 924 | | 8374 | 9251 | 10201 | 10565 | | | | | |
| Rs7327 | Rs4084 | YN2018C | | Rs4084 | RsSHC014 | LY | RsSHC014 | | | | |
| | 924 | | 8374 | 9251 | 10201 | 10565 | | | | | |
| Rs9401 | Rs4084 | YN2018C | | Rs4084 | RsSHC014 | LY | RsSHC014 | | | | |
| | | 4541 | | 8789 | 9251 | 10201 | 10565 | | | | |
| YN2018B | Rs7327 | | YN2018D | | 7327RsSHC014 | | LY | RsSHC014 | | | |
| | | 3536 | | 8840 | 9780 | 10233 | | | | | |
| WIV16 | Rs4081 | | Rs4231 | SARS | 014 | WIV1 | | | | | |
| | | 3536 | | 8840 | 9780 | 10233 | | | | | |
| Rs4874 | Rs4081 | | Rs4231 | SARS | 014 | WIV1 | | | | | |
| | 1368 | 4240 | | 9027 | | 10869 | | | | | |
| LYRa11 | Cp/Yunnan2011 | Anlong-112 | Cp/Yunnan2011 | SARS | | Cp/Yunnan2011 | | | | | |

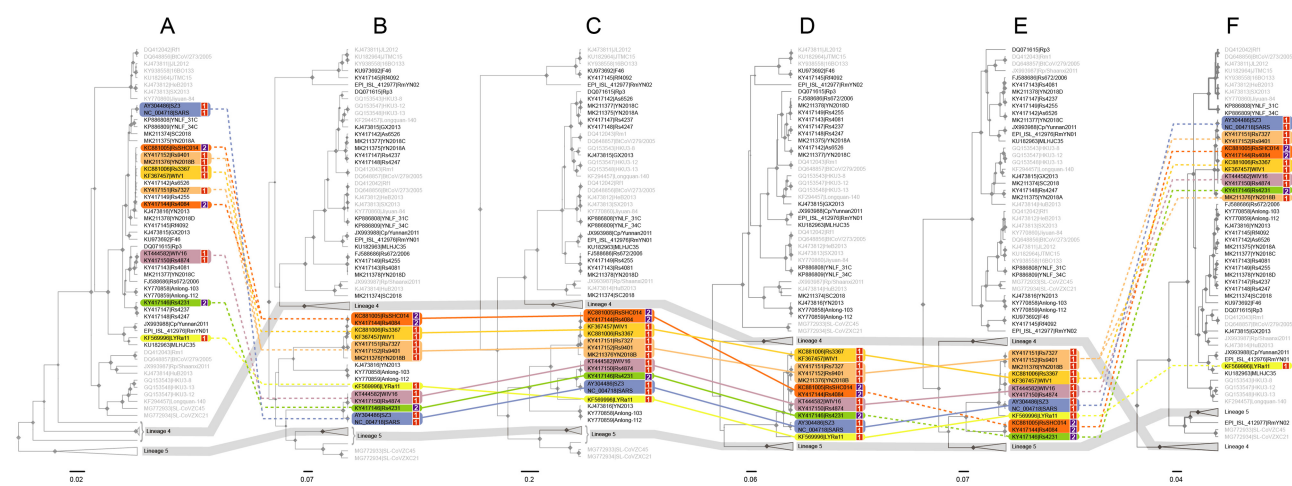


Figure 6. Recombination breakpoints detected in Lineage 1 ACE2-using sequences. The top of this figure illustrates that the recombination suggested by the change in topology in [Fig. 2](#) for 13 Lineage 1 viruses is supported by formal breakpoint analysis. The breakpoints detected for each of the 13 recombinant Lineage 1 sequences with ACE2-using structure (no deletions) are shown. Sequences that are nearly identical are colored the same for simplicity. The bars represent the sequence of genome beginning 750 bp before RdRp spanning through the end of S2 (SARS-CoV-2 nucleotides 12,681 through 25,176) and each box within represents a recombinant section within the sequence. The breakpoints correspond to those identified in [Table 2](#). Numbering is relative to the alignment. The parental sequence is shown within each box. Sequences identified as the minor parent by 3SEQ were labeled within the breakpoint margins and the major parent outside. Six regions where these sequences appear to be free of recombination are labeled A-F and a corresponding phylogeny for each region is shown below. Regions A and E were further tested for recombination breakpoints in all sequences, not just the 13 Lineage 1 viruses, and were found to be breakpoint-free. The topology of regions A and E is not different enough from [Fig. 2](#) to suggest that recombination within RdRp or RBD significantly changed the interpretation of our results. For each region, sequences were tracked with connecting lines of corresponding color to identify where recombination may have occurred between Lineage 1 and Lineage 5 and hypothesized events are specifically marked with dotted lines. This highlights the secondary recombination of Rs4084 and RsSHC014 in region E on top of the primary recombination in regions B through E. Sequence names of Lineage 2 and 3 viruses are greyed out and Lineages 4 and 5 are collapsed and highlighted in darker grey to make the changes in topology between the trees more visible.

4. Discussion

4.1 ACE2 usage in lineage 1 viruses was acquired via recombination

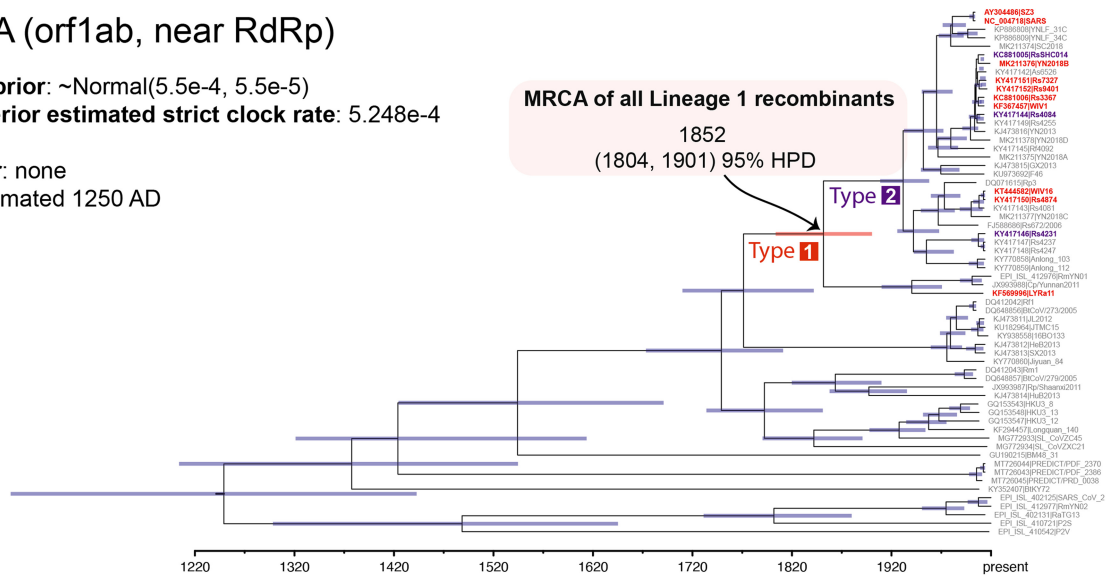
At first glance, ACE2 usage does not appear to be phylogenetically conserved among sarbecoviruses, especially since many phylogenies are built using RdRp. This naturally leads to the hypothesis that ACE2 usage arose independently in SARS-CoV-1 and SARS-CoV-2 via convergent evolution. This has been suggested previously for another ACE2-using human coronavirus, NL63 ([Chen et al. 2020](#)). However, a phylogeny constructed using

the RBD perfectly separates viruses that have been shown to utilize ACE2 from those that do not ([Fig. 2](#)). Viruses that cannot utilize ACE2 have significant differences in their RBDs, including large deletions in critical interfacial residues and low amino acid identity with viruses that do use ACE2 ([Fig. 5](#)). Notably, in addition to the large deletions, viruses that cannot use ACE2 deviate considerably at the interacting surface, including positions that play fundamental roles dictating binding and cross-species transmission ([Li et al. 2005](#); [Li et al. 2005](#); [Wan et al. 2020](#); [Shang et al. 2020](#)). It is unknown whether viruses that cannot use hACE2 are utilizing bat ACE2 or an entirely different receptor

Region A (orf1ab, near RdRp)

Clock rate prior: ~Normal(5.5e-4, 5.5e-5)
 Mean posterior estimated strict clock rate: 5.248e-4

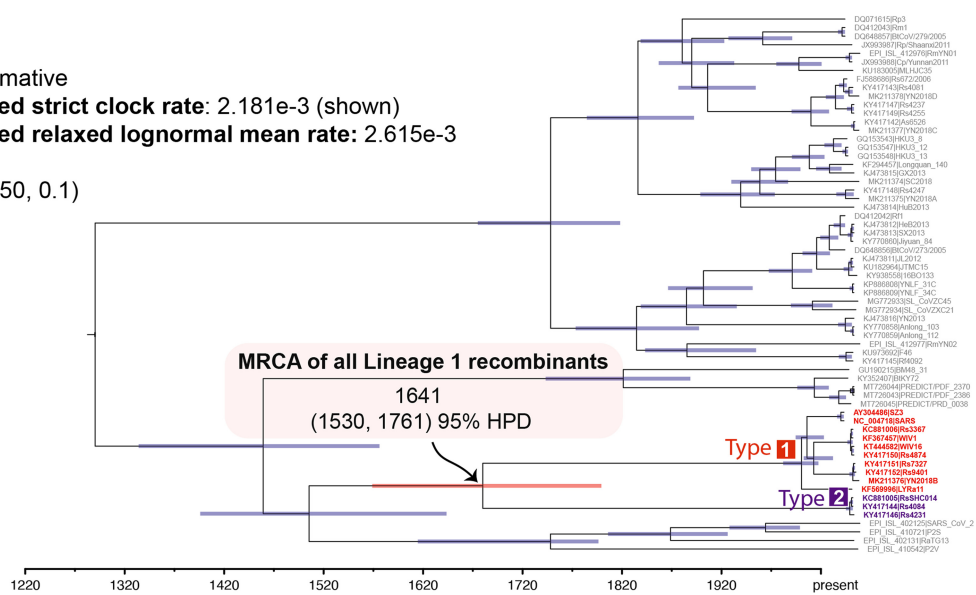
MRCA prior: none
 tMRCA: estimated 1250 AD



Region E (RBD)

Clock rate prior: uninformative
 Mean posterior estimated strict clock rate: 2.181e-3 (shown)
 Mean posterior estimated relaxed lognormal mean rate: 2.615e-3

MRCA prior: Laplace(1250, 0.1)
 tMRCA: fixed 1250 AD



Posterior distributions of rate estimates

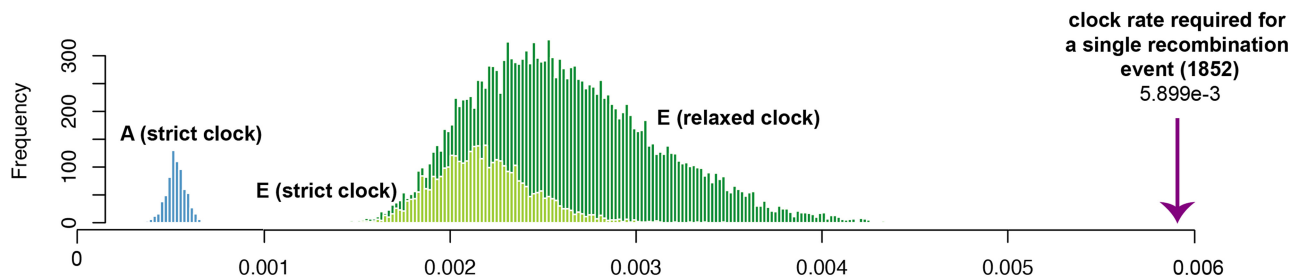


Figure 7. Time-calibrated phylogenies for recombination-free regions of the genome. Breakpoint-free regions A and E from Fig. 6 were chosen for time calibration since evidence of recombination was found in both RdRp and RBD. Both regions A and E were free of recombination for all sequences included in the tree, ensuring the best possible dating estimates. The MRCA of all Lineage 1 recombinants and its corresponding divergence date are labeled on each tree, demonstrating that the MRCA in region E (within the RBD) is much older than the MRCA in region A (proxy for RdRp, see Fig. 6). This suggests that there would not have been enough time for the RBDs of the recombinants to diversify to the extent shown here if only a single recombination event occurred between Lineage 5 and Lineage 1. The MRCAs of each type are labeled in red (Type 1) and purple (Type 2). Posterior distributions of rate estimates are also shown for each model as well as for a relaxed clock model of region E. For the observed sequence divergence in region E to have accumulated since the MRCA of the 13 recombinants in region A (1852), a clock rate of 5.899e-3 would be required, which is well outside the posterior distributions estimated by both our strict and relaxed clock models.

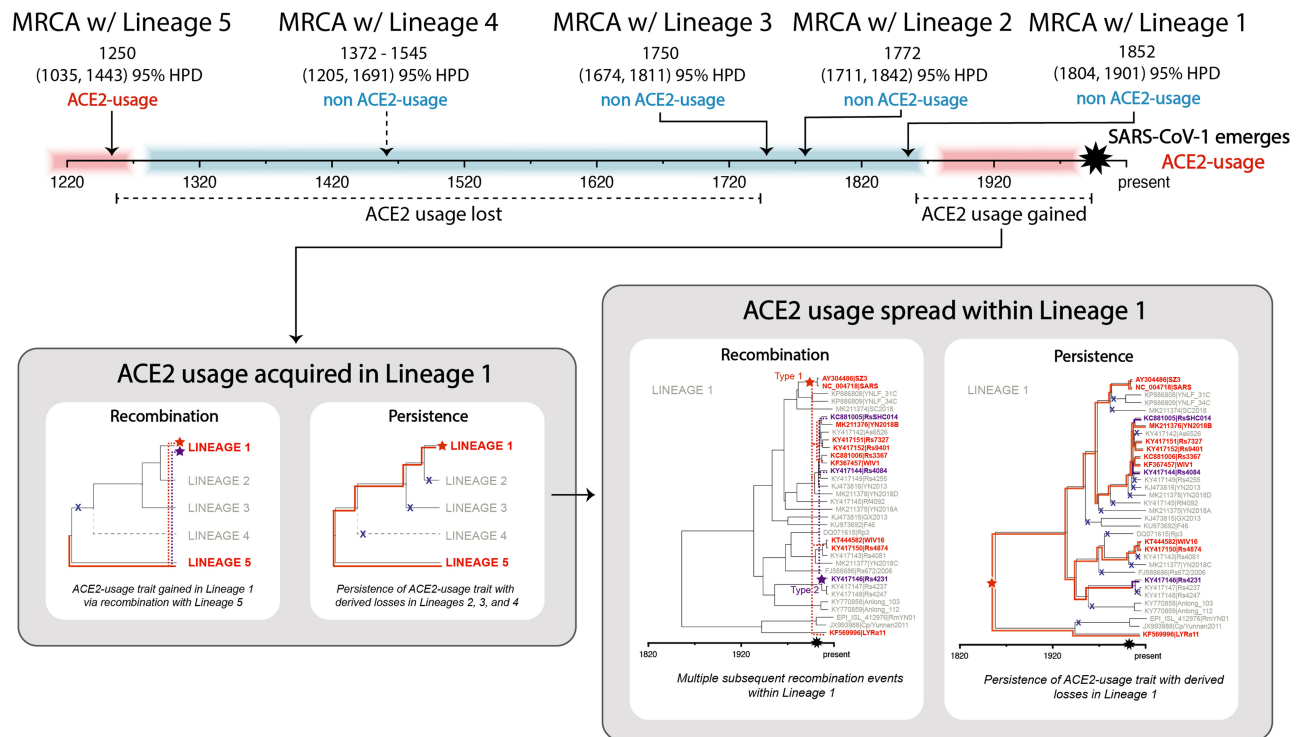


Figure 8. Proposed timeline of deletion and recombination events. The timeline demonstrates the sequence of events that led to loss of ACE2 usage in Lineages 2, 3, and 4 and gain of ACE2 usage within Lineage 1, leading to the emergence of SARS-CoV-1. Events are dated with MRCA age estimates; however, the exact intention is less to provide exact dates and more to suggest a particular order of events, which is strongly supported by the posterior probabilities of the time-calibrated phylogenies. The arrow for the Lineage 4 event is again dashed to demonstrate uncertainty in its positioning. We illustrate two hypotheses for the acquisition and subsequent spread of ACE2 usage in Lineage 1: recombination and persistence. The recombination hypothesis is much more parsimonious, as persistence would require multiple independent deletion events to generate the observed pattern of ACE2 usage.

altogether, but since mammalian ACE2 is so conserved (Damas et al. 2020; Lam et al. 2020) and ACE2-using viruses demonstrate broad host tropism (Li 2008; Hou et al. 2010; Zheng et al. 2020; Zhao et al. 2020), we hypothesize that there is likely a different receptor involved for the non-ACE2 users (see Supplementary File S1).

The difference in topology, specifically in the positioning of ACE2-using Lineage 1 viruses, between RdRp and RBD trees suggests that the ability to use ACE2 was introduced into Lineage 1 by recombination between a recent ancestor of the ACE2-using Lineage 1 viruses (including SARS-CoV-1) and an undiscovered Lineage 5 virus in the RBD. As there are two types of closely related RBD sequences in the recombinant Lineage 1 viruses (Fig. 2) with incompatible divergence dates (Fig. 7), we suggest that two such recombination events occurred between Lineage 1 and Lineage 5 (Fig. 8) independently introducing ACE2 usage into Lineage 1. The non-monophyletic nature of ACE2 usage within Lineage 1 can then be most parsimoniously explained by secondary intra-lineage recombination events (Fig. 8). It is possible that both hypotheses are partially true and that both intra-lineage recombination and the persistence of this trait alongside sister Lineage 1 viruses without the trait gave rise to the observed patterns of Type 1 and Type 2 ACE2 usage within Lineage 1. It is also very possible that further sampling may illuminate that some of the events proposed here have been distorted by sampling bias. We have estimated that these events may have occurred roughly within the last two centuries, though this estimate will likely change with further sampling as well. Our intention is not necessarily to date these events exactly, but rather to infer their order relative to each other and to make

hypotheses based on this order of events. Confidence intervals for many node dates overlap, but high posterior probabilities on internal nodes indicate that events most likely occurred in a certain order.

Our conclusion that ACE2 usage originated in Lineage 5 and was introduced into Lineage 1 by recombination is based on phylogenetics; however, studies of recombination using phylogenetics are often limited in their ability to definitively determine the direction of recombination. Nonetheless, there are several lines of evidence that support the direction having occurred from Lineage 5 to Lineage 1. First, recombination is notoriously more frequent in spike compared to orf1ab (Fu and Baric 1994; Boni et al. 2020; Ulferts et al. 2010). Second, Lineage 5 constitutes the base of the tree and has the oldest MRCA, meaning it likely shares more ancestral traits with the MRCA of all sarbecoviruses. Third, phylogenetic topology in orf1ab before the recombinant region of the genome mirrors that of S2 after the recombinant region (Fig. 6A), orienting orf1ab/S2 as sequence from the major parent of the recombination event. And finally, that spike is the recombinant region as opposed to RdRp is also supported by numerous studies that have provided evidence that SARS-CoV-1 is recombinant and SARS-CoV-2 is not (Hon et al. 2008; Lau et al. 2010; Zhou et al. 2020; Wu et al. 2020).

In order for recombination to have occurred between Lineage 1 and Lineage 5, these viruses must have had the opportunity to coinfect the same host cell. We demonstrate that recombination is possible given that viruses related to SARS-CoV-1 and SARS-CoV-2 appear to share both geographic and host space in southwestern China and in *R. sinicus* and *R. affinis* bats. Highlighting that this previously known recombination

event (i.e. SARS-CoV-1) occurred with a previously unknown group of viruses that are related to SARS-CoV-2 is an important finding of this study and demonstrates that recombination is an important driver of spillover for sarbecoviruses.

4.2 A series of deletion events most likely resulted in the ancestral loss of ACE2 usage in Lineages 1–4

Using the RdRp tree as the evolutionary history to which to compare because of its stability and relative lack of recombination, sequences without the deletions in the RBD most likely represent the ancestral state, as the SARS-CoV-2 Lineage 5 viruses at the base of the tree do not show this trait (Fig. 2). This is in accordance with the findings of Boni et al. (2020). Alternatively, it is possible that the deletion state is the ancestral state, and that this ancestral deletion state was conserved in Lineages 1, 2, and 3; however, insertions acquired during the evolution of Lineages 4 and 5 would have had to have occurred independently, which is less parsimonious. Persistence of the ACE2 usage trait from the MRCA of Lineage 5 all the way to Lineage 1 is also not parsimonious, as the RBD deletions would have had to have been lost many times independently (Fig. 8).

Further, the viruses from bats in Africa and Europe have one of the two deletions, which may indicate that these are descendant from an evolutionary intermediate and support a stepwise deletion hypothesis; however, this hypothesis hinges completely on the uncertain positioning of Lineage 4 on the phylogeny, which may support independent deletion within region 2 in Lineage 4 instead. Since ACE2-using Lineage 1 viruses including SARS-CoV-1 are nested within a clade of viruses that all have both deletions, this implies that both deletions arose before the diversification of Lineages 1, 2, and 3 viruses (Figs. 5 and 8). According to the branching order shown here, the smaller deletion in region 2 was likely acquired earliest, before the diversification of the clades into Africa and Europe, since it is shared by all clades with the exception of SARS-CoV-2 Lineage 5 at the base of the tree (Fig. 5). These large deletions in the RBD-ACE2 interface and the similarity of Rhinolophid and hACE2 also suggest that non-ACE2-using viruses, including Lineages 1, 2, 3, and 4, are using at least one receptor other than ACE2 (Letko, Marzi, and Munster 2020; Zhou et al. 2020).

4.3 ACE2 usage is not well explained by convergent evolution

Under a hypothetical convergent evolution scenario, large insertions would have had to be reacquired in precisely the same regions from which they were lost within the RBD independently in ACE2-using Lineage 1 viruses. The most parsimonious argument is that ACE2-using Lineage 1 viruses are descendant from at least two recombinant viruses (containing Types 1 and 2 RBDs) and that recombination best explains the non-monophyletic pattern of ACE2 usage within the *Sarbecovirus* subgenus. In contrast, human coronavirus NL63 is an alphacoronavirus that is also a hACE2 user but most likely represents a true case of convergent evolution. The RBD of SARS-CoV-1 and SARS-CoV-2 is structurally identical, while NL63 has a different structural fold, suggesting that they are not evolutionarily homologous (Chen et al. 2020). Nonetheless, NL63 also binds to hACE2 in the same region—suggesting all of the ACE2-using viruses have converged towards this interaction mode (Chen et al. 2020).

Additional evidence supports a recombination scenario over convergent evolution, including (1) the detection of statistically supported recombination breakpoints in all ACE2-using Lineage 1 viruses between RdRp and the RBD, and (2) a growing number of reports identifying recombination in the spike gene of other CoVs (Regan et al. 2012; Terada et al. 2014; Boniotti et al. 2016; Anthony et al. 2017a; Tao et al. 2017). We also highlight an additional unreported recombination event between Lineage 5 and Lineage 1 giving rise to RmYN02 that further demonstrates the importance of this evolutionary mechanism. We observed that the Lineage 5 bat virus RmYN02, which is highly similar to SARS-CoV-2 within the RdRp, actually has a RBD with the Lineage 1 deletion trait associated with the inability to use ACE2. This indicates a recombination in the opposite direction, from Lineage 1 to Lineage 5, and is again consistent with their overlapping host and geographic ranges. The RmYN02 virus was sequenced from a pooled sample that also contained a second strain, RmYN01, so the possibility that the assembled RmYN02 sequence is chimeric cannot be ruled out. However, both RmYN01 and RmYN02 have deletions in the RBD, so whether or not the sequence is chimeric, it is most likely still recombinant. Again, recombination is a much more parsimonious explanation for the loss of ACE2 usage in RmYN02 rather than convergence, which would require independent and identical deletions in the interfacial residues of the RBD.

4.4 Differences in receptor usage within Sarbecoviruses would explain observed phylogeographic patterns

Lineage 1 and Lineage 5 viruses appear to occupy the same geographic space, which is necessary for the opportunity to recombine to exist. However, the co-circulation of these distantly phylogenetically related viruses is a notable deviation from previous observations that show sarbecovirus phylogeny mirrors geography. It is unknown why Lineages 1–4 show strong phylogeographic clustering. Isolation by distance is one ecological mechanism that could explain concordance between phylogeny and geography; however, this would not explain why Lineage 5 deviates from this pattern and overlaps geographically with Lineage 1. Instead, we hypothesize that immune cross-reactivity between closely related viruses within hosts results in indirect competitive exclusion and priority effects, and that this explains the phylogeographic signal of Lineages 1–3. Antibodies against the spike protein are critical components of the immune response against CoVs (Buchholz et al. 2004; Lu et al. 2004; Prabakaran et al. 2006). Hosts that have been infected by one sarbecovirus may be immunologically resistant to infection from a related sarbecovirus, leading to geographic exclusion of closely related strains and a pattern of evolution that is concordant with geography despite the fact that species and individuals are not strictly confined (Fig. 1). It is unlikely that this pattern is caused by differing competencies amongst *Rhinolophus* bats, as host-switching of these viruses appears to be common. The co-circulation of Lineage 5 viruses (including SARS-CoV-2 and related viruses) in the same species and the same geographic location as Lineage 1 viruses may suggest a release in the competitive interactions maintaining geographic specificity. This would preclude recognition by cross-reactive antibodies, such as those produced against the spike protein, and may be evolutionarily advantageous for the recombinant virus. Furthermore, if these two groups of viruses utilize

different receptors, antibodies against one would be ineffective at excluding the other, potentially allowing both viral groups to infect the same hosts. If competitive release has indeed occurred among these viruses, it is likely that the SARS-CoV-2 clade is potentially much more diverse and geographically widespread than currently understood.

4.5 Implications for future research

Here, we highlight the critical need for further surveillance specifically in southwestern China and surrounding regions in Southeast Asia given that all ACE2-using bat viruses discovered to date were isolated from bats in Yunnan Province. If this holds true, it would support the hypothesis that SARS-CoV-2 originated in Yunnan or the surrounding regions of southwest China before the initial epidemic then amplified in Wuhan. Southeast Asia and parts of Europe and Africa have been previously identified as hotspots for sarbecoviruses (Anthony et al. 2017b), but increased surveillance will help characterize the true range of ACE2-using sarbecoviruses in particular. The receptors for viruses from northern China and other regions such as Europe and Africa remain unknown, and may not pose a threat to human health if they cannot utilize hACE2, though their potential to acquire hACE2 usage by recombination should be considered along with the potential for their existing spike proteins to use other human receptors for cell entry. It is unclear whether the lack of hACE2 binding for sarbecoviruses from Uganda and Rwanda is due to the small deletion in region 2 or to the numerous amino acid changes in other interfacial residues. It is possible that sarbecoviruses in Africa with different residues in these interfacial regions could potentially still use hACE2. It is also unknown whether the sarbecoviruses from Africa in particular use a different receptor altogether, or whether sarbecoviruses with the potential to utilize hACE2 without the region 2 deletion have also diversified into Africa or Europe. If competitive release between groups of viruses utilizing different receptors has indeed occurred, further surveillance is needed to determine the true extent of Lineage 5 viruses. In addition, experimental evidence to support or refute a competitive release hypothesis should be prioritized.

This study highlights that hACE2 usage is unpredictable using phylogenetic proximity to SARS-CoV-1 or SARS-CoV-2 in the RdRp gene. This is due to vastly different evolutionary histories in different parts of the viral genome due to recombination. Phylogenetic relatedness in the RdRp gene is not an appropriate proxy for pandemic potential among CoVs (the 'nearest neighbor' hypothesis). By extension, the consensus PCR assays most commonly used for surveillance and discovery, which mostly generate a small fragment of sequence from within this gene (De Souza Luna et al. 2007; Quan et al. 2010; Watanabe et al. 2010), are insufficient to predict hACE2 usage. Using phylogenetic distance in RdRp as a quantitative metric to predict the potential for emergence is tempting because of the large amount of data available, but this approach is unlikely to capture the biological underpinnings of emergence potential compared to more robust data sources such as full viral genome sequences. The current collection of full-length sarbecovirus genomes is heavily weighted toward China and *Rhinolophus* hosts, despite evidence of sarbecoviruses prevalent outside of China (such as in Africa) and in other mammalian hosts (such as pangolins). Further, investigations into determinants of pathogenicity and transmission for CoVs and the genomic signatures of such features will be an important step towards the prediction of viruses

with spillover potential, and distinguishing those with pandemic potential.

Finally, these findings reiterate the importance of recombination as a driver of spillover and emergence, particularly in the spike gene. If SARS-CoV-1 gained the ability to use hACE2 through recombination, other non-ACE2-using viruses could become human health threats through recombination as well. We know that recombination occurs much more frequently than just this single event with SARS-CoV-1, as the RdRp phylogeny does not mirror host phylogeny and the RBD tree has significantly different topology across all geographic lineages. In addition, the bat virus RmYN02 appears to be recombinant in the opposite direction (Lineage 5 backbone with Lineage 1 RBD) (Zhou et al. 2020), again supporting the hypothesis that recombination occurs between these lineages. Our analyses support two hypotheses: first, that sarbecoviruses frequently undergo recombination in this region of the genome, resulting in this pattern, and second, that sarbecoviruses are commonly shared amongst multiple host species, resulting in a lack of concordance with host species phylogeny and a reasonable opportunity for coinfection and recombination. Bats within the family *Rhinolophidae* have also repeatedly shown evidence of introgression between species (Mam et al. 2010; Mao et al. 2013a; Mao et al. 2013b; Mao et al. 2014; Mao, Zhang, and Rossiter 2016; Dool et al. 2016), supporting the hypothesis that many species in this family have close contact with one another which may facilitate viral host switching. Given that we have shown that ACE2-using viruses are co-occurring with a large diversity of non-ACE2-using viruses in Yunnan Province and in a similar host landscape, recombination poses a significant threat to the emergence of novel sarbecoviruses (Hu et al. 2017).

With recombination constituting such an important variable in the emergence of novel CoVs, understanding the genetic and ecological determinants of this process is a critical avenue for future research. Here we have shown not only that recombination was involved in the emergence of SARS-CoV-1, but also demonstrated how knowledge of the evolutionary history of these viruses can be used to infer the potential for other viruses to spillover and emerge. Understanding this evolutionary process is highly dependent on factors influencing viral co-occurrence and recombination, such as the geographic range of these viruses and their bat hosts, competitive interactions with co-circulating viruses within the same hosts, and the range of host species these viruses are able to infect. Our understanding depends on the data we have available—the importance of generating more data for such investigations cannot be understated. Investing effort now into further sequencing these viruses and describing the mechanisms that underpin their circulation and capacity for spillover will have important payoffs for predicting and preventing sarbecovirus pandemics in the future.

Acknowledgements

We thank three anonymous reviewers who provided thoughtful and robust suggestions that substantially improved this manuscript. The research reported in this publication was supported by the National Institute of Allergy and Infectious Diseases of the National Institutes of Health under Award Number R01AI149693 (PI S.J.A.). M.L. and V.J.M. are supported by the Intramural Research Program of the National Institute of Allergy and Infectious Diseases (NIAID), National Institutes of Health (NIH). G.L. and K.C. are

supported by National Institutes of Health (NIH) grant U19AI142777. This study was also made possible by the support of the American people through the United States Agency for International Development (USAID) Emerging Pandemic Threats PREDICT project, GHN-A-00-09-00010-00 (PI J.A.K.M.) and AID-OAA-A-14-00102 (PI J.A.K.M.). The content is solely the responsibility of the authors and does not necessarily represent the official views of the U.S. Government.

Data availability

All sequences have been submitted to GenBank (accession numbers MT726043-MT726045 MT738926-MT738928, MT732776, MW183243) and alignments used for phylogenetics are included as [supplementary materials](#).

Supplementary data

[Supplementary data](#) are available at *Virus Evolution* online.

Conflict of interest: None declared.

References

- Agnarsson, I. et al. (2011) 'A Time-Calibrated Species-Level Phylogeny of Bats (Chiroptera, Mammalia)', *PLoS Currents*, 3: RRN1212.
- Anthony, S. J. et al. (2017a) 'Further Evidence for Bats as the Evolutionary Source of Middle East Respiratory Syndrome Coronavirus'. *MBio*, 8: e00373–17.
- , PREDICT Consortium et al. (2017b) 'Global Patterns in Coronavirus Diversity', *Virus Evolution*, 3: vex012.
- Ar Gouilh, M. et al. (2018) 'SARS-CoV Related Betacoronavirus and Diverse Alphacoronavirus Members Found in Western Old-World', *Virology*, 517: 88–97.
- Boni, M. F. et al. (2020) 'Evolutionary Origins of the SARS-CoV-2 Sarbecovirus Lineage Responsible for the COVID-19 Pandemic', *Nature Microbiology*, 5: 1408–17.
- Boniotti, M. B. et al. (2016) 'Porcine Epidemic Diarrhea Virus and Discovery of a Recombinant Swine Enteric Coronavirus, Italy', *Emerging Infectious Diseases*, 22: 83–7.
- Bouckaert, R. et al. (2019) 'BEAST 2.5: An Advanced Software Platform for Bayesian Evolutionary Analysis', *PLoS Computational Biology*, 15: e1006650.
- Buchholz, U. J. et al. (2004) 'Contributions of the Structural Proteins of Severe Respiratory Syndrome Coronavirus to Protective Immunity', *Proceedings of the National Academy of Sciences of Sciences*, 101: 9804–9.
- Challender, D. et al. (2014) '*Manis javanica*'. The IUCN Red List of Threatened Species 2019: e.T12763A123584856.
- Chen, Y. et al. (2020) 'Structure Analysis of the Receptor Binding of 2019-NCoV', *Biochemical and Biophysical Research Communications*, 525: 135–40.
- Cui, J. et al. (2007) 'Evolutionary Relationships between Bat Coronaviruses and Their Hosts', *Emerging Infectious Diseases*, 13: 1526–32.
- Damas, J. et al. (2020) 'Broad Host Range of SARS-CoV-2 Predicted by Comparative and Structural Analysis of ACE2 in Vertebrates', *Proceedings of the National Academy of Sciences*, 117: 22311–22.
- Dool, S. E. et al. (2016) 'Nuclear Introns Outperform Mitochondrial DNA in Inter-Specific Phylogenetic Reconstruction: Lessons from Horseshoe Bats (Rhinolophidae: Chiroptera)', *Molecular Phylogenetics and Evolution*, 97: 196–212.
- Drexler, J. F. et al. (2010) 'Genomic Characterization of Severe Acute Respiratory Syndrome-Related Coronavirus in European Bats and Classification of Coronaviruses Based on Partial RNA-Dependent RNA Polymerase Gene Sequences', *Journal of Virology*, 84: 11336–49.
- Fu, K., and Baric, R. S. (1994) 'Map Locations of Mouse Hepatitis Virus Temperature-Sensitive Mutants: Confirmation of Variable Rates of Recombination', *Journal of Virology*, 68: 7458–66.
- Ge, X. Y. et al. (2013) 'Isolation and Characterization of a Bat SARS-like Coronavirus That Uses the ACE2 Receptor', *Nature*, 503: 535–8.
- Gorbalenya, A E. et al. (2020) 'The Species Severe Acute Respiratory Syndrome-Related Coronavirus: Classifying 2019-NCoV and Naming It SARS-CoV-2', *Nature Microbiology*, 5: 536–544.
- Graham, R. L., and Baric, R. S. (2010) 'Recombination, Reservoirs, and the Modular Spike: Mechanisms of Coronavirus Cross-Species Transmission', *Journal of Virology*, 84: 3134–46.
- He, B. et al. (2014) 'Identification of Diverse Alphacoronaviruses and Genomic Characterization of a Novel Severe Acute Respiratory Syndrome-Like Coronavirus from Bats in China', *Journal of Virology*, 88: 7070–82.
- Hon, C.-C. et al. (2008) 'Evidence of the Recombinant Origin of a Bat Severe Acute Respiratory Syndrome (SARS)-like Coronavirus and Its Implications on the Direct Ancestor of SARS Coronavirus', *Journal of Virology*, 82: 1819–26.
- Hou, Y. et al. (2010) 'Angiotensin-Converting Enzyme 2 (ACE2) Proteins of Different Bat Species Confer Variable Susceptibility to SARS-CoV Entry', *Archives of Virology*, 155: 1563–9.
- Hu, B. et al. (2017) 'Discovery of a Rich Gene Pool of Bat SARS-Related Coronaviruses Provides New Insights into the Origin of SARS Coronavirus', *PLoS Pathogens*, 13: e1006698.
- Lam, H. M., Ratmann, O., and Boni, M. F. (2018) 'Improved Algorithmic Complexity for the 3SEQ Recombination Detection Algorithm', *Molecular Biology and Evolution*, 35: 247–51.
- Lam, S. D. et al. (2020) 'SARS-CoV-2 Spike Protein Predicted to Form Complexes with Host Receptor Protein Orthologues from a Broad Range of Mammals', *Scientific Reports*, 10: 16471.
- Lam, T., et al. (2020) 'Identifying SARS-CoV-2 Related Coronaviruses in Malayan Pangolins', *Nature*, 583: 282–285.
- Lan, J. et al. (2020) 'Structure of the SARS-CoV-2 Spike Receptor-Binding Domain Bound to the ACE2 Receptor', *Nature*, 581: 215–220.
- Lau, S. K. P. et al. (2010) 'Ecoepidemiology and Complete Genome Comparison of Different Strains of Severe Acute Respiratory Syndrome-Related Rhinolophus Bat Coronavirus in China Reveal Bats as a Reservoir for Acute, Self-Limiting Infection That Allows Recombination Events', *Journal of Virology*, 84: 2808–19.
- et al. (2005) 'Severe Acute Respiratory Syndrome Coronavirus-like Virus in Chinese Horseshoe Bats', *Proceedings of the National Academy of Sciences of the United States of America*, 102: 14040–5.
- Lecis, R. et al. (2019) 'Molecular Identification of Betacoronavirus in Bats from Sardinia (Italy): First Detection and Phylogeny', *Virus Genes*, 55: 60–7.
- Lefkowitz, E.J. et al. (2018) 'Virus taxonomy: the database of the International Committee on Taxonomy of Viruses (ICTV)', *Nucleic Acids Research*, 46: D708–D717.

- Lelli, D. et al. (2013) 'Detection of Coronaviruses in Bats of Various Species in Italy', *Viruses*, 5: 2679–89.
- Leopardi, S. et al. (2018) 'Interplay between Co-Divergence and Cross-Species Transmission in the Evolutionary History of Bat Coronaviruses', *Infection, Genetics and Evolution*, 58: 279–89.
- Letko, M., Marzi, A., and Munster, V. (2020) 'Functional Assessment of Cell Entry and Receptor Usage for SARS-CoV-2 and Other Lineage B Betacoronaviruses', *Nature Microbiology*, 5: 562–9.
- et al. (2018) 'Adaptive Evolution of MERS-CoV to Species Variation in DPP4', *Cell Reports*, 24: 1730–7.
- Li, D. et al. (2016) 'MEGAHIT v1.0: A Fast and Scalable Metagenome Assembler Driven by Advanced Methodologies and Community Practices', *Methods*, 102: 3–11.
- Li, F. (2008) 'Structural Analysis of Major Species Barriers between Humans and Palm Civets for Severe Acute Respiratory Syndrome Coronavirus Infections', *Journal of Virology*, 82: 6984–91.
- et al. (2005) 'Structural Biology: Structure of SARS Coronavirus Spike Receptor-Binding Domain Complexed with Receptor', *Science*, 309: 1864–8.
- Li, W. et al. (2005) 'Bats Are Natural Reservoirs of SARS-like Coronaviruses', *Science*, 310: 676–9.
- et al. (2003) 'Angiotensin-Converting Enzyme 2 is a Functional Receptor for the SARS Coronavirus', *Nature*, 426: 450–4.
- et al. (2005) 'Receptor and Viral Determinants of SARS-Coronavirus Adaptation to Human ACE2', *The Embo Journal*, 24: 1634–43.
- Liu, P., et al. (2019) 'Viral Metagenomics Revealed Sendai Virus and Coronavirus Infection of Malayan Pangolins (*Manis javanica*)', *Viruses*, 11: 979.
- Lu, G., Wang, Q., and Gao, G. F. (2015) 'Bat-to-Human: Spike Features Determining 'Host Jump' of Coronaviruses SARS-CoV, MERS-CoV, and Beyond', *Trends in Microbiology*, 23: 468–78.
- Lu, L. et al. (2004) 'Immunological Characterization of the Spike Protein of the Severe Acute Respiratory Syndrome Coronavirus', *Journal of Clinical Microbiology*, 42: 1570–6.
- Luk, H. K. H. et al. (2019) 'Molecular Epidemiology, Evolution and Phylogeny of SARS Coronavirus', *Infection, Genetics and Evolution*, 71: 21–30.
- MAO, X. et al. (2010) 'Historical Male-Mediated Introgression in Horseshoe Bats Revealed by Multilocus DNA Sequence Data', *Molecular Ecology*, 19: 1352–66.
- Mao, X. et al. (2013a) 'Lineage Divergence and Historical Gene Flow in the Chinese Horseshoe Bat (*Rhinolophus sinicus*)', *PLoS ONE*, 8: e56786.
- et al. (2013b) 'Multiple Cases of Asymmetric Introgression among Horseshoe Bats Detected by Phylogenetic Conflicts across Loci', *Biological Journal of the Linnean Society*, 110: 346–61.
- , Zhang, S., and Rossiter, S. J. (2016) 'Differential Introgression Suggests Candidate Beneficial and Barrier Loci between Two Parapatric Subspecies of Pearson's Horseshoe Bat *Rhinolophus*', *Current Zoology*, 62: 405–12.
- et al. (2014) 'Differential Introgression among Loci across a Hybrid Zone of the Intermediate Horseshoe Bat (*Rhinolophus affinis*)', *BMC Evolutionary Biology*, 14: 154.
- Menachery, V. D., Graham, R. L., and Baric, R. S. (2017) 'Jumping Species—a Mechanism for Coronavirus Persistence and Survival', *Current Opinion in Virology*, 23: 1–7.
- et al. (2015) 'A SARS-like Cluster of Circulating Bat Coronaviruses Shows Potential for Human Emergence', *Nature Medicine*, 21: 1508–13.
- Pettersen, E. F. et al. (2004) 'UCSF Chimera - A Visualization System for Exploratory Research and Analysis', *Journal of Computational Chemistry*, 25: 1605–12.
- Pieper, U. et al. (2011) 'ModBase, a Database of Annotated Comparative Protein Structure Models, and Associated Resources', *Nucleic Acids Research*, 39: D465–D474.
- Prabakaran, P. et al. (2006) 'Structure of Severe Acute Respiratory Syndrome Coronavirus Receptor-Binding Domain Complexed with Neutralizing Antibody', *Journal of Biological Chemistry*, 281: 15829–36.
- , Xiao, X., and Dimitrov, D. S. (2004) 'A Model of the ACE2 Structure and Function as a SARS-CoV Receptor', *Biochemical and Biophysical Research Communications*, 314: 235–41.
- Quan, P. L., et al. (2010) 'Identification of a Severe Acute Respiratory Syndrome Coronavirus-like Virus in a Leaf-Nosed Bat in Nigeria', *MBio*, 1: e00208–10.
- Regan, A. D. et al. (2012) 'Characterization of a Recombinant Canine Coronavirus with a Distinct Receptor-Binding (S1) Domain', *Virology*, 430: 90–9.
- Ren, W. et al. (2008) 'Difference in Receptor Usage between Severe Acute Respiratory Syndrome (SARS) Coronavirus and SARS-Like Coronavirus of Bat Origin', *Journal of Virology*, 82: 1899–907.
- Rihtarič, D. et al. (2010) 'Identification of SARS-like Coronaviruses in Horseshoe Bats (*Rhinolophus hipposideros*) in Slovenia', *Archives of Virology*, 155: 507–14.
- Shang, J. et al. (2020) 'Structural Basis of Receptor Recognition by SARS-CoV-2', *Nature*, 581: 221–4.
- Souza, L. et al. (2007) 'Generic Detection of Coronaviruses and Differentiation at the Prototype Strain Level by Reverse Transcription-PCR and Nonfluorescent Low-Density Microarray', *Journal of Clinical Microbiology*, 45: 1049–52.
- Su, S. et al. (2016) 'Epidemiology, Genetic Recombination, and Pathogenesis of Coronaviruses', *Trends in Microbiology*, 24: 490–502.
- Takada, A. et al. (1997) 'A System for Functional Analysis of Ebola Virus Glycoprotein', *Proceedings of the National Academy of Sciences of Sciences*, 94: 14764–9.
- Tang, X. et al. (2009) 'Differential Stepwise Evolution of SARS Coronavirus Functional Proteins in Different Host Species', *BMC Evolutionary Biology*, 9: 52.
- Tao, Y. et al. (2017) 'Surveillance of Bat Coronaviruses in Kenya Identifies Relatives of Human Coronaviruses NL63 and 229E and Their Recombination History', *Journal of Virology*, 91: e01953–16.
- , and Tong, S. (2019) 'Complete Genome Sequence of a Severe Acute Respiratory Syndrome-Related Coronavirus from Kenyan Bats', *Microbiology Resource Announcements*, 8: e00548–19.
- Terada, Y. et al. (2014) 'Emergence of Pathogenic Coronaviruses in Cats by Homologous Recombination between Feline and Canine Coronaviruses', *PLoS ONE*, 9: e106534.
- Ulferts, R. et al. (2010) 'Expression and Functions of SARS Coronavirus Replicative Proteins', in *Molecular Biology of the SARS-Coronavirus*, 75–98. Berlin, Heidelberg: Springer.
- Wan, Y. et al. (2020) 'Receptor Recognition by Novel Coronavirus from Wuhan: An Analysis Based on Decade-Long Structural Studies of SARS', *Journal of Virology*, 94: e00127–20.
- Watanabe, S. et al. (2010) 'Bat Coronaviruses and Experimental Infection of Bats, the Philippines', *Emerging Infectious Diseases*, 16: 1217–23.
- Woo, P. C. Y. et al. (2009) 'Coronavirus Diversity, Phylogeny and Interspecies Jumping', *Experimental Biology and Medicine*, 234: 1117–27.

- Wrobel, A. G. et al. (2021) 'Structure and Binding Properties of Pangolin-CoV Spike Glycoprotein Inform the Evolution of SARS-CoV-2', *Nature Communications*, 12: 837.
- Wu, F. et al. (2020) 'A New Coronavirus Associated with Human Respiratory Disease in China', *Nature*, 579: 265–9.
- Wu, K. et al. (2012) 'Mechanisms of Host Receptor Adaptation by Severe Acute Respiratory Syndrome Coronavirus', *Journal of Biological Chemistry*, 287: 8904–11.
- Yang, X.-L. et al. (2016) 'Isolation and Characterization of a Novel Bat Coronavirus Closely Related to the Direct Progenitor of Severe Acute Respiratory Syndrome Coronavirus', *Journal of Virology*, 90: 3253–6.
- Yu, P. et al. (2019) 'Geographical Structure of Bat SARS-Related Coronaviruses', *Infection, Genetics and Evolution*, 69: 224–9.
- Yuan, J. et al. (2010) 'Intraspecies Diversity of SARS-like Coronaviruses in *Rhinolophus Sinicus* and Its Implications for the Origin of SARS Coronaviruses in Humans', *Journal of General Virology*, 91: 1058–62.
- Zhao, X. et al. (2020) 'Broad and Differential Animal Angiotensin-Converting Enzyme 2 Receptor Usage by SARS-CoV-2', *Journal of Virology*, 94: e00940–20.
- Zheng, M. et al. (2020) 'Bat SARS-Like WIV1 Coronavirus Uses the ACE2 of Multiple Animal Species as Receptor and Evades IFITM3 Restriction via TMPRSS2 Activation of Membrane Fusion', *Emerging Microbes & Infections*, 9: 1567–79.
- Zhou, H. et al. (2020) 'A Novel Bat Coronavirus Closely Related to SARS-CoV-2 Contains Natural Insertions at the S1/S2 Cleavage Site of the Spike Protein', *Current Biology*, 30: 3896.
- Zhou, P. et al. (2020) 'A Pneumonia Outbreak Associated with a New Coronavirus of Probable Bat Origin', *Nature*, 588: E6–E6.