

# Evidence for adaptive evolution in the receptor-binding domain of seasonal coronaviruses

Kathryn E. Kistler<sup>\*1,2</sup>, Trevor Bedford<sup>1,2</sup>

<sup>1</sup>*Vaccine and Infectious Disease Division, Fred Hutchinson Cancer Research Center, Seattle, WA USA*

<sup>2</sup>*Molecular and Cellular Biology Program, University of Washington, Seattle, WA USA*

\* To whom correspondence should be addressed.

## Abstract

Seasonal coronaviruses (OC43, 229E, NL63 and HKU1) are endemic to the human population, regularly infecting and reinfecting humans while typically causing asymptomatic to mild respiratory infections. It is not known to what extent reinfection by these viruses is due to waning immune memory or antigenic drift of the viruses. Here, we address the influence of antigenic drift on immune evasion of seasonal coronaviruses. We provide evidence that at least two of these viruses, OC43 and 229E, are undergoing adaptive evolution in regions of the viral spike protein that are exposed to human humoral immunity. This suggests that reinfection may be due, in part, to positively-selected genetic changes in these viruses that enable them to escape recognition by the immune system. It is possible that, as with seasonal influenza, these adaptive changes in antigenic regions of the virus would necessitate continual reformulation of a vaccine made against them.

## Introduction

Coronaviruses were first identified in the 1960s and, in the decades that followed, human coronaviruses (HCoV) received a considerable amount of attention in the field of infectious disease research. At this time, two species of HCoV, OC43 and 229E were identified as the causative agents of roughly 15% of common colds (McIntosh 1974; Heikkinen and Järvinen 2003). Infections with these viruses were shown to exhibit seasonal patterns, peaking in January-March in the Northern Hemisphere, as well as yearly variation, with the greatest incidence occurring every 2-4 years (Monto and Lim 1974; Hamre and Beem 1972). Subsequently, two additional seasonal HCoVs, HKU1 and NL63, have entered the human population. These 4 HCoVs endemic to the human population usually cause mild respiratory infections, but occasionally result in more severe disease in immunocompromised patients or the elderly (Liu, Liang, and Fung 2020). In the past 20 years, three additional HCoVs (SARS-CoV-1, MERS-CoV and SARS-CoV-2) have emerged, which cause more severe respiratory illness. At the writing of this paper, amidst the SARS-CoV-2 pandemic, no vaccine for any HCoV is currently available, though many candidate SARS-CoV-2 vaccines are in production and clinical trials (Krammer 2020).

Coronaviruses are named for the ray-like projections of spike protein that decorate their surface. Inside these virions is a positive-sense RNA genome of roughly 30kB (Li 2016). This large genome size can accommodate more genetic variation than a smaller genome (Woo et al. 2009). Genome flexibility, coupled with a RNA virus error-prone polymerase (Drake 1993) and a high rate of homologous recombination (Pasternak, Spaan, and Snijder 2006), creates genetic diversity that is acted upon by evolutionary pressures that select for viral replication. This spawns much of the diversity within and between coronaviruses species (Woo et al. 2009; Hon et al. 2008), and can contribute to the virus' ability to jump species-barriers, allowing a previously zoonotic CoV to infect and replicate in humans.

The battle between virus and host results in selective pressure for mutations that alter viral antigens in a way that evades immune recognition. Antigenic evolution, or antigenic drift, leaves a characteristic mark of positively selected epitopes within the viral proteins most exposed to the host immune system (Smith et al. 2004). For CoVs, this is the spike protein, exposed on the surface of the virion to human humoral immunity. Some human respiratory illnesses caused by RNA viruses, like seasonal influenza (Smith et al. 2004), evolve antigenically while others, like measles, do not (Fulton et al. 2015). Because of this, seasonal influenza vaccines must be reformulated on a nearly annual basis, while measles vaccines typically provide lifelong protection. Whether HCoVs undergo antigenic drift is relevant not only to understanding HCoV evolution and natural immunity against HCoVs, but also to predicting the duration of a vaccine's effectiveness.

Early evidence that closely-related HCoVs are antigenically diverse comes from a 1980s human challenge study in which subjects were infected and then reinfected with a variety of 229E-related strains (Reed 1984). All subjects developed symptoms and shed virus upon initial virus inoculation. After about a year, subjects who were re-inoculated with the same strain did not show symptoms or shed virus. However, the majority of subjects who were re-inoculated with a heterologous strain developed symptoms and shed virus. This suggests that immunity mounted against 229E viruses provides protection against some, but not all, other 229E strains. This is a result that would be expected of an antigenically evolving virus.

More recent studies have identified 8 OC43 genotypes and, in East Asian populations, certain genotypes were shown to temporally replace other genotypes (Lau et al. 2011; Zhang et al. 2015; Zhu et al. 2018). Whether certain genotypes predominate due to antigenic differences that confer a fitness advantage is not known. However, evidence for selection in the spike protein of one of these dominant OC43 genotypes has been provided by  $dN/dS$ , a standard computational method for detecting positive selection (Ren et al. 2015). This method has also been used to suggest positive selection in the spike protein of 229E (Chibo and Birch 2006). Additionally, two genetically distinct groupings (each of which include multiple of the aforementioned 8 genotypes) of OC43 viruses have been shown to alternate in prevalence within a Japanese community, meaning that the majority of OC43 infections are caused by one group for about 2-4 years at which point the other group begins to account for the bulk of

infections. It has been suggested that antigenic differences between these groups contribute to this epidemic switching (Komabayashi et al. 2020).

Here, we use a variety of computational approaches to detect adaptive evolution in spike and comparator proteins in HCoV. These methods were designed as improvements to  $dN/dS$  with the intention of identifying positive selection within a serially-sampled RNA virus population. We focus on the seasonal HCoVs that have been continually circulating in humans: OC43, 229E, HKU1 and NL63. Our analyses of nonsynonymous divergence, rate of adaptive substitutions, and Time to Most Recent Ancestor (TMRCA) provide evidence that the spike protein of OC43 and 229E is under positive selection. Though we conduct these analyses on HKU1 and NL63, we do not observe evidence for adaptive evolution in the spike protein of these viruses. For HKU1, there is not enough longitudinal sequencing data available for us to confidently make conclusions as to whether or not this lack of evidence reflects an actual lack of adaptive evolution.

## Results

We constructed time-resolved phylogenies of the OC43 and 229E using publicly accessible sequenced isolates. A cursory look at these trees confirms previous reports that substantial diversity exists within each viral species (Zhang et al. 2015; Komabayashi et al. 2020; Lau et al. 2011). The phylogeny of OC43 bifurcates immediately from the root (Fig. 1), indicating that OC43 consists of multiple, co-evolving lineages. Because of the distinct evolutionary histories, it is appropriate to conduct phylogenetic analyses separately for each lineage. We have arbitrarily labeled these lineages 'A' and 'B' (Fig. 1).

Because recombination is common amongst coronaviruses (Pasternak, Spaan, and Snijder 2006; Hon et al. 2008; Lau et al. 2011), we built separate phylogenies for each viral gene. In the absence of recombination, each tree should show the same evolutionary relationships between viral isolates. A dramatic difference in a given isolate's position on one tree versus another is strongly indicative of recombination (Kosakovsky Pond et al. 2006). Comparing the RNA-dependent RNA polymerase (RdRp) and spike trees reveals this pattern of recombination in some isolates (Fig. 1 Supplement 1A). A comparison of the trees of the S1 and S2 sub-domains of spike shows more limited evidence for intragenic recombination (Fig. 1 Supplement 1B), which is consistent with the fact that the distance between two genetic loci is inversely-related to the chance that these loci remain linked during a recombination event. Though intragenic recombination likely does occur occasionally, analyzing genes, rather than isolates, greatly reduces the contribution of recombination to genetic variation in our analyses. Because of this, we designate the lineage of each gene separately, based on that gene's phylogeny. Though most isolates contain all genes from the same lineage, some isolates have, say, a lineage A spike gene and a lineage B RdRp gene. This allows us to consider the evolution of each gene separately, and interrogate the selective pressures acting on them. Because of its essential role in viral replication and lack of antibody exposure, we expect RdRp to be under purifying selection to maintain its structure and function. If HCoVs evolve antigenically, we expect to see adaptive evolution in spike, and particularly in the S1 domain of

Spike (Hofmann et al. 2006; Hulswit et al. 2019), due to its exposed location at the virion's surface and interaction with the host receptor.

Using phylogenies constructed from the spike gene, we tallied the number of independent amino acid substitutions at each position within spike. The average number of mutations per site is higher in S1 than S2 for HCoV lineages in OC43 and 229E (Fig. 2A). A greater occurrence of repeated mutations is expected if some mutations within S1 confer immune avoidance. Not only should S1 contain more repeated mutations, but we would also expect these mutations to spread widely after they occur due to their selective advantage. Additionally, we expect sites within S1 to experience diversifying selection due to the ongoing arms race between virus and host immune system. This is visible in the distribution of genotypes at the most repeatedly-mutated sites in OC43 lineage A (Fig. 2B and 2C).

An adaptively evolving gene, or region of the genome, should exhibit a high rate of nonsynonymous substitutions. For each seasonal HCoV lineage, we calculated nonsynonymous and synonymous divergence as the average Hamming distance from that lineage's common ancestor (Zanini et al. 2015). The rate of nonsynonymous divergence is markedly higher within spike versus RdRp of 229E and OC43 lineage A (Fig. 3A). While nonsynonymous divergence increases steadily over time in spike, it remains roughly constant at 0.0 in RdRp, while rates of synonymous evolution are similar between spike and RdRp. These results suggest that there is predominantly positive selection on OC43 and 229E spike, but predominantly purifying selection on RdRp. Separating spike into the S1 (receptor-binding) and S2 (membrane-fusion) domains reveals that the majority of nonsynonymous divergence in spike occurs within S1 (Fig. 3B). In fact, the rates of nonsynonymous divergence in S2 are similar to those seen in RdRp, suggesting S2 evolves under purifying selection while S1 evolves adaptively.

As a complement to the divergence analysis, we implemented an alternative to the  $dN/dS$  method that was specifically designed to detect positive selection within RNA virus populations (Bhatt, Holmes, and Pybus 2011). Compared with traditional  $dN/dS$  methods, the Bhatt method has the advantages of: 1) measuring the strength of positive selection within a population given sequences collected over time, 2) higher sensitivity to identifying mutations that occur only once and sweep through the population, and 3) correcting for deleterious mutations (Bhatt, Katzourakis, and Pybus 2010; Bhatt, Holmes, and Pybus 2011). We adapted this method to detect adaptive substitutions in seasonal HCoVs and compare these rates to H3N2, the canonical example of antigenic evolution (Rambaut et al. 2008; Yang 2000). As shown in Figure 4, OC43 lineage A has continuously amassed adaptive substitutions in spike over the past >30 years while RdRp has accrued few, if any, adaptive substitutions. These adaptive substitutions are located within the S1, and not the S2, domain of spike (Fig. 4). We observe a largely linear accumulation of adaptive substitutions in spike and S1 through time, although the method does not dictate a linear increase.

We estimate that OC43 lineage A accumulates roughly  $0.6 \times 10^{-3}$  adaptive substitutions per codon per year (or 0.45 adaptive substitutions each year) in the S1 domain of spike while the

rate of adaptation in OC43 lineage B is slightly higher and is estimated to result in an average 0.56 adaptive substitutions in S1 per year (Fig. 5). The S1 domain of 229E is estimated to accrue 0.26 adaptive substitutions per year. A benefit of the Bhatt method is the ability to calculate the strength of selection, which allows us to compare these seasonal HCoVVs to other viruses. We used our implementation of the Bhatt method to calculate the rate of adaptation for influenza H3N2, which is known to undergo rapid antigenic evolution, and measles, which does not. We estimate that the receptor-binding domain of influenza H3N2 accumulates adaptive substitutions about 3 times faster than the HCoVVs OC43 and 229E (Fig. 6). We detect no adaptive substitutions in the measles receptor-binding protein. These results put the evolution of the S1 domain of OC43 and 229E in context, indicating that the S1 domain is under positive selection, and that this positive selection generates new variants in the putative antigenic regions of these HCoVVs at about a third of the rate of the canonical example of antigenic evolution, the HA1 domain of influenza H3N2.

Because coronaviruses are known to recombine, and recombination has the potential to impact evolutionary analyses of selection, we sought to verify that our results are not swayed by the presence of recombination. To do this, we simulated the evolution of OC43 lineage A spike and RdRp genes under varying levels of recombination and positive selection and used our implementation of the Bhatt method to identify adaptive substitutions. As the strength of positive selection increases, we detect more adaptive substitutions, regardless of the level of recombination (Fig. 7). This demonstrates that our estimates of adaptive evolution are not biased by recombination events.

Finally, we know that strong directional selection skews the shape of phylogenies (Volz, Koelle, and Bedford 2013). In influenza H3N2, immune selection causes the genealogy to adopt a ladder-like shape where the rungs are formed by viral diversification and each step is created by the appearance of new, antigenically-superior variants that replace previous variants. This ladder-like shape can also be seen in the phylogenies of the OC43 and 229E (Fig. 1). In this case, selection can be quantified by the timescale of population turnover as measured by the Time to Most Recent Common Ancestor (TMRCA), with the expectation that stronger selection will result in more frequent steps and therefore a smaller TMRCA measure (Bedford, Cobey, and Pascual 2011). We computed average TMRCA values from phylogenies built on Spike, S1, S2 or RdRp sequences of OC43 and 229E (Table 1). We observe that, for both OC43 lineage A and 229E, the average TMRCA is lower in spike than RdRp and lower in S1 versus S2. These results suggest strong directional selection in S1, likely driven by pressures to evade the humoral immune system. The difference in TMRCA between S1 and S2 is indicative not only of differing selective pressures acting on these two spike domains, but also of intra-spike recombination, which emphasizes the importance of using methods that are robust to recombination to detect adaptive evolution.

Because HKU1 was identified in the early 2000's, there are fewer longitudinally-sequenced isolates available for this HCoV compared to 229E and OC43 (Fig. 1 Supplement 2). Consequently, the phylogenetic reconstructions and divergence analysis of HKU1 have a higher

level of uncertainty. To begin with, it is less clear from the phylogenies whether HKU1 represents a single HCoV lineage like 229E or, instead, should be split into multiple lineages like OC43 (Fig. 1). Because of this, we completed all antigenic analyses for HKU1 twice: once considering all isolates to be members of a single lineage, and again after splitting isolates into 2 separate lineages. These lineages are arbitrarily labeled 'A' and 'B' as was done for OC43. When HKU1 is considered to consist of just one lineage, there is no signal of antigenic evolution by divergence analysis (Fig. 3 Supplement 1B) or by the Bhatt method of estimating adaptive evolution (Fig. 5 Supplement 1A). However, when HKU1 is assumed to consist of 2 co-circulating lineages, HKU1 lineage A has a markedly higher rate of adaptive substitutions in S1 than in S2 or RdRp (Fig. 5 Supplement 1B).

To demonstrate the importance of having a well-sampled longitudinal series of sequenced isolates for our antigenic analyses, we returned to our simulated OC43 S1 datasets. We mimicked shorter longitudinal series by truncating the dataset to only 24, 14, 10, or 7 years of samples and ran the Bhatt analysis on these sequentially shorter time series (Fig. 7 Supplement). This simulated data reveals a general trend that less longitudinal data reduces the ability to detect adaptive evolution and increases the uncertainty of the analysis. Given the dearth of longitudinal data for HKU1, we do not feel that it is appropriate to make strong conclusions about adaptive evolution, or lack thereof, in this HCoV.

Despite being identified at roughly the same time as HKU1, substantially more NL63 isolates have been sequenced (Fig. 1 Supplement 2) making the phylogenetic reconstruction and evolutionary analyses of this virus correspondingly more reliable. We do not observe evidence for adaptive evolution in NL63 (Fig. 3 Supplement 1A and Fig. 5 Supplement 1A) and this lack of support for adaptive evolution in the NL63 spike gene is more likely to reflect an actual lack of adaptive evolution in this virus.

## Discussion

Using several corroborating methods, we provide evidence that the seasonal HCoVs OC43 and 229E undergo adaptive evolution in S1, the region of the spike protein exposed to human humoral immunity (Figs. 3, 4 and 5). We additionally confirm that RdRp and S2 do not show signals of adaptive evolution. We observe that S1 accumulates between 0.3 (229E) and 0.5 (OC43) adaptive substitutions per year. We infer that these viruses accumulate adaptive substitutions at roughly a third of the rate of influenza H3N2 (Fig. 6). The most parsimonious explanation for the observation of substantial adaptive evolution in S1 is that antigenic drift is occurring in which mutations that escape from human population immunity are selectively favored in the viral population leading to repeated adaptive changes. However, it is formally possible that the adaptive evolution we detect is a result of selective pressures other than evasion of the adaptive immune system. Showing that this is truly antigenic evolution could involve a serological comparison of isolates that differ at S1 residues under positive selection. We do not observe evidence of antigenic evolution in NL63 or HKU1 (Figs. 3 and 5 Supplements). For NL63, this likely represents a true lack of marked adaptive evolution in S1. There is much less longitudinal sequencing data available for HKU1 and it is possible that a

more completely sampled time series of genome sequences could alter the result for this virus (Fig. 7 Supplement 1).

Our conclusions of adaptive evolution in S1, arrived at through computational analyses of sequencing data, agree with studies that observe reinfection of subjects by heterologous isolates of 229E (Reed 1984), sequential dominance of specific genotypes of OC43 (Lau et al. 2011; Zhang et al. 2015), and common reinfection by seasonal HCoV from longitudinal serological data (Edridge et al. 2020). In this latter study, HCoV infections were identified from longitudinal serum samples by assaying for increases in antibodies against the nucleocapsid (N) protein of representative OC43, 229E, HKU1, and NL63 viruses. This study concluded that the average time between infections was 1.5–2.5 years, depending on the HCoV (Edridge et al. 2020). In comparison, influenza H3N2 reinfects people roughly every 5 years (Kucharski et al. 2015). Thus, frequent reinfection by seasonal HCoVs is likely due to a combination of factors and suggests waning immune memory, and/or incomplete immunity against reinfection, in addition to antigenic drift.

Human coronaviruses are a diverse grouping split, phylogenetically, into two genera: NL63 and 229E are alphacoronaviruses, while OC43, HKU1, MERS, SARS, and SARS-CoV-2 are betacoronaviruses. Transmissibility and pathology do not seem to correlate with genus, nor does the method of cell-entry. Coronaviruses bind to a remarkable range of host-cell receptors including peptidases, cell adhesion molecules and sugars. Amongst the seasonal HCoVs, OC43 and HKU1 both bind 9-O-acetylsialic acid (Hulswit et al. 2019) while 229E binds human aminopeptidase N (hAPN) and NL63 binds angiotensin-converting enzyme 2 (ACE2) (Liu, Liang, and Fung 2020). Despite a relatively large phylogenetic distance and divergent S1 structures, NL63 and SARS-CoV-1 and SARS-CoV-2 bind to the same host receptor using the same virus-binding motifs (VBMs) (Li 2016). This VBM is located in the C-terminal domain of S1 (S1-CTD), which fits within the trend of S1-CTD receptor-binding in CoVs that bind protein receptors (Hofmann et al. 2006; Li 2016). This is opposed to the trend amongst CoVs that bind sugar receptors, where receptor-binding is located within the S1 N-terminal domain (S1-NTD) (Li 2016). This localization roughly aligns with our observations that the majority of the repeatedly-mutated sites occur toward the C-terminal end of 229E S1 and the N-terminal end of OC43 S1 (Fig. 2).

Here, we have provided support that at least 2 of the 4 seasonal HCoVs evolve adaptively in the region of spike that is known to interact with the humoral immune system. These two viruses span both genera of HCoVs, though due to the complexity of HCoV receptor-binding and pathology mentioned above, it is not clear whether or not this suggests that other HCoVs, such as SARS-CoV-2, will also evolve adaptively in S1. This is important because, at the time of writing of this manuscript, many SARS-CoV-2 vaccines are in production and most of these exclusively include spike (Krammer 2020). If SARS-CoV-2 evolves adaptively in S1 as the closely-related HCoV OC43 does, it is possible that the SARS-CoV-2 vaccine would need to be frequently reformulated to match the circulating strains, as is done for seasonal influenza vaccines.

## Materials and methods

All data, source code and analyses can be found at

<https://github.com/blab/seasonal-cov-adaptive-evolution>. All phylogenetic trees constructed and analyzed in this manuscript can be viewed interactively at

<https://nextstrain.org/community/blab/seasonal-cov-adaptive-evolution>.

## Sequence data

All viral sequences are publicly accessible and were downloaded from ViPR ([www.viprbrc.org](http://www.viprbrc.org)) under the “Coronaviridae” with host “human” (Pickett et al. 2012). Sequences labeled as “OC43”, “229E”, “HKU1” and “NL63” were pulled out of the downloaded FASTA file into 4 separate data files. Additionally, a phylogeny of all downloaded human coronaviruses was made and unlabeled isolates that clustered within clades formed by labeled OC43, 229E, HKU1 or NL63 isolates were marked as belonging to that HCoV type and added to our data files. Code for these data-parsing steps is located in

`data-wrangling/postdownload_formatting_for_rerun.ipynb`.

## Phylogenetic inference

For each of the 4 HCoV datasets, full-length sequences were aligned to a reference genome using the `augur align` command (Hadfield et al. 2018) and MAFFT (Katoh et al. 2002). Individual gene sequences were then extracted from these alignments if sequencing covered 50% or more of the gene using the code in

`data-wrangling/postdownload_formatting_for_rerun.ipynb`. Sequence files for each gene are located in the `data/` directory within each HCoV parent directory (ex: `oc43/data/oc43_spike.fasta`). A Snakemake file (Köster and Rahmann 2012) within each HCoV directory was then used to align each gene to a reference strain and a time-resolved phylogeny was built with IG-Tree (Nguyen et al. 2015) and TimeTree (Sagulenko, Puller, and Neher 2018). Phylogenies were viewed to identify the distribution of genotypes throughout the tree, different lineages, and signals of recombination using the `nextstrain view` command (Hadfield et al. 2018). The clock rate of the phylogeny based on spike sequences for each isolate (as shown in Fig. 1 and Fig. 1 Supplement 2) was 0.0005 for OC43, 0.0006 for 229E, 0.0007 for NL63, and 0.0062 for HKU1. All NL63 and HKU1 trees were rooted on an outgroup sequence. For NL63, the outgroup was 229e/AF304460/229e\_ref/Germany/2000 and for HKU1 the outgroup was mhv/NC\_048217\_1/mhv/2006. Clock rates for the phylogenies built on each individual gene can be found within the `results/` directory within each HCoV parent directory (ex: `oc43/results/branch_lengths_oc43_spike.json`).

## Mutation counting

Amino acid substitutions at each position in spike were tallied from the phylogeny using code in `antigenic_evolution/site_mutation_rank.ipynb`.

## Divergence analysis



For each HCoV lineage and each gene, synonymous and nonsynonymous divergence was calculated at all timepoints as the average Hamming distance between each sequenced isolate and the consensus sequence at the first timepoint (founder sequence). The total number of observed differences between the isolate and founder nucleotide sequences that result in nonsynonymous (or synonymous) substitutions is divided by the number of possible nucleotide mutations that result in nonsynonymous (or synonymous) substitutions, weighted by kappa, to yield an estimate of divergence. Kappa is the ratio of rates of transitions:transversions, and was calculated by averaging values from spike and RdRp trees built by BEAST 2.6.3 (Bouckaert et al. 2019) using the HKY+gamma4 model with 2 partitions and “coalescent constant population”. All BEAST results are found in `.log` files in gene- and HCoV-specific subdirectories within `beast/`. Divergence is calculated from nucleotide alignments. Sliding 3-year windows were used and only timepoints that contained at least 2 sequences were considered. The concept for this analysis is from (Zanini et al. 2015) and code for our adaptation is in `antigenic_evolution/divergence_weighted.ipynb`.

### Implementation of the Bhatt method

The rate of adaptive evolution was computed using an adaptation of the Bhatt method (Bhatt, Holmes, and Pybus 2011; Bhatt, Katzourakis, and Pybus 2010). Briefly, this method defines a class of neutrally-evolving nucleotide sites, then identifies other classes with higher rates of nonsynonymous nucleotide fixations and high-frequency polymorphisms. This method compares nucleotide sequences at each timepoint (the ingroup) to the consensus nucleotide sequence at the first time point (the outgroup) and yields an estimate of the number of adaptive substitutions within a given genomic region at each of these timepoints. Eight estimators (silent fixed, replacement fixed, silent high frequency, replacement high frequency, silent mid-frequency, replacement mid-frequency, silent low frequency and replacement low-frequency) are then calculated by the site-counting method (Bhatt, Katzourakis, and Pybus 2010). In the site-counting method, each estimator is the product of the fixation or polymorphism score times the silent or replacement score, summed for each site in that frequency class. Fixation and polymorphism scores depend on the number of different nucleotides observed at the site and whether the outgroup base is present in the ingroup. Selectively neutral sites are assumed to contain the classes of silent polymorphisms and replacement polymorphisms occurring at a frequency between 0.15 and 0.75. A class of nonneutral, adaptive sites is then identified as having an excess of replacement fixations or polymorphisms (Bhatt, Holmes, and Pybus 2011). Sliding 3-year windows were used and only timepoints that contained at least 2 sequences were considered. For each lineage and gene, 100 bootstrap alignments and ancestral sequences were generated and run through the Bhatt method to assess the statistical uncertainty of our estimates of rates of adaptation (Bhatt, Holmes, and Pybus 2011). The rate of adaptation (per codon per year) shown in Fig. 5 is calculated by linear regression of the time series values of adaptive substitutions per codon (Fig. 4). Our code for implementing the Bhatt method is at `antigenic_evolution/bhatt_bootstrapping.ipynb`.

### Estimation of rates of adaptation of H3N2 and measles

Influenza H3N2 and measles sequencing data was downloaded from <https://github.com/nextstrain/seasonal-flu> and <https://github.com/nextstrain/measles>, respectively. The rates of adaptation of different genes was calculated using our implementation of the Bhatt method described above. The receptor-binding domain used for H3N2 was HA1, for measles was the H protein, and for the HCoV was S1. The membrane fusion protein used for H3N2 was HA2, for measles was the F protein, and for the HCoV was S2. The polymerase for H3N2 was PB1, for measles was the P protein, and for the HCoV was RbRd (nsp12). Our code for this analysis is at `antigenic_evolution/measles_h3n2_bhatt.ipynb`.

### Simulation of evolving OC43 sequences

The evolution of OC43 lineage A Spike and RdRp genes was simulated using SANTA-SIM (Jariani et al. 2019). The OC43 lineage A root sequence was used as a starting point and the simulation was run for 500 generations and 10 simulated sequences were sampled every 50 generations. Evolution was simulated in the absence of recombination and with moderate and high levels of recombination during replication. Under each of these recombination paradigms, we simulated evolution in the absence of positive selection within spike and with moderate and high levels of positive selection. Positive selection was simulated at a subset of Spike sites proportional to the number of epitope sites in H3N2 HA (Luksza and Lässig 2014). All simulations were run with a nucleotide mutation rate of  $1 \times 10^{-4}$  (Vijgen et al. 2005). Config files, results and source code for these simulations can be at `santa-sim_oc43a/`.

### Estimation of TMRCA

Mean TMRCA values were estimated for each gene and each HCoV using PACT (Bedford, Cobey, and Pascual 2011). The PACT config files and results for each run are in the directory `antigenic_evolution/pact/`. The TMRCA estimations and subsequent analyses are executed by code in `antigenic_evolution/tmrca_pact.ipynb`.

### Acknowledgments

We thank Jesse Bloom and members of the Bedford lab for useful feedback. KEK was supported by the National Science Foundation Graduate Research Fellowship Program under Grant No. DGE-1762114. TB is a Pew Biomedical Scholar and is supported by NIH R35 GM119774-01.

### References

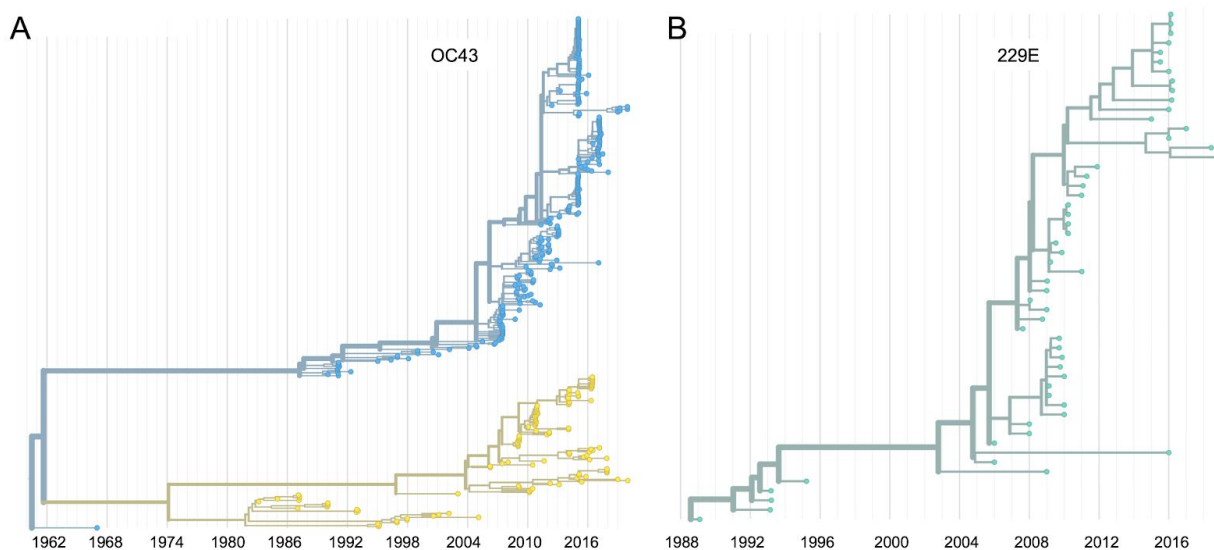
- Bedford, Trevor, Sarah Cobey, and Mercedes Pascual. 2011. "Strength and Tempo of Selection Revealed in Viral Gene Genealogies."
- Bhatt, Samir, Edward C. Holmes, and Oliver G. Pybus. 2011. "The Genomic Rate of Molecular Adaptation of the Human Influenza A Virus." *Molecular Biology and Evolution* 28 (9): 2443–51.
- Bhatt, Samir, Aris Katzourakis, and Oliver G. Pybus. 2010. "Detecting Natural Selection in RNA Virus Populations Using Sequence Summary Statistics." *Infection, Genetics and Evolution: Journal of Molecular Epidemiology and Evolutionary Genetics in Infectious Diseases* 10 (3): 421–30.

- Bouckaert, Remco, Timothy G. Vaughan, Joëlle Barido-Sottani, Sebastián Duchêne, Mathieu Fourment, Alexandra Gavryushkina, Joseph Heled, et al. 2019. "BEAST 2.5: An Advanced Software Platform for Bayesian Evolutionary Analysis." *PLoS Computational Biology* 15 (4): e1006650.
- Chibo, Doris, and Chris Birch. 2006. "Analysis of Human Coronavirus 229E Spike and Nucleoprotein Genes Demonstrates Genetic Drift between Chronologically Distinct Strains." *The Journal of General Virology* 87 (Pt 5): 1203–8.
- Drake, J. W. 1993. "Rates of Spontaneous Mutation among RNA Viruses." *Proceedings of the National Academy of Sciences of the United States of America* 90 (9): 4171–75.
- Edridge, Arthur W. D., Joanna Kaczorowska, Alexis C. R. Hoste, Margreet Bakker, Michelle Klein, Katherine Loens, Maarten F. Jebbink, et al. 2020. "Seasonal Coronavirus Protective Immunity Is Short-Lasting." *Nature Medicine*, September. <https://doi.org/10.1038/s41591-020-1083-1>.
- Fulton, Benjamin O., David Sachs, Shannon M. Beaty, Sohui T. Won, Benhur Lee, Peter Palese, and Nicholas S. Heaton. 2015. "Mutational Analysis of Measles Virus Suggests Constraints on Antigenic Variation of the Glycoproteins." *Cell Reports* 11 (9): 1331–38.
- Hadfield, James, Colin Megill, Sidney M. Bell, John Huddleston, Barney Potter, Charlton Callender, Pavel Sagulenko, Trevor Bedford, and Richard A. Neher. 2018. "Nextstrain: Real-Time Tracking of Pathogen Evolution." *Bioinformatics* 34 (23): 4121–23.
- Hamre, D., and M. Beem. 1972. "Virologic Studies of Acute Respiratory Disease in Young Adults. V. Coronavirus 229E Infections during Six Years of Surveillance." *American Journal of Epidemiology* 96 (2): 94–106.
- Heikkinen, Terho, and Asko Järvinen. 2003. "The Common Cold." *The Lancet* 361 (9351): 51–59.
- Hofmann, Heike, Graham Simmons, Andrew J. Rennekamp, Chawaree Chaipan, Thomas Gramberg, Elke Heck, Martina Geier, et al. 2006. "Highly Conserved Regions within the Spike Proteins of Human Coronaviruses 229E and NL63 Determine Recognition of Their Respective Cellular Receptors." *Journal of Virology* 80 (17): 8639–52.
- Hon, Chung-Chau, Tsan-Yuk Lam, Zheng-Li Shi, Alexei J. Drummond, Chi-Wai Yip, Fanya Zeng, Pui-Yi Lam, and Frederick Chi-Ching Leung. 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 (4): 1819–26.
- Hulswit, Ruben J. G., Yifei Lang, Mark J. G. Bakkers, Wentao Li, Zeshi Li, Arie Schouten, Bram Ophorst, et al. 2019. "Human Coronaviruses OC43 and HKU1 Bind to 9-O-Acetylated Sialic Acids via a Conserved Receptor-Binding Site in Spike Protein Domain A." *Proceedings of the National Academy of Sciences of the United States of America* 116 (7): 2681–90.
- Jarhani, Abbas, Christopher Warth, Koen Deforche, Pieter Libin, Alexei J. Drummond, Andrew Rambaut, Frederick A. Matsen Iv, and Kristof Theys. 2019. "SANTA-SIM: Simulating Viral Sequence Evolution Dynamics under Selection and Recombination." *Virus Evolution* 5 (1): vez003.
- Katoh, Kazutaka, Kazuharu Misawa, Kei-Ichi Kuma, and Takashi Miyata. 2002. "MAFFT: A Novel Method for Rapid Multiple Sequence Alignment Based on Fast Fourier Transform." *Nucleic Acids Research* 30 (14): 3059–66.
- Komabayashi, Kenichi, Yohei Matoba, Shizuka Tanaka, Junji Seto, Yoko Aoki, Tatsuya Ikeda, Yoshitaka Shimotai, Yoko Matsuzaki, Tsutomu Itagaki, and Katsumi Mizuta. 2020. "Longitudinal Epidemiology of Human Coronavirus OC43 in Yamagata, Japan, 2010–2017: Two Groups Based on Spike Gene Appear One after Another." *Journal of Medical Virology*

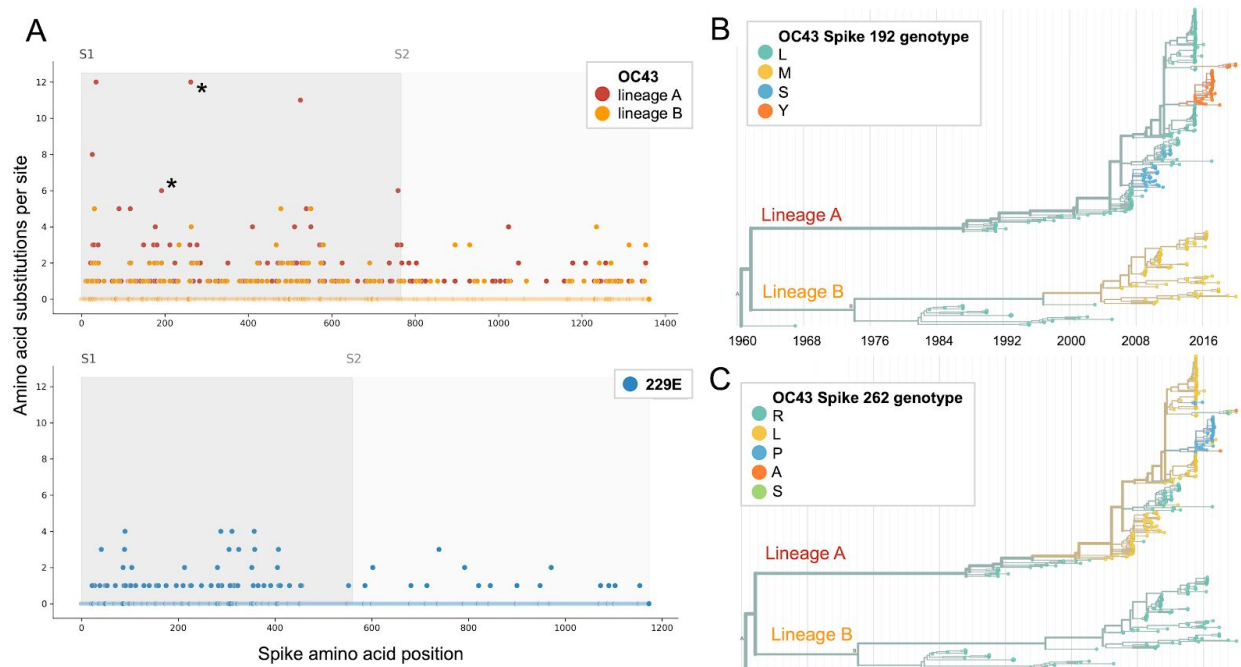
7 (August): 825.

- Kosakovsky Pond, Sergei L., David Posada, Michael B. Gravenor, Christopher H. Woelk, and Simon D. W. Frost. 2006. "Automated Phylogenetic Detection of Recombination Using a Genetic Algorithm." *Molecular Biology and Evolution* 23 (10): 1891–1901.
- Köster, Johannes, and Sven Rahmann. 2012. "Snakemake—a Scalable Bioinformatics Workflow Engine." *Bioinformatics* 28 (19): 2520–22.
- Krammer, Florian. 2020. "SARS-CoV-2 Vaccines in Development." *Nature*, September. <https://doi.org/10.1038/s41586-020-2798-3>.
- Kucharski, Adam J., Justin Lessler, Jonathan M. Read, Huachen Zhu, Chao Qiang Jiang, Yi Guan, Derek A. T. Cummings, and Steven Riley. 2015. "Estimating the Life Course of Influenza A(H3N2) Antibody Responses from Cross-Sectional Data." *PLOS Biology*. <https://doi.org/10.1371/journal.pbio.1002082>.
- Lau, Susanna K. P., Paul Lee, Alan K. L. Tsang, Cyril C. Y. Yip, Herman Tse, Rodney A. Lee, Lok-Yee So, et al. 2011. "Molecular Epidemiology of Human Coronavirus OC43 Reveals Evolution of Different Genotypes over Time and Recent Emergence of a Novel Genotype due to Natural Recombination." *Journal of Virology* 85 (21): 11325–37.
- Li, Fang. 2016. "Structure, Function, and Evolution of Coronavirus Spike Proteins." *Annual Review of Virology* 3 (1): 237–61.
- Liu, Ding X., Jia Q. Liang, and To S. Fung. 2020. "Human Coronavirus-229E, -OC43, -NL63, and -HKU1." *Reference Module in Life Sciences*. <https://doi.org/10.1016/b978-0-12-809633-8.21501-x>.
- Luksza, Marta, and Michael Lässig. 2014. "A Predictive Fitness Model for Influenza." *Nature* 507 (7490): 57–61.
- McIntosh, Kenneth. 1974. "Coronaviruses: A Comparative Review." In *Current Topics in Microbiology and Immunology / Ergebnisse Der Mikrobiologie Und Immunitätsforschung*, 85–129. Springer Berlin Heidelberg.
- Monto, Arnold S., and Sook K. Lim. 1974. "The Tecumseh Study of Respiratory Illness. VI. Frequency of and Relationship between Outbreaks of Coronavims Infection." *The Journal of Infectious Diseases* 129 (3): 271–76.
- Nguyen, Lam-Tung, Heiko A. Schmidt, Arndt von Haeseler, and Bui Quang Minh. 2015. "IQ-TREE: A Fast and Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies." *Molecular Biology and Evolution* 32 (1): 268–74.
- Pasternak, Alexander O., Willy J. M. Spaan, and Eric J. Snijder. 2006. "Nidovirus Transcription: How to Make Sense...?" *The Journal of General Virology* 87 (6): 1403–21.
- Pickett, Brett E., Eva L. Sadat, Yun Zhang, Jyothi M. Noronha, R. Burke Squires, Victoria Hunt, Mengya Liu, et al. 2012. "ViPR: An Open Bioinformatics Database and Analysis Resource for Virology Research." *Nucleic Acids Research* 40 (Database issue): D593–98.
- Rambaut, Andrew, Oliver G. Pybus, Martha I. Nelson, Cecile Viboud, Jeffery K. Taubenberger, and Edward C. Holmes. 2008. "The Genomic and Epidemiological Dynamics of Human Influenza A Virus." *Nature* 453 (7195): 615–19.
- Reed, Sylvia E. 1984. "The Behaviour of Recent Isolates of Human Respiratory Coronavirus in Vitro and in Volunteers: Evidence of Heterogeneity among 229E-Related Strains." *Journal of Medical Virology* 13 (2): 179–92.
- Ren, Lili, Yue Zhang, Jianguo Li, Yan Xiao, Jing Zhang, Ying Wang, Lan Chen, Gláucia Paranhos-Baccalà, and Jianwei Wang. 2015. "Genetic Drift of Human Coronavirus OC43 Spike Gene during Adaptive Evolution." *Scientific Reports* 5 (June): 11451.
- Sagulenko, Pavel, Vadim Puller, and Richard A. Neher. 2018. "TreeTime: Maximum-Likelihood Phylodynamic Analysis." *Virus Evolution* 4 (1): vex042.

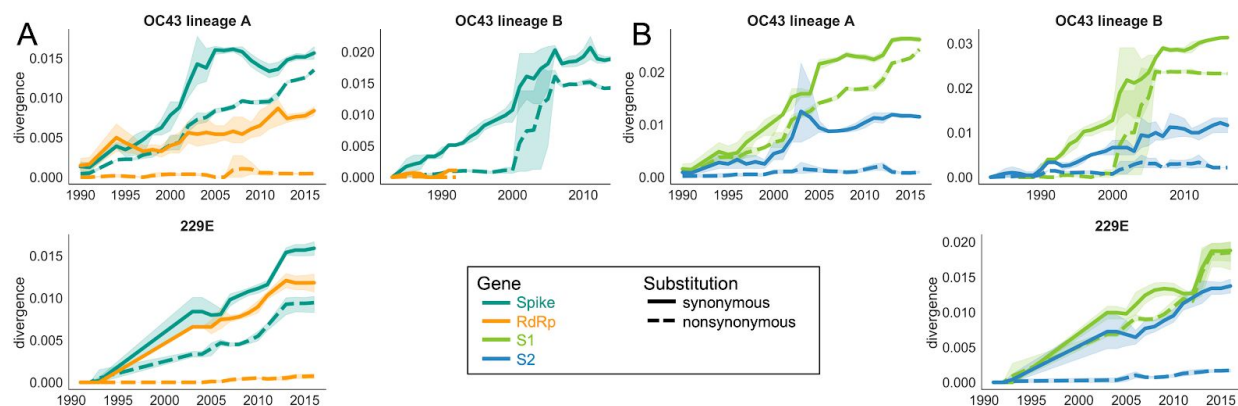
- Smith, Derek J., Alan S. Lapedes, Jan C. de Jong, Theo M. Bestebroer, Guus F. Rimmelzwaan, Albert D. M. E. Osterhaus, and Ron A. M. Fouchier. 2004. "Mapping the Antigenic and Genetic Evolution of Influenza Virus." *Science* 305 (5682): 371–76.
- Vijgen, Leen, Philippe Lemey, Els Keyaerts, and Marc Van Ranst. 2005. "Genetic Variability of Human Respiratory Coronavirus OC43." *Journal of Virology*.
- Volz, Erik M., Katia Koelle, and Trevor Bedford. 2013. "Viral Phylodynamics." *PLoS Computational Biology* 9 (3): e1002947.
- Woo, Patrick C. Y., Susanna K. P. Lau, Yi Huang, and Kwok-Yung Yuen. 2009. "Coronavirus Diversity, Phylogeny and Interspecies Jumping." *Experimental Biology and Medicine* 234 (10): 1117–27.
- Yang, Z. 2000. "Maximum Likelihood Estimation on Large Phylogenies and Analysis of Adaptive Evolution in Human Influenza Virus A." *Journal of Molecular Evolution* 51 (5): 423–32.
- Zanini, Fabio, Johanna Brodin, Lina Thebo, Christa Lanz, Göran Bratt, Jan Albert, and Richard A. Neher. 2015. "Population Genomics of Inpatient HIV-1 Evolution." *eLife* 4 (December). <https://doi.org/10.7554/eLife.11282>.
- Zhang, Yue, Jianguo Li, Yan Xiao, Jing Zhang, Ying Wang, Lan Chen, Gláucia Paranhos-Baccalà, Lili Ren, and Jianwei Wang. 2015. "Genotype Shift in Human Coronavirus OC43 and Emergence of a Novel Genotype by Natural Recombination." *Journal of Infection*. <https://doi.org/10.1016/j.jinf.2014.12.005>.
- Zhu, Yun, Changchong Li, Li Chen, Baoping Xu, Yunlian Zhou, Ling Cao, Yunxiao Shang, et al. 2018. "A Novel Human Coronavirus OC43 Genotype Detected in Mainland China." *Emerging Microbes & Infections* 7 (1): 173.



**Figure 1. Phylogenetic trees for spike gene of seasonal HCoVs OC43 and 229E.** Phylogenies built from A: OC43 spike sequences from 389 isolates over 53 years, and B: 229E spike sequences from 54 isolates over 31 years. HCoVs that bifurcate immediately after the root are split into blue and yellow lineages. 229E and contains just one lineage (teal). For the analyses in this paper, the evolution of each gene (or genomic region) is considered separately, so phylogenies are built for each viral gene and those phylogenies are used to split isolates into lineages for each gene. These are temporally resolved phylogenies with year shown on the x-axis. The clock rate estimate is  $5 \times 10^{-4}$  for OC43 and  $6 \times 10^{-4}$  for 229E.

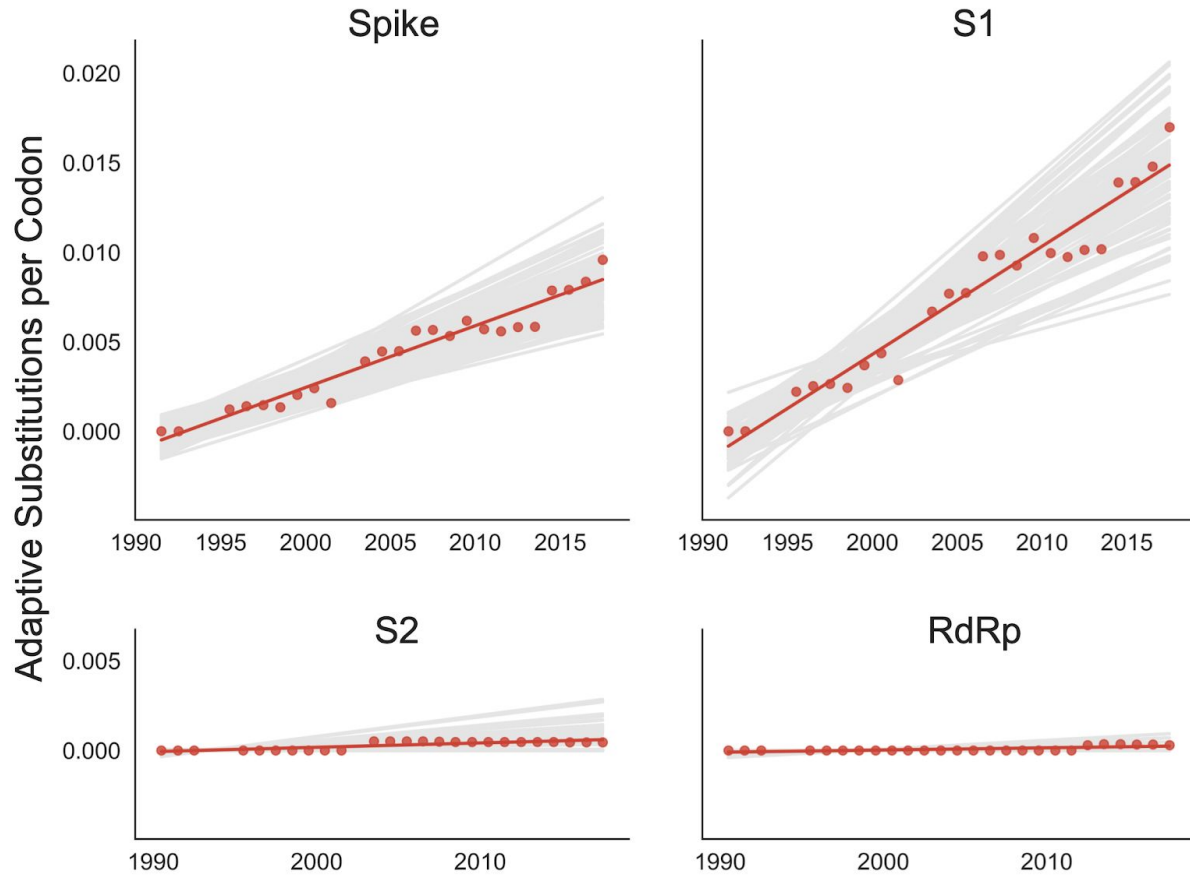


**Figure 2. More sites mutate repeatedly within spike S1 versus S2.** A: Number of mutations observed at each position in the spike gene. S1 (darker gray) and S2 (light gray) are indicated by shading and the average number of mutations per site is indicated by a dot and color-coded by HCoV lineage. Asterisks indicate positions 192 and 262, which mutate repeatedly throughout the OC43 lineage A phylogeny. The OC43 phylogeny built from spike sequences and color-coded by genotype at position 192 and 262 is shown in B) and C), respectively.

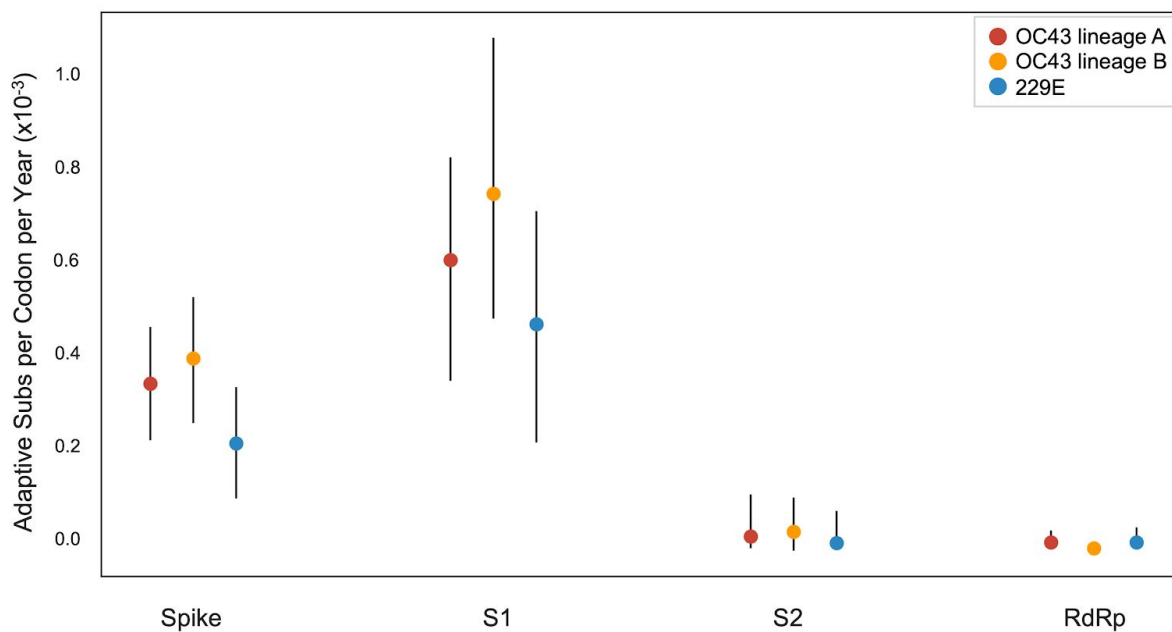


**Figure 3. Nonsynonymous divergence is higher in OC43 and 229E Spike S1 versus S2 or RdRp.** A: Nonsynonymous (dashed lines) and synonymous divergence (solid lines) of the spike (teal) and RdRp (orange) genes of all 229E and OC43 lineages over time. Divergence is the average Hamming distance from the ancestral sequence, computed in sliding 3-year windows which contain at least 2 sequenced isolates. Shaded region shows 95% confidence intervals. B: Nonsynonymous and synonymous divergence within the S1 (light green) and S2 (blue) domains of spike. Year is shown on the x-axis. Note that x- and y-axis scales are not shared between plots.

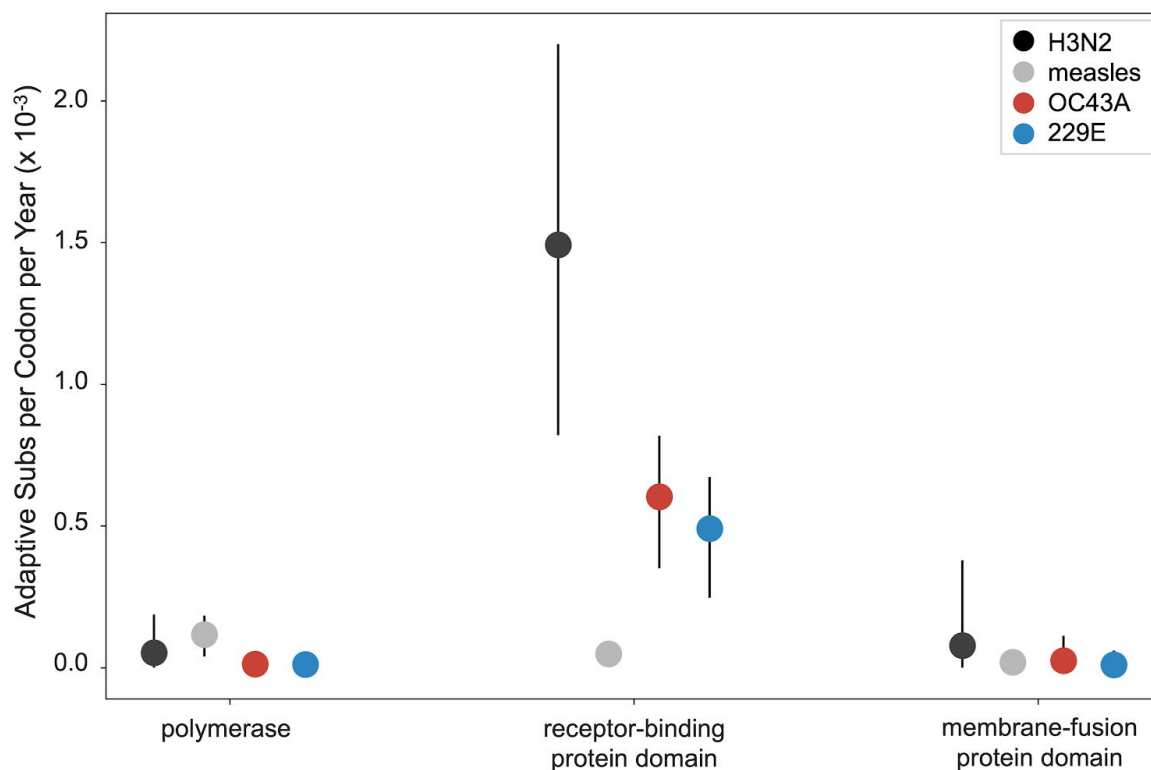




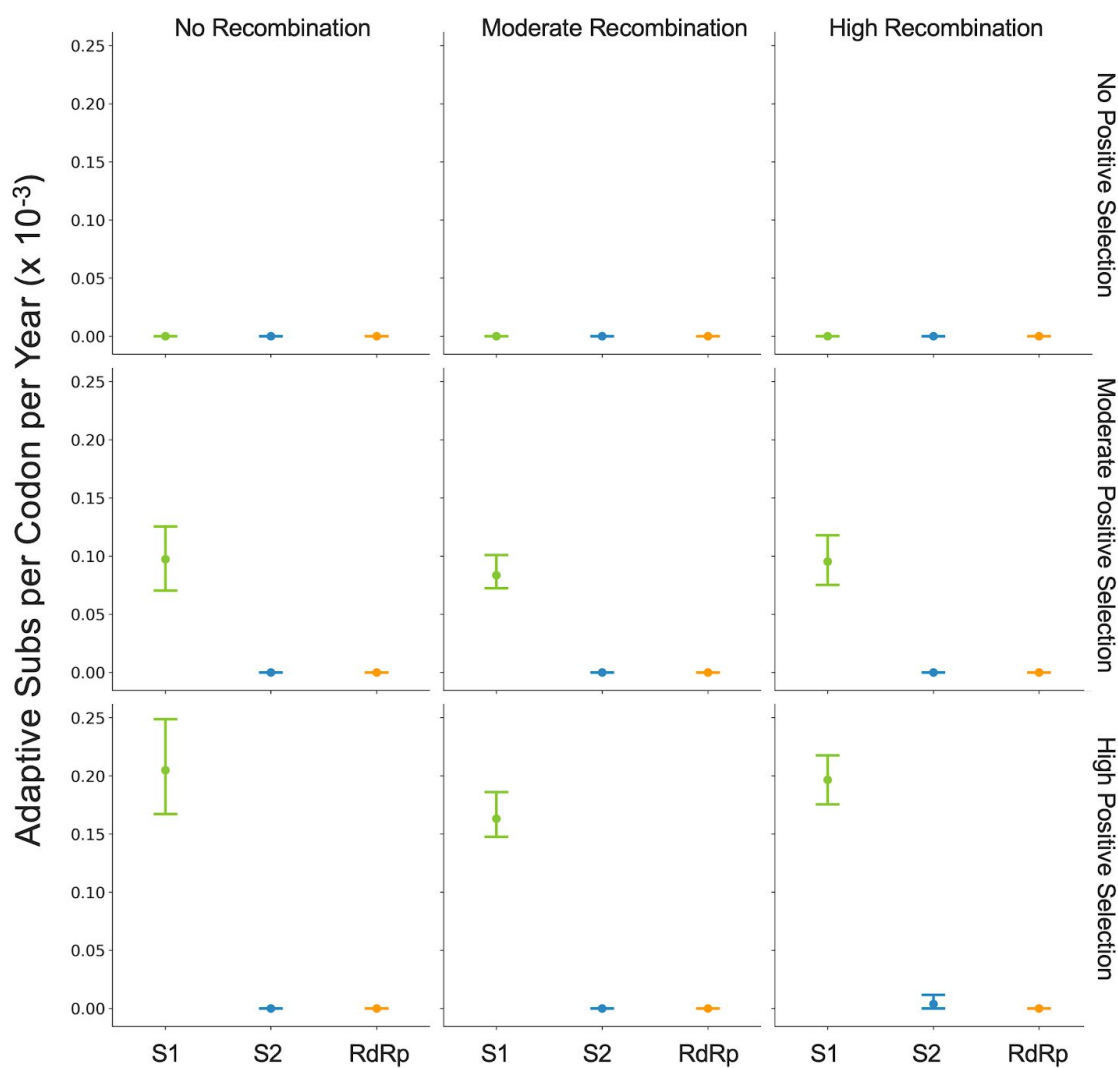
**Figure 4. Adaptive substitutions accumulate over time in OC43 lineage A spike S1.** Adaptive substitutions per codon within OC43 lineage A spike, S1, S2 and RdRp as calculated by our implementation of the Bhatt method. Adaptive substitutions are computed in sliding 3-year windows, and only for timepoints that contain 3 or more sequenced isolates. Red dots display estimated values calculated from the empirical data and red lines show linear regression fit to these points. Grey lines show the distribution of regressions fit to the computed number of adaptive substitutions from 100 bootstrapped datasets. Year is shown on the x-axis.



**Figure 5. The rate of adaptive substitution is highest in spike S1.** Adaptive substitutions per codon per year as calculated by our implementation of the Bhatt method. Rates are calculated within Spike, S1, S2 and RdRp for 229E and OC43 lineages. Error bars show 95% bootstrap percentiles from 100 bootstrapped datasets.



**Figure 6. OC43 and 229E spike S1 accumulates adaptive substitutions faster than measles but slower than influenza H3N2.** Comparison of adaptive substitutions per codon per year between influenza H3N2 (black), measles (gray), OC43 lineage A (red), and 229E (orange). The polymerase, receptor binding domain and membrane fusion domain for H3N2 are PB1, HA1 and HA2. For both HCoVs, they are RdRp, S1 and S2, respectively. For measles, the polymerase is the P gene, the receptor-binding protein is the H gene and the fusion protein is the F gene. Error bars show 95% bootstrap percentiles from 100 bootstrapped datasets.



**Figure 7. Detection of positive selection is not biased by recombination.** OC43 lineage A sequences were simulated with varying levels of recombination and positive selection. The Bhatt method was used to calculate the rate of adaptive substitutions per codon per year for S1 (light green), S2 (blue) and RdRp (orange). The mean and 95% confidence interval of 5 independent simulations is plotted.

	<b>Spike</b>	<b>S1</b>	<b>S2</b>	<b>RdRp</b>
<b>OC43A</b>	4.67	3.45	13.05	17.39
<b>229E</b>	4.19	2.23	5.08	4.86

**Table 1. Mean TMRCA is lower in S1 than RdRp or S2.** Average TMRCA values (in years) for OC43 lineage A and 229E.