



Human pathogenic RNA viruses establish noncompeting lineages by occupying independent niches

Pascal Mutz^{a,1}, Nash D. Rochman^{a,1,2}, Yuri I. Wolf^a, Guilhem Faure^b, Feng Zhang^{b,c,d,e,f,2}, and Eugene V. Koonin^{a,2}

Contributed by Feng Zhang; received November 23, 2021; accepted April 1, 2022; reviewed by Sergei Maslov and Claus Wilke

Many pathogenic viruses are endemic among human populations and can cause a broad variety of diseases, some potentially leading to devastating pandemics. How virus populations maintain diversity and what selective pressures drive population turnover is not thoroughly understood. We conducted a large-scale phylodynamic analysis of 27 human pathogenic RNA viruses spanning diverse life history traits, in search of unifying trends that shape virus evolution. For most virus species, we identify multiple, cocirculating lineages with low turnover rates. These lineages appear to be largely noncompeting and likely occupy semi-independent epidemiological niches that are not regionally or seasonally defined. Typically, intralinear mutational signatures are similar to interlineage signatures. The principal exception are members of the family *Picornaviridae*, for which mutations in capsid protein genes are primarily lineage defining. Interlineage turnover is slower than expected under a neutral model, whereas intralinear turnover is faster than the neutral expectation, further supporting the existence of independent niches. The persistence of virus lineages appears to stem from limited outbreaks within small communities, so that only a small fraction of the global susceptible population is infected at any time. As disparate communities become increasingly connected through globalization, interaction and competition between lineages might increase as well, which could result in changing selective pressures and increased diversification and/or pathogenicity. Thus, in addition to zoonotic events, ongoing surveillance of familiar, endemic viruses appears to merit global attention with respect to the prevention or mitigation of future pandemics.

effective population size | pandemic | globalization | human RNA viruses | influenza

Viruses, ubiquitous across the tree of life, occupy an astounding diversity of ecological niches (1–3). Viral niches are primarily defined by the behavior and immunity of the respective hosts and are often the subject of deep but narrow investigation (4, 5). In this work, we sought to uncover common trends at relatively short evolutionary distances by studying the microevolution of human pathogenic RNA viruses. The devastating COVID-19 pandemic has made it abundantly clear that understanding these microevolutionary features is of vital importance not only to forward our understanding of virology, in general, but to inform appropriate public health measures during a pandemic (6, 7).

Viral populations explore their viable sequence space defined by both intrinsic constraints and those imposed by host behavior (8) through the accumulation of mutations, potentially leading to diversification (9). A single host species can offer multiple independent niches that are explored by distinct virus subpopulations. Niches can be formed and maintained through regional or seasonal separation. Regional separation of subpopulations has been demonstrated, for example, for yellow fever virus (YFV) (10, 11). At sufficiently long evolutionary distances, niches can be defined by immunological differences, which enable a viral subpopulation to overcome immune cross-protection, allowing the same host to be infected by multiple subpopulations largely independent of prior infections. This phenomenon has been demonstrated for enteroviruses with many cocirculating serotypes (12). These immunological niches do not need to be spatially or temporally segregated.

Generally, niches are not necessarily static entities and can overlap or merge depending on dissemination rates, transmission modes, and other life history traits (2). When outbreaks are limited to small communities so that only a small fraction of the global susceptible population is infected at any time, niches can form that are not regionally or seasonally defined but are still maintained through a combination of spatial and temporal separation at a local scale. Thus, the maintenance of these viral niches is highly sensitive to changes in host behavior. The number and sequence diversity of such lineages depends on constraints intrinsic to viral biology as well as host behavior (13, 14). For example, there is a sharp contrast between the emergence of immunological niches among measles morbilli virus

Significance

Numerous pathogenic viruses are endemic in humans and cause a broad variety of diseases, but what is their potential for causing new pandemics? We show that most human pathogenic RNA viruses form multiple, cocirculating lineages with low turnover rates. These lineages appear to be largely noncompeting and occupy distinct epidemiological niches that are not regionally or seasonally defined, and their persistence appears to stem from limited outbreaks in small communities so that only a small fraction of the global susceptible population is infected at any time. However, due to globalization, interaction and competition between lineages might increase, potentially leading to increased diversification and pathogenicity. Thus, endemic viruses appear to merit global attention with respect to the prevention of future pandemics.

Author contributions: P.M., N.D.R., Y.I.W., G.F., F.Z., and E.V.K. analyzed data; and P.M., N.D.R., and E.V.K. wrote the paper.

Reviewers: S.M., University of Illinois Urbana-Champaign; and C.W., University of Texas at Austin.

The authors declare no competing interest.

Copyright © 2022 the Author(s). Published by PNAS. This open access article is distributed under Creative Commons Attribution License 4.0 (CC BY).

¹P.M. and N.D.R. contributed equally to this work.

²To whom correspondence may be addressed. Email: nash.rochman@nih.gov, zhang@broadinstitute.org, or koonin@ncbi.nlm.nih.gov.

This article contains supporting information online at <http://www.pnas.org/lookup/suppl/doi:10.1073/pnas.2121335119/-DCSupplemental>.

Published May 31, 2022.

(MMV) and influenza A virus (IAV) strain H3N2 (named H3N2 here). Through rapid antigenic drift and shift that involve a nonhuman host reservoir, IAV is able to overcome adaptive-immune protection, despite infecting a substantial fraction of the susceptible host population each year (15). As a consequence, H3N2 goes through phases of stasis, in which neutral evolution and purifying selection are dominant, parallel lineages are established, and population diversity grows. Once the pool of naïve hosts shrinks, the competition between lineages intensifies, resulting in a short phase of strong positive selection that favors one lineage to replace all others (16–18). In contrast, no such antigenic drift has been observed for MMV, and parallel MMV lineages do not replace each other but, rather, stably coexist (17, 19).

The persistence of multiple, coexisting viral lineages implies minimal interlineage competition. When such lineages are maintained through spatial or temporal separation, increased host–host or host–vector contact can result in the merger between and competition among multiple lineages. Climate change can support the spread and mixing of previously separated vectors, which could carry distinct viral lineages. With more vectors, the dissemination rate can rise, decreasing the number of susceptible hosts, and increasing competition globally (2). This can result in accelerated lineage turnover of human and agricultural pathogens, with the potential for substantial epidemiological and economic impact (20).

We sought to identify unifying trends of lineage emergence, persistence, and turnover among human pathogenic RNA viruses and to characterize the niches occupied by these lineages through phylogenetic analysis (21). Taking advantage of the substantial recent progress in virus genome sequencing (22), we constructed phylogenetic trees for the genomes of all monopartite human pathogenic RNA viruses for which extensive genome sequence information was available. These phylogenies were employed to assess the selection pressures affecting the evolution of these viruses through an analysis of the ratio of nonsynonymous to synonymous substitution rates (dN/dS), and to estimate the effective population size (N_e) and the census population size (N) for each. The viruses studied here are of clear epidemiological relevance, span a broad variety of life history traits (2, 23), and thus seem suitable to reveal unifying trends in the microevolution of RNA viruses. Our analysis of these viruses indicates that most form multiple, coexisting, noncompeting lineages which appear to occupy independent niches.

Results

Data Aggregation. Despite the substantial progress of the past several years (22), the available numbers of (nearly) complete genome sequences of human pathogenic RNA viruses differ widely among viral species, with few sequences available for several viruses. In the dataset for the present analysis, we included only those species for which 200 or more (nearly) complete genome sequences, with at least 50 isolated from a human host, were available in the National Center for Biotechnology Information (NCBI) virus database (24) or Global Initiative on Sharing All Influenza Data (GISAID) for severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) (25) (Fig. 1). These criteria excluded viruses which are widespread, for example, lyssa rhabdovirus and rubella virus, but for which few (nearly) complete genomes were available, as well as comparatively rare, even if highly pathogenic, viruses including some Ebola virus species (Zaire, but not Reston or Sudan, was included) and Marburg virus.

Only monopartite RNA viruses were considered, in order to exclude potential effects of segment reassortment and enable

the construction of a single, unambiguous phylogeny. This restriction excluded six species with many genomes available: three influenza viruses (A, B, and C), Reovirus, Lassa mammarenavirus, and Dabie bandavirus. We further omitted HIV, given its retrotranscribing replication strategy. Altogether, our dataset included 26 monopartite virus species. We added IAV H3N2 to this group (with a phylogeny constructed from hemagglutinin) as a thoroughly studied reference virus (14, 16, 17). The 27 viruses analyzed here cover a broad variety of viral lifestyles and ecological constraints and have been subject to varied countermeasures, including vaccination (*SI Appendix, Table S1*). This diversity enables the exploration of potential unifying trends of viral lineage turnover and niche formation.

Sequences were aligned using MAFFT (26), and, with the exception of SARS-CoV-2 and H3N2, for which the large number of sequences necessitated an iterative strategy, phylogenetic trees were constructed using IQ-TREE (27) (see *Brief Methods* and *SI Appendix, Extended Methods* for details). For most of viruses, the resulting trees included several large, clearly distinguishable clades (Fig. 1) that, in some cases, corresponded to known serotypes or genotypes (for example, Dengue virus, DENV, serotypes 1 to 4).

Low Rates of Lineage Turnover among Human RNA Viruses.

The major virus clades and the smaller lineages contained within them are subject to turnover whereby an older lineage goes extinct, being gradually replaced by individuals from a newer lineage. Trees with high turnover rates are often described as “cactus”- or “ladder”-like, and, in the limit of extreme turnover, as “caterpillar” trees, whereas those with low turnover are often described as “bush”-like, with ultrametric trees representing the limit of no turnover (14, 17). In an effort to explicitly measure lineage turnover [without relying solely on global measures such as coalescence rate (17), which is also estimated], we first sought to establish how many isolates, and distributed on the tree in what way, constitute a lineage. This information is important, in large part, because varying substitution rates across the tree complicate the estimation of global lineage turnover (9, 28). We defined lineages as monophyletic groups of sequences separated by periods (branches of the tree) with apparently different substitution rates and within which the sequencing date and the distance to the tree root are significantly correlated (see *Brief Methods* and *SI Appendix, Extended Methods* for details). Lineages cannot be defined in this way to encompass all sequences, and *SI Appendix, Fig. S1* shows the fraction of sequences included in correlated lineages for each virus. Arguably, significantly different substitution rates mark different selective environments and may reflect movement into distinct epidemiological niches. Because there are no apparent periods of different substitution rates within each lineage and, consequently, high-confidence date-constrained genealogical trees with a single substitution rate could be fit for each (see below), we denote these “genealogical lineages” (GLs).

Multiple GLs were identified in this manner for most viruses (*SI Appendix, Figs. S2–S7*), as illustrated in Fig. 2 *A* and *B* for the three lineages of enterovirus D (EVD). For human betacoronavirus 1 (BCoV1), Ebola Zaire (Ebola), MERS (Middle East respiratory syndrome-related coronavirus), H3N2, SARS-CoV-2, and Zika virus (ZIKV), the majority of the phylogeny comprised a single GL. Thus, the entire population of each of these viruses might occupy a single epidemiological niche at any point in time (which may be subject to rapid lineage turnover as is the case for H3N2). For Mumps rubulavirus (MRV) and YFV, the GLs identified were not large enough for subsequent

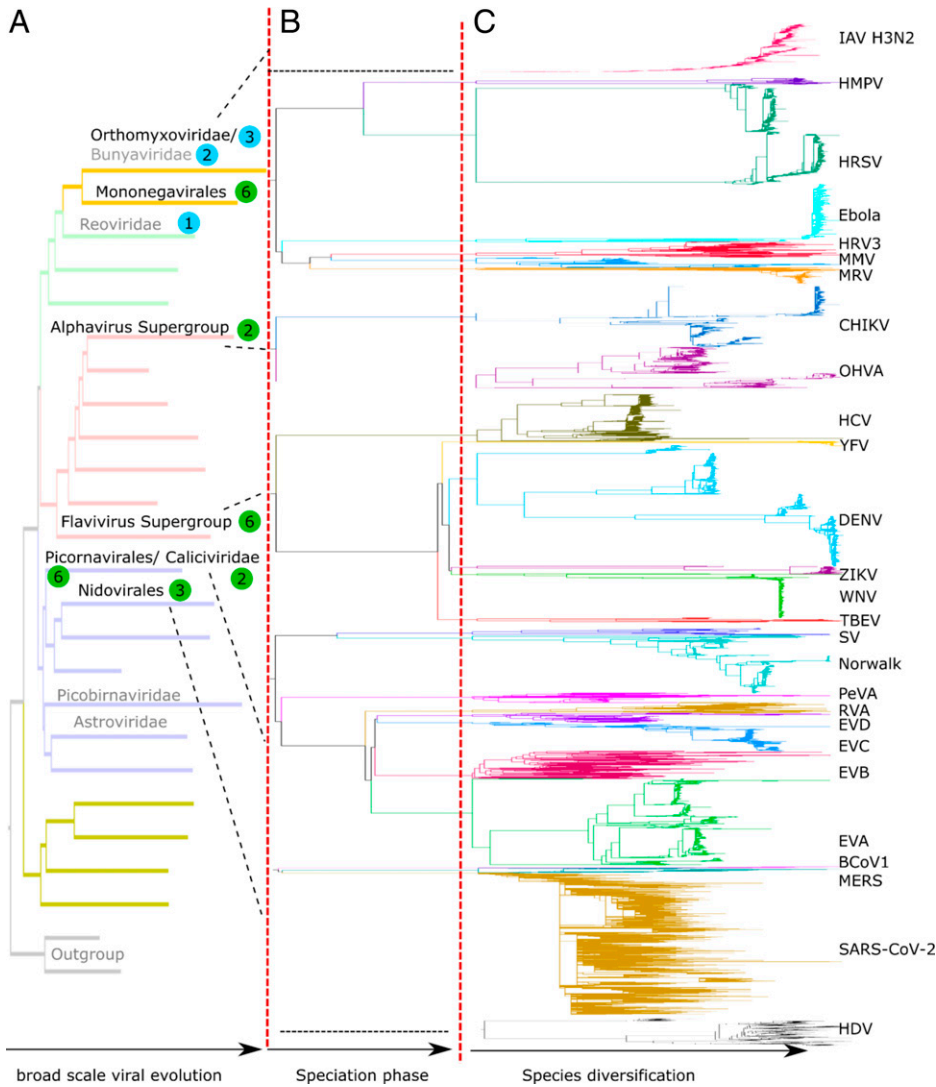


Fig. 1. Phylogenies of human pathogenic RNA viruses. Schematic depicting the origins and phylogenetic tree topologies of 27 human pathogenic RNA viruses. (A) Placement of each virus in the global phylogeny of RNA viruses (realm Riboviria). The tree topology is from ref. 76. Viral groups containing human pathogenic viruses are named in black if containing viruses analyzed in this work, and gray otherwise. The numbers of viral species, for which at least 200 nearly complete genome sequences were available, at least 50 of these isolated from humans, are shown in colored circles (green: monopartite viruses; blue: segmented viruses). (B) Speciation of viral families or orders. (C) Diversification within species. Trees for species are scaled to the same distance from the root to the most distal leaf and are grafted on the tree scaffold with arbitrary branch lengths for speciation but respecting topology.

analysis. When interpreting our observation of a single GL for Ebola, it should be noted that more than half of the Ebola isolates stem from the 2014–2016 outbreak in Sierra Leone, Liberia, and Guinea (29), the common assumption being that each individual Ebola outbreak stems from an individual zoonotic spillover event (30).

Having identified the virus GLs, we quantified the extent of lineage turnover using the Shannon entropy of the GL distribution over time, S_p , as well as traveling up the tree from root to leaf, S_d . These measures are independent of the timescale over which individual GLs persist and allow us to compare whether the phylogenetic structure of GLs (S_d) reflects their distribution over time (S_p). For this analysis, only sequences included within a GL were considered. First, sliding windows (indexed over j) containing the closest 5% of all isolates to the specified date, $w_{t,i}^j$, (from the date of the earliest isolated sequence to the latest) or distance to the tree root, $w_{t,d}^j$, were established, and the GL distribution within each window was obtained (Fig. 2C). Next, the probability that a sequence, x , within each window belongs to the i th GL was computed: $P_{t,d}^j(x \in GL_i | x \in w_{t,d}^j)$. The Shannon entropy of the GL distribution was then calculated using log base N equal to the number of GLs identified within the tree (and yielding a maximum value of one),

$$S_{t,d}^j = - \sum_{i=1}^N P_{t,d}^j(x \in GL_i | x \in w_{t,d}^j) \log_N \left(P_{t,d}^j(x \in GL_i | x \in w_{t,d}^j) \right).$$

Finally, the mean over all windows was computed for each tree: $S_{t,d}^j$ (Fig. 2D). A mean entropy near zero corresponds to a phylogeny composed of clades that rapidly displace one another (although effects of sampling bias cannot be excluded). A mean entropy near one corresponds to a phylogeny in which all clades are uniformly distributed at every time point. We observed $S_{t,j}^j > S_{t,d}^j$ in all but one case (hepatitis D virus [HDV]), indicating that entropy is greater than that expected from analysis of the tree structure with no known dates of isolation. This observation, coupled with the generally high mean entropy, suggests that most of the analyzed viruses evolve with low rates of lineage turnover.

To further quantify lineage turnover, we constructed date-constrained genealogical trees. As suggested above, GLs are separated by periods (branches of the tree) with apparently different substitution rates. These branches are often deep within the tree and are sparsely populated with leaves (if at all), making the assignment of a global model for substitution rates statistically dubious and highlighting the importance of rates inferred for individual GLs. Date-constrained trees were produced using a least-square distance approach based on the date of isolation

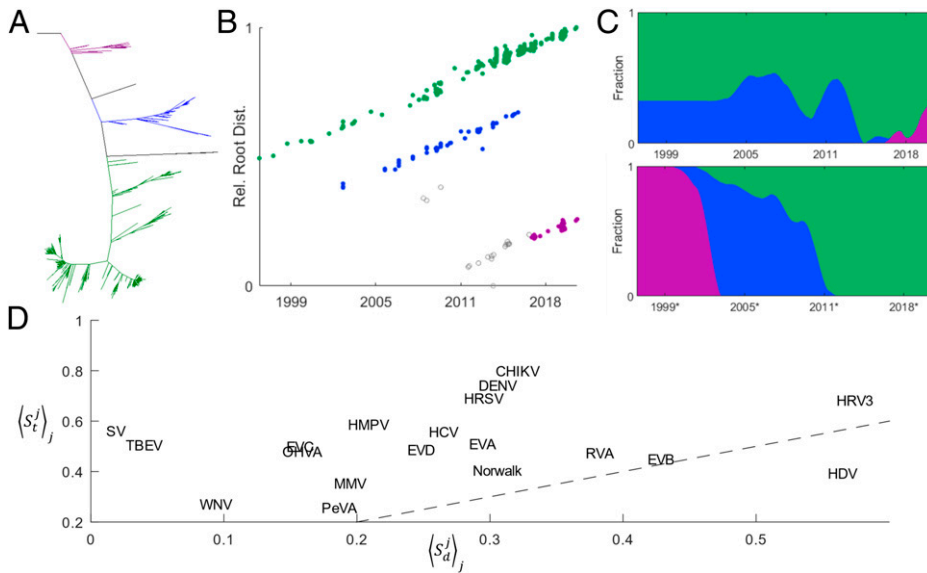


Fig. 2. Stable coexistence of lineages among human pathogenic RNA viruses. (A) EVD tree colored to represent the location of the three GLs. (B) Distance to the tree root vs. date of isolation for EVD. Distance is scaled by the maximum for any sequence within a correlated clade, and the x axis is bounded by the minimum/maximum date of isolation for any such sequence. (C) The fraction of isolates within each correlated clade (and excluding isolates that did not belong to any correlated clade) computed over a sliding window containing the nearest 5% of all isolates indexed by sequencing date (*Top*) and root distance (*Bottom*, where date* represents the date of isolation corresponding to each sequence in an alternative phylogeny where the date of isolation for each sequence exactly corresponds to the distance to the root for that sequence). (D) The mean Shannon entropy for the correlated-clade distribution respecting the sequencing date (y axis, S_{ij}^t) and root distance (x axis, S_{ij}^d). Dashed line displays $S_{ij}^t = S_{ij}^d$.

for each sequence (31). A mutation rate and the date of the last common ancestor (LCA) were estimated for all global trees and each GL individually (Fig. 3A). Samples without a known date of isolation introduce additional uncertainty into the calculation and can result in future-dated portions of the genealogical tree (see *Brief Methods* and *SI Appendix, Extended Methods* for details). GLs tend to accumulate over time, with few if any extinction events (Fig. 3B). Note that the apparent decline of parallel GLs for many viruses in recent years (~2016 to present) is most likely a sampling artifact, although, in principle, such decline could also point to a change in GL dynamics. These trends are further indicative of low lineage turnover and suggestive of minimal competition among GLs. Thus, each GL is likely to occupy a distinct niche, and we sought to identify the factors that shape and maintain these niches. *SI Appendix, Table S12* lists the number of GLs per virus species and other key parameters obtained during the analysis described below.

Most Virus Lineages Are Not Regionally or Seasonally Defined.

Perhaps the simplest explanation for the existence of distinct niches would be regional separation. To assess the role played by regionality, we examined whether isolates within a single GL clustered by region. The great circle map distance between pairs of isolates within each GL and between pairs of GLs was computed to retrieve the mean intra- and inter-GL distance, respectively (see *Brief Methods* and *SI Appendix, Extended Methods* for details). Given that GLs do not span the entire phylogeny for all viruses, were defined algorithmically without the incorporation of metadata beyond the date of isolation, and typically include a small number of isolates, we additionally examined the regionality of larger clades, usually defined by serotype or genotype (“manual lineages” [MLs]; see *Brief Methods* and *SI Appendix, Extended Methods* for details).

The ratio of the interlineage to intralinear map distances is expected to be greater than unity for regionally defined GLs or MLs, and near or below unity for those lineages that are not regionally defined (*SI Appendix, Fig. S8*). For most viruses analyzed, niches do not appear to be regionally defined, with a few notable exceptions (*SI Appendix, Fig. S8*). In particular, YFV is known to split into three regional lineages (East/Central Africa, West Africa, and South America), although the underlying mechanisms for this separation, especially between the African lineages, are not well understood (10, 11). Similarly, Chikungunya virus

(CHIKV) displays regionality, although, in this case, the separation seems to be incomplete (32). West Nile virus (WNV) lineage 1 can be found globally, whereas all other lineages are regional (33). However, evidence of the local coexistence of multiple WNV subtypes (34) indicates that additional WNV niches not linked to regionality might exist. HDV displays weak regionality that might be determined by its helper virus, hepatitis B virus (HBV), on which HDV depends for reproduction. The interplay between HDV and HBV genotypes is not yet well understood (35). Ebola outbreaks show a clear regional structure, which is due to de novo spillover events for each outbreak as well as successful containment measures (30). Some, but not all, GLs for Dengue virus (DENV), enterovirus A (EVA), and enterovirus B (EVB) may be regionally defined (*SI Appendix, Fig. S8*). Thus, while common, regionality does not appear to explain the existence of most apparent niches, although this does not imply an absence of spatial separation of localized outbreaks.

Similar to regionality, seasonality could potentially support niche formation. Although the temporal resolution of our analysis was limited by the amount of metadata available and the precision with which dates of isolation are specified, we found no evidence that seasonality plays a role in lineage maintenance within viral species. We observed no biannual or longer global periodicity of any GLs, but rather a continuous distribution of lineages through time (*SI Appendix, Figs. S2–S7*), although shorter temporal patterns are likely for respiratory viruses (36).

GLs could represent localized outbreaks (phases of enhanced virus spread) whereby a virus infects only a minute fraction of the global susceptible population at any given time. Under these conditions, even lineages which do not form distinct immunological niches and do not admit near-simultaneous infection can coexist within short distances of one another. Infection or vaccination leading to lifelong immunity, as observed, for example, in the case of MMV or MRV (14), can support the emergence of localized outbreaks. In these cases, naïve hosts are born and are not vaccinated, so that a local community of susceptible hosts emerges. Given sufficient evolutionary distance, lineages can become so diverse antigenically that they form different serotypes, which induce weak to no cross-immunity against each other and thus admit near-simultaneous infection. This pattern has been reported for some picornaviruses (12). In the case of zoonotic viruses, distinct lineages can originate when the same virus species is introduced from different animal reservoirs, which

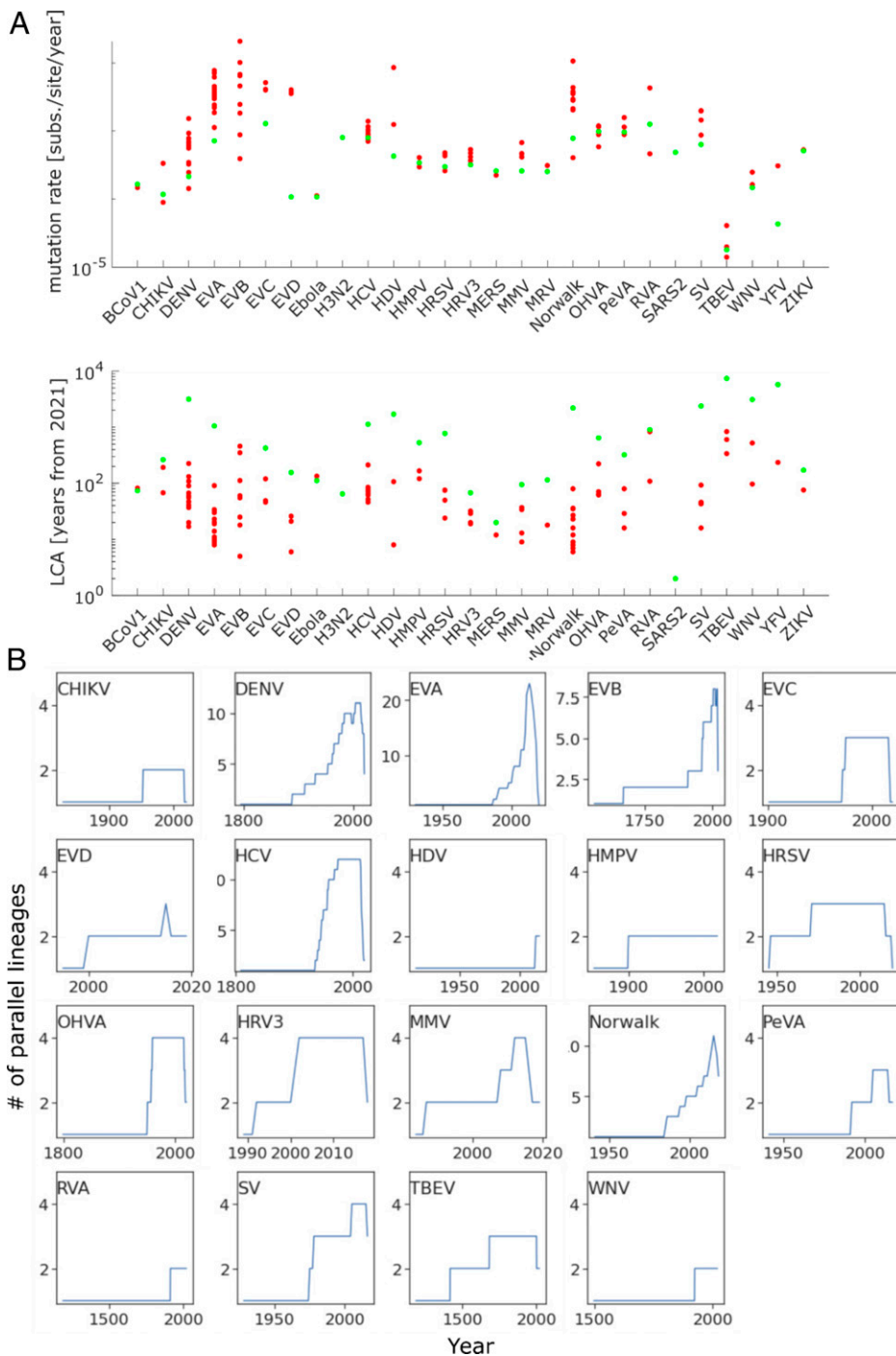


Fig. 3. GLs of human pathogenic RNA viruses throughout time. (A) Mutation rates (substitutions per site per year) for all main and GL phylogenetic trees used to construct genealogical trees (Top). Time to LCA in years from 2021 for all main and GL populations used to construct genealogical trees (Bottom). (B) Number of GLs per virus species circulating at the same time, based on the genealogical trees for each GL. For unannotated samples within a GL, the sampling date was estimated based on the date of the MRCA and root distance.

could support ongoing diversification and lineage turnover not observed in the human population. This is how some orthohepavirus A (OHVA) GLs (37) and possibly some TBEV (Tick Borne Encephalitis Virus) and WNV GLs (33, 38) could originate. However, even in this case, the maintenance of multiple niches with low turnover within human populations requires spatiotemporal or immunological separation. Regardless of the specific mechanisms underlying the apparent coexistence of noncompeting GLs, we sought to explore lineage-defining mutational signatures and to establish whether significant differences existed between the distributions of mutations within and between lineages.

Selective Pressures Acting on Human RNA Viruses. Selective processes are often categorized as diversifying, positive, or purifying, in contrast to neutral evolution via genetic drift (39–41).

We sought to probe the selective pressures involved in human pathogenic virus evolution by estimating the ratio of nonsynonymous to synonymous substitution rates (dN/dS), a gauge of protein-level selection (42, 43). Given that different genes are subject to distinct selective constraints and pressures, the dN/dS value was estimated separately for each viral protein-coding gene (44). Seeking to identify defining features of lineage emergence and maintenance, we would ideally estimate dN/dS across deep and shallow portions of each phylogeny separately. However, because most GLs antedated modern sequencing technologies, and, therefore, few samples located near the root were available, this approach was not feasible. To partially compensate for this lack of data, we compared dN/dS ratios for whole trees, which include deep branches connecting GLs, with those computed over each GL and ML (which are typically larger)

individually (see *Brief Methods* and *SI Appendix, Extended Methods* for details). The dN/dS estimates for whole trees ranged between 0.02 and 0.5 for most virus protein-coding genes, which is indicative of strong to moderate purifying selection, in line with previous results (45) (*SI Appendix, Fig. S9*). The few virus genes with elevated dN/dS ratios encode proteins that are either presented on the virion surface, such as human respiratory syncytial virus (HRSV) glycoprotein (G) and M-2 ($\sim 3.5\times$ above the species mean dN/dS), or human metapneumovirus (HMPV) SH and G ($\sim 3\times$ and $5\times$ above the species mean, respectively) (*SI Appendix, Figs. S9 and S10*), or are involved in interactions with the host immune system, for example, MMV V protein (46) ($\sim 4\times$ above the species mean) (*SI Appendix, Fig. S10*). These proteins are likely to experience positive selection, as described, for example, for HMPV G, where sites under positive selection were identified in the putative ectodomain (47). Elevated dN/dS values were also observed for some very short proteins, for example, the 6k peptide of DENV (*SI Appendix, Fig. S11*). However, such observations are sensitive to statistical artifacts and should be interpreted with caution. For OHVA ORF3, the dN/dS estimate was $\sim 4\times$ above the species mean (0.3; *SI Appendix, Fig. S13*), suggesting that this gene, which encodes an ion channel, plays a role in host adaptation following zoonosis (48).

Next, we computed dN/dS for each GL and ML individually (*SI Appendix, Figs. S10–S13 and S14–S17*, respectively). Despite considerable differences in size, generally, the results for GLs and MLs were comparable. For 12 of the 27 viruses studied (members of the order *Mononegavirales*, HMPV, HRSV, human respirovirus 3, MMV, and MRV; some flaviviruses ZIKV, YFV, TBEV, and YFV; HDV; MERS; and CHIKV), the dN/dS estimates for individual proteins as well as the mean for the whole tree differed little relative to the respective estimates for individual lineages (*SI Appendix, Figs. S10–S17*), with no indication of how selective pressures might have varied over time for any genes. In contrast, the GLs of enteroviruses (EVA-D) show elevated dN/dS , mainly among capsid proteins (*SI Appendix, Fig. S16*). Although frequent recombination among enteroviruses necessitates interpreting these results with caution (49), this finding, coupled with the observation that mutations in enterovirus capsid protein genes appear to be the primary lineage-defining features (see below), suggests a substantial change in the selective pressure acting on the capsid proteins between the periods of lineage emergence and subsequent maintenance. Notably, OHVA lineages show similarly elevated dN/dS for domain of unknown function 3729 (DUF3729) (up to 0.4) and ORF3 dN/dS (up to 0.3, which is also elevated relative to the species mean as discussed above) (*SI Appendix, Fig. S17*). Both these genes are likely to be involved in host adaptation following zoonosis (37, 48). Further, a twofold to fivefold increase in mean dN/dS was detected for DENV, WNV, and hepatitis C virus (HCV) GLs across most genes relative to the complete phylogeny (*SI Appendix, Fig. S15*). The interpretation of genome-wide elevation of dN/dS in GLs is more challenging and depends on whether the GL is newly emergent, possibly reflecting a period of rapid host adaptation and intense positive selection (45, 50). Given the distant dates predicted for the LCAs for these GLs and lack of lineage turnover, reduced selective pressure moving from stronger purifying selection toward neutral drift appears more likely. Overall, dN/dS analysis revealed little about potentially differing selective pressures acting within and between GLs, despite the apparent differences in substitution rates critical to the definition of the GLs themselves (as discussed above).

Intralineage and Interlineage Mutational Signatures. Gene-scale dN/dS analysis is often unable to uncover positive selection acting on specific sites or neighborhoods, which can occur in widely different backgrounds, from neutral drift to strong purifying selection (51). Identification of individual positively selected mutations can provide additional insight into differences between the evolutionary contexts of GL emergence and subsequent maintenance. Multiple, parallel nonsynonymous mutations comprise the most obvious indication of site-wise positive selection. With the prominent exception of SARS-CoV-2, for which we have previously identified up to 100 sites with recurrent amino acid replacements that are likely subject to positive selection (52), too few recurrent amino acid substitutions were detected for comparable analysis in the remaining viruses analyzed here, despite being the species with the largest number of genome sequences available.

Given the infeasibility of the direct, site-specific approach, we performed a genomic neighborhood analysis to compare interlineage and intralineage mutational signatures. First, amino acid sites were labeled according to three categories of amino acid substitutions (Fig. 4A; see *Brief Methods* and *SI Appendix, Extended Methods* for details): 1) multiple, deep (MD) substitutions which are “lineage defining,” being conserved in at least 90% of the samples within at least two GLs, but represented by different amino acid residues in each of these GLs—for example, consider the third amino acid of the CHIKV ORF gp1, 97% of the sequences in GL 1 contain a serine in that site, whereas 96% of the sequences in GL2 contain a proline in that site; 2) multiple, shallow (MS) substitutions occurring on multiple, independent occasions across GLs; and 3) all shallow (AS) substitutions occurring at least once within a GL, representing all “recent” events. We then computed site densities for each of

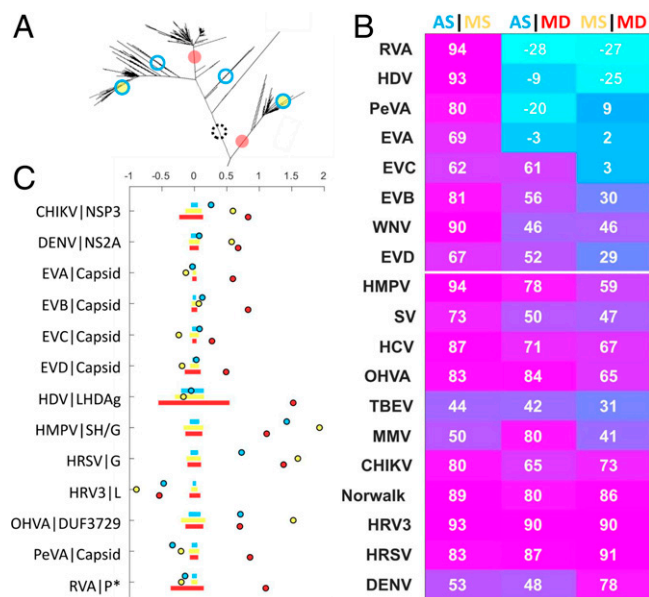


Fig. 4. Mutational signatures in human pathogenic RNA viruses vary little with tree depth. (A) Illustration of three amino acid site classes considered: 1) MD (red), 2) MS (yellow), and 3) AS (blue). Dashed circle represents deep, singular mutations which are excluded from this analysis. (B) Pearson correlation coefficient between site densities for all pairs of site classes across the genome computed over a moving average of 101 amino acids respecting mature peptide boundaries. Rows are sorted by the first column subtracted from the third column. (C) Log ratio of the mean site density across the specified peptide relative to the whole genome for select peptides. Bars are bounded by the 25th and 75th percentiles of simulated data drawn from the binomial distribution with n , total number of sites across the genome, trials of probability P , length of peptide/length of genome. RVA|P* represents the union of peptides P2-C, P-3A, and protease-3C.

these three categories over a sliding window of 101 amino acid sites, respecting protein boundaries.

We then examined the correlation between the site densities in these categories of amino acid substitutions across the genome (Fig. 4B and *SI Appendix*, Figs. S18–S24). In most cases, there was a strong, positive correlation between all three categories of amino acid substitutions, indicating that most genomic regions are subject to similar selection pressures during interlineage and intralineage evolution, with several notable exceptions (Fig. 4C). For enteroviruses, we observed an elevated MD site density within capsid proteins (VP1–4; *SI Appendix*, Figs. S18–S20), suggesting that capsid mutations are primarily lineage defining, but become less frequent once an epidemiological niche is established and occupied. This trend is consistent with the historical classification of enterovirus lineages by serotype, which is determined by the antigenic properties of the capsid proteins (53). As mentioned above, frequent intraspecies and interspecies recombination among all four types of enteroviruses requires caution when interpreting these results (49). This trend was similarly observed for the capsid proteins of the picornavirus Parechovirus A (PeVA) (VP0, VP1C, and VP1D; *SI Appendix*, Fig. S23).

Elevated MD and MS substitution densities were observed for HRSV G, potentially suggesting multiple residues evolving under positive selection throughout the entire course of evolution (including both lineage emergence and maintenance) of this immunologically exposed protein (54). Both MD and MS substitution densities are also increased in DENV NS2A and CHIKV nsp3. DENV NS2A is involved in virus replication and assembly and shows viroporin-like properties (55, 56). A detailed functional understanding of CHIKV nsp3 is lacking, although this protein is known to be part of the replication complex and is also involved in modulating the host cell's antiviral response (57). In line with the observation of elevated dN/dS in for OHVA DUF3729 in GLs compared to the whole population, MS substitution densities are elevated in this gene (*SI Appendix*, Fig. S22), suggesting that this poorly characterized protein contains multiple positively selected residues. These residues might have played a role in relatively recent host adaptation, but were not necessarily involved in the emergence of multiple lineages. The high MD substitution density observed for large human delta antigen (LHDAG) might result from statistical fluctuations given the short length (20 aa) of this peptide and should be interpreted with caution. As observed for the dN/dS analysis, we found few mutational signatures, which would shed light on different selection pressures acting within and between GLs. These observations seem to suggest that, although the tempo is variable, the mode of molecular evolution is broadly conserved from the deep to the shallow portions of each phylogeny, thus spanning considerable evolutionary distances.

Effective Population Size of Human Pathogenic RNA Viruses.

Another tool to indirectly assess the selective pressures shaping a phylogeny is to estimate the effective population size (N_e), which defines the timescale of population turnover across generations and thus can reveal major evolutionary events, including population bottlenecks (58). Assuming an evolutionary model, such as Wright–Fisher (59) or Moran (60), one can estimate the number of individuals per generation (that is, N_e) required for the observed rate of turnover in an idealized population. In what follows, we refer to “selection” as the sum of evolutionary pressures that promote lineage turnover. Although the background could vary from strong purifying selection to neutral drift, the occurrence of lineage turnover implicitly assumes some degree of positive selection in most scenarios.

In the context of lineage turnover, under strong selection, N_e is small, whereas lack of competition leads to larger N_e values over time (17). N_e can be inferred from the coalescence rate (Cr) estimated for any genealogical tree (17, 61, 62): $N_e * t \approx 1/Cr$ where t is the viral generation time (the time in days a virus needs to complete a transmission cycle from human to human). This expression enables a measurement of diversity and strength of selection among phylogenies represented by a single GL (e.g., H3N2, SARS-CoV-2), as well. Further, the census population size N (individuals present at each generation) can be estimated as $N = D * t$, where D is the number of yearly cases estimated. The N/N_e ratio may be used to quantify lineage turnover, where $N/N_e \gg 1$ indicates population bottlenecks, and $N/N_e \approx 1$ suggests stable population diversity (17). As tree topology, and hence N_e estimates, depend on selection strength and sampling effort (17, 58), we directly assessed the effect of sampling by randomly drawing up to either 10 or 100 samples per year, with three replicates each, for H3N2 and EVA as representative viruses with fast and slow turnover, respectively. We then used these reduced ensembles of isolates for genealogical tree construction (see *Brief Methods* and *SI Appendix, Extended Methods* for details; mutation rates and time of LCA are shown in *SI Appendix*, Fig. S25). We constructed two additional ensembles for each virus, composed of the same number of isolates selected above, this time maximizing the sequence diversity (as measured by the hamming distance between alignment rows; see *Brief Methods* and *SI Appendix, Extended Methods* for details). As a result, we obtained six trees evenly sampled over time (3×, e10, and e100) as well as two maximally diverse subtrees of the same size (d10 and d100) for both H3N2 and EVA.

We calculated the coalescence rates for all complete trees and for the H3N2 and EVA subtrees using the PACT package (<http://www.trevorbedford.com/pact>) (17) and estimated N_e (Fig. 5A). The N_e estimation was not performed for some zoonotic viruses, including TBEV and WNV, for which the generation time could not be reliably estimated. N/N_e ratios were calculated based on the best estimates of N_e for complete trees and those obtained after even and diverse sampling (H3N2 and EVA) (Fig. 5A and B). The estimated N_e values span more than two orders of magnitude, with H3n2 “even” 10 (H3N2e10) and EVD representing the extremes (Fig. 5A; N_e of around 400 and 270,000, respectively). Sampling was found to have a major effect on the estimates. The N_e estimates for EVA d100, d10, e100, and the whole tree were similar, whereas EVA e10 was about an order of magnitude lower. It should be noted that EVA e100 contained most sequences present in the whole tree. This trend was even more pronounced for H3N2, where the N_e estimates for d100, d10, and the whole tree were similar, whereas those for e10 and e100 were several orders of magnitude lower. It is first important to note that the estimates were not sensitive to the number of sequences present in the phylogeny, as illustrated by the equivalency of d10 and the whole phylogeny for both viruses, indicating that the differences observed between e10 and e100 and the whole phylogeny are not merely methodological artifacts. However, the estimate is sensitive to sampling, and, as could be expected, this sensitivity is more pronounced for viruses with fast lineage turnover. While perhaps unsurprising, this finding implies that N_e has been, and likely continues to be in this work, underestimated due to data limitations for most viruses and H3N2 in particular.

To illustrate the potential equivalency of reduced sampling and increased selection on N_e estimation, we simulated an ensemble of genealogical trees under a simple phenomenological model. Trees were iteratively constructed through the

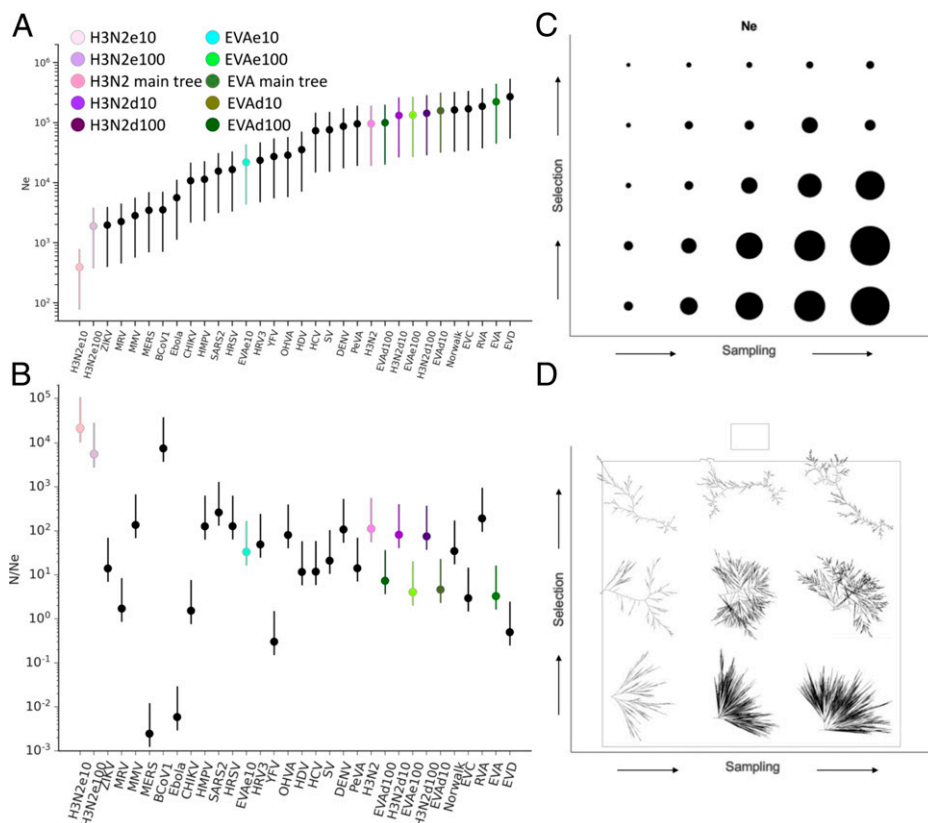


Fig. 5. Estimation of effective population sizes for human pathogenic RNA viruses. (A) N_e estimated for genealogical trees. Bars represent varying generation time t for each virus ranging between 0.5 and 5 times the value corresponding to the filled circle. For H3N2 and EVA, N_e estimates for evenly sampled trees (up to 10 or 100 samples per year, e10 and e100, respectively) and diverse sampled trees (d10 and d100) are also displayed. (B) N/N_e ratios, where N is the census population. Bars represent varying N between 0.5 and 5 times the value corresponding to the filled circle. N and t estimates are shown in *SI Appendix, Table S10*. As in C, N/N_e estimates are shown for evenly and diverse sampled trees for EVA and H3N2. Color code is as in C. (C) N_e (with generation time t fixed to one for N_e calculation for all trees) for simulated trees with varied selection strength and sampling density (N_e is represented by the circle area). Sampling density (as multifurcation multiplicity) varied from two to six. Selection strength was simulated by positioning new leaves each iteration at the 0th, 25th, 50th, 75th, or 95th percentile of the root distance distribution for the entire tree (see *SI Appendix, Extended Methods* for details). (D) Examples of simulated trees with varied selection strength and sampling density.

addition of clades representing local sequencing efforts. Increased sequencing efforts were modeled by changing the number of sequences in each clade from two to six. Increased selection strength was modeled by changing the placement of these clades on the tree, relative to the root, from the prior iteration. A root distance threshold was set to be the 0th, 25th, 50th, 75th, or 95th percentile of the root distance distribution for all leaves at the prior iteration, with higher thresholds corresponding to new isolates being placed farther from the tree root and representing increased selection (Fig. 5C and *SI Appendix, Fig. S26 A–C*). Although selection and undersampling result in qualitatively different tree topologies (Fig. 5D), their effects cannot be disentangled from N_e analysis alone. Furthermore, sensitivity to undersampling is more pronounced under high selection than under low selection (Fig. 5D and *SI Appendix, Fig. S26D*). These effects must be considered when evaluating the expectation that genetic diversity (and hence N_e) plateaus earlier in a growing census population N when selection is strong (17). This challenge is reflected in the damped increase of N_e from H3N2 e10 to e100 when compared to the increase from EVA e10 to e100 (approximately fourfold and eightfold, respectively; Fig. 5A).

Keeping these sensitivities in mind, we proceeded to examine the N/N_e ratios. High N/N_e ratios can be an indicator of population bottlenecks. The highest N/N_e ratio was observed for H3N2e10 (e100 was similar); in contrast, the estimate for the whole phylogeny was about 200-fold lower (Fig. 5D), within the range of the majority of the other viruses. Thus, sampling efforts can substantially affect the interpretation of the N/N_e estimation, moving H3N2 from an outlier associated with extreme bottlenecks to typical behavior. As discussed above, it has been well established that H3N2 is subject to pronounced population bottlenecks as a result of alternating periods of stasis and rapid host adaptation (16–18). However, the results presented here

emphasize that, on shorter timescales, the transmission dynamics of local outbreaks play a larger role in determining the extent of the diversity of the H3N2 population (63), as was the case for the majority of viruses studied in this work. BCoV1 also demonstrated a high N/N_e ratio, within the range of H3N2 e10 and e100 (Fig. 5B). Although this observation could simply result from insufficient sampling, given the high incidence of this virus (*SI Appendix, Fig. S27D*), it might point to pronounced population bottlenecks during the evolution of the BCoV1 population (although not of comparable magnitude to those for H3N2; see below). In contrast, other viruses seem to experience less severe bottlenecks and maintain greater genetic diversity (e.g., MRV and enteroviruses A, C, D). The low N/N_e values for MERS, ZIKV, and CHIKV likely result from an underestimation of N due to large animal reservoirs that might impact estimates of N for H3N2 as well. Furthermore, we observed a positive correlation, albeit not statistically significant, between N/N_e estimates and the extent of lineage turnover as quantified by the Shannon entropy (*SI Appendix, Fig. S28A*). This trend is likely perturbed by uncertainty in estimating N , and removing CHIKV slightly improves the correlation (*SI Appendix, Fig. S28B*).

Whereas insufficient sampling can lead to an underestimation of both N and N_e , the complete unavailability of genomes from premodern periods can, perhaps counterintuitively, lead to an overestimation of N_e . As discussed above, GLs are separated by periods (branches of the tree) with apparently different substitution rates. These branches are often deep within the tree topology and sparsely populated with leaves (if at all), making the assignment of a global model for substitution rates statistically dubious. This can result in inaccurate deep branch lengths for genealogical trees and substantially change the predicted date for the LCA. This date, as well as the predicted dates of other deep nodes, is used to estimate the effective population size. Given these limitations, we sought to establish a lower bound

for the effective population size, which still preserves all GLs, through the construction of truncated global genealogical trees or “grafted trees.” The LCA of each grafted tree is set to the LCA of the oldest GL, and the remaining GLs are connected to this (multifurcated) root through branches preserving the LCA of each respective GL (see *SI Appendix, Fig. S27A* and *Brief Methods* or *SI Appendix, Extended Methods* for details). We proceeded to estimate N_e for each grafted tree as well as for each GL separately (*SI Appendix, Fig. S27B*). By construction, N_e estimates for grafted trees are generally significantly smaller than those for complete trees and larger than those for individual GLs. Notably, the N_e estimate for H3N2 (which is represented by a single GL) is the second-highest value observed (after rhinovirus A) among the viruses studied when this lower bound is considered. This counterintuitive finding emphasizes another facet of the sensitivity of N_e estimation to data availability.

These sensitivities are evidently greater within individual GLs, which represent only a subset of the viral diversity for each species. These limitations notwithstanding, in an effort to characterize lineage turnover within individual GLs, we analyzed skyline plots representing the time to the most recent common ancestor (TMRCA) of all clades present at a given time point and diversity within the population over time (as measured by the average time for any two isolates to coalesce, that is, to find their most recent common ancestor) for individual GLs and the complete population (*SI Appendix, Fig. S29*). The average population diversity can be displayed as the mean diversity per year (*SI Appendix, Fig. S29A*). Populations with high turnover, such as H3N2, show a low average diversity per year, whereas those with low turnover are characterized by high diversity. In general, the skyline plots and mean diversity values for complete phylogenies correspond well to N_e and N/N_e estimations, supporting slow lineage turnover for most of the viruses analyzed. For example, H3N2 displays $\sim 4\times$ and $\sim 8\times$ lower mean diversity per year compared to BCoV1 and ZIKV, respectively. This observation suggests BCoV1, despite having high N_e and N/N_e values, has a slower population turnover compared to H3N2. Of note is that evidence of high intra-GL turnover was obtained for a few GLs (as demonstrated by a mean diversity in

the range of H3N2 or BCoV1). The two principal examples are HRSV GL2 and Norwalk GL3 (*SI Appendix, Fig. S29 C and D*). The majority of GLs show mean diversity within the range of viral populations with low turnover, and individual GLs generally display lower diversity than each respective whole population (*SI Appendix, Fig. S29A*). Although it is expected that any sub-population has a lower diversity than the larger population from which it is selected, a higher intra-GL turnover rate relative to the inter-GL turnover might additionally lower diversity within individual GLs. In the next section, we demonstrate that inter-GL turnover is substantially nonneutral, further supporting our conclusion on the existence of independent environmental and epidemiological niches among pathogenic RNA viruses.

Significant Deviations from Neutrality Confirm Noncompeting Niches.

In principle, neutral evolution can yield multiple, coexisting lineages. Although, per the definition proposed above, GLs cannot emerge within a phylogeny resulting from an evolutionary process with a constant substitution rate, neutral processes could result in trees qualitatively similar to those of the genealogical trees inferred for the viruses analyzed here. To explicitly demonstrate the deviation from neutrality of these viral evolutionary histories, we simulated trees under a neutral branching process [Yule–Harding (64, 65)] and shallowly subsampled these, to reflect the fact that efficient sequencing technologies were unavailable (for most viruses) for the majority of the period between the present and the predicted date of the most recent common ancestor (MRCA) (see *SI Appendix, Extended Methods* for details). For each virus, and each real and simulated tree (which has the same number of leaves as the respective real tree), the number of parallel branches from the root to the tree tips was tabulated, and the date corresponding to the deepest leaves was estimated (Fig. 6A). In order to probe the coalescence rate among, rather than within, extant lineages, the period from the root to the deepest leaves was further examined. The number of branches from the deepest leaves back toward the root was then fit to a power law model (Fig. 6B; see *SI Appendix, Extended Methods* for details).

The number of parallel branches was generally greater among simulated trees than among the respective viral genealogical trees

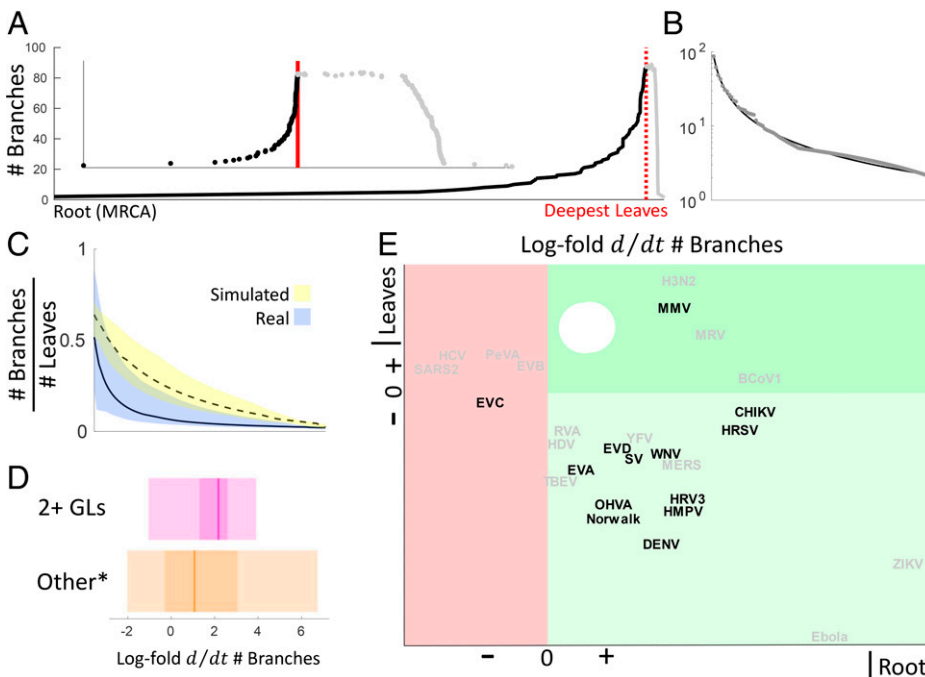


Fig. 6. Probing pathogenic RNA virus evolution for neutrality. (A) The number of parallel branches in the genealogical tree for the (sub-sampled) Norwalk virus from the root (188 BCE) to the tree tips. The red line indicates the position of the deepest leaves. (Inset) Rescaled so the periods from the root to the deepest leaves and that from the deepest leaves to the tips are evenly displayed; branch count is displayed at each node. (B) The number of parallel branches of Norwalk virus moving from the deepest leaves back to the root overlaid with fitted power-law. (C) The number of parallel branches from deepest leaves to root normalized by the number of leaves in each tree for all simulated and real trees fit to power-laws. Shading identifies median 50% of data and curves identify median. (D) Log-fold change (simulated/real) of the time derivatives for the number of parallel branches evaluated at the root for those viruses with at least two GLs vs. all others. Shading identifies minima, maxima, and the median 50% of data. Curves identify median. Viruses with at least two GLs, but for which the threshold correlation was reduced (from the default 0.8), appear in “Other*.” (E) Log-fold change of the time derivative at the deepest leaves vs. that at the root. Viruses containing two or more GLs with a correlation coefficient of 0.8 or above are displayed in black, and others are shown in gray.

for the duration of this period (Fig. 6C). The log-fold difference (simulated/real) in the derivative evaluated at each tree root, which measures the coalescence rate among (rather than within) extant lineages, was positive, indicating that lineage turnover was slower than that under a neutral evolution model. This deviation from neutrality was greater among viruses containing two or more GLs with a correlation coefficient of 0.8 (and disregarding those with a poorer correlation; see *SI Appendix, Extended Methods* for details) than among other viruses (Fig. 6D). For most viruses, this log-fold difference was negative when evaluated at the deepest leaves but became positive moving toward the root (Fig. 6E and *SI Appendix, Fig. S41*). Of the five viruses (EVB, enterovirus C [EVC], HCV, PeVA, and SARS-CoV-2) for which the log-fold difference evaluated at the root was negative, only one (EVC) contained multiple GLs. Both the EVC and HCV phylogenies are additionally complicated by the need to remove laboratory-related, vaccine-related, or clinical sequences, which do not circulate globally. Recombination among PeVA genotypes (66) (each with few sequences, leading to poorly correlated GLs) complicates the inference of a complete genealogical tree for this virus as well. As mentioned above, EVB lacks a global genealogical tree and is represented here by a single GL. SARS-CoV-2 has recently emerged and constitutes a single GL in its entirety. The log-fold difference for these and other viruses represented by a single GL should be interpreted with caution. For these trees, the deepest leaf is relatively nearer the root, and the portion of the tree studied is much smaller than in other viruses. Furthermore, interlineage turnover (between the single extant lineage and prior, extinct lineages) may not be measurable with these methods.

Thus, among the viruses studied here, interlineage turnover is typically slower than the neutral expectation, and intralinear turnover is faster than the neutral expectation (as described for select GLs above; *SI Appendix, Fig. S29 C and D*). These results are consistent with our interpretation that human pathogenic RNA virus populations are largely organized into stable, coexisting lineages (GLs). Each GL occupies a distinct epidemiological niche, within which competition leads to increased lineage turnover compared to neutrality. In contrast, among the GLs, there is little competition, which leads to decreased lineage turnover relative to neutrality.

Discussion

Here we present a comprehensive phylodynamic analysis of monopartite human pathogenic RNA viruses (and H3N2 hemagglutinin) in an effort to establish global trends in viral evolution in human populations. Despite data limitations, the viruses studied in this work span a wide variety of viral life history characteristics. This lends considerable generality to the study, while making it outside the scope of this work to investigate many features specific to individual lifestyles [for example, intrahost diversity, symptom characteristics, or acute vs. chronic infection (2)]. Given this diversity, the commonalities we observe among the virus phylogenies constructed are notable. Consistent with the conclusions of previous efforts (45, 67), we observed moderate to strong purifying selection among all viruses.

Nearly all virus populations are characterized by low rates of lineage turnover, and most consist of multiple, coexisting GLs, monophyletic groups of sequences separated by periods with apparently different substitution rates. Despite these differing substitution rates, dN/dS and genomic neighborhood analysis revealed little about how selective pressures might have differed between the early period of GL formation and the subsequent period defined by persistent coexistence. This lack of resolution

seems to suggest that, although the tempo is variable, the mode of molecular evolution is broadly conserved from the deep to the shallow portions of each phylogeny, spanning considerable evolutionary distances. The distribution of lineage-defining mutations across the virus genome is similar to that of shallow, repeated mutations for almost all viruses, indicating that positive selection affects sites in the same neighborhoods during both periods (Fig. 4). The lineage-defining role of enterovirus capsid proteins was the principle exception observed, in line with the traditional serotype classification (12). Other virus proteins with different intralinear and interlineage mutational signatures, which might provide insight into ongoing host and/or vector adaptive evolution, are OHVA DUF3729, DENV NS2A, and CHIKV NSP3. In the case of CHIKV E1-A226V, NSP3 has been demonstrated to play an important role in the adaptation to the vector *Aedes albopictus* (68).

The low inter-GL turnover, below the neutral expectation, and broadly stable mutational signatures appear to be indicative of weak, if any, competition among GLs, suggesting that each GL occupies an independent epidemiological niche (Figs. 5 and 6). Such niches could be maintained in a variety of ways, the most obvious possibility being regionality and/or seasonality. Although these factors can explain the persistence of some GLs identified in this work, the majority do not show regional localization, and none display biannual (or more coarse-grained) temporal trends (the limit of time resolution we can reliably detect). At sufficiently large evolutionary distances, niches can be defined by immunological differences, which can overcome immune cross-protection, allowing the same host to be infected by multiple subpopulations largely independent of prior infections, as seems to be the case for picornaviruses and HRSV (54). These effects are also insufficient to account for the stability of most GLs. We suggest that, in many if not most cases, niches are maintained through a series of localized outbreaks such that only a small fraction of the global susceptible population is infected at any given time. Under this scenario, even lineages that do not overcome immune cross-protection can coexist within short distances of one another. Furthermore, extensive environmental transfer or fragmented animal reservoir populations could play a role. Virions that can persist for extended periods of time outside of the (identified) host or vector might maintain the genetic diversity of a lineage during time periods when no active infections from that lineage occur.

As a result of globalization, disparate communities are becoming increasingly connected, which might lead to increased interaction between previously separated lineages, enhancing between-lineage competition within the viral population. This effect has been demonstrated already for DENV in Thailand, where multiple lineages typically coexist throughout the country, with a well-defined pattern of dissemination. However, within densely populated areas of Bangkok, genomic analysis pointed to increased competition and lineage turnover (69).

Conclusions

Phylodynamic analysis revealed multiple cocirculating lineages (GLs) for the majority of human pathogenic RNA viruses, separated by periods of apparently different substitution rates within the phylogeny. The dN/dS and genomic neighborhood analysis yielded surprisingly little evidence of different selection pressures acting within and between GLs, suggesting that, whereas the tempo of molecular evolution is variable, its mode is broadly conserved. This slow lineage turnover, below the neutral expectation, suggests each GL occupies an independent epidemiological niche, with little inter-GL competition. No pronounced patterns of

regional or temporal separation of the GLs were detected, suggesting that the stability of the GLs primarily stems from limited outbreaks within small communities, so that only a small fraction of the global susceptible population is infected at any time. These results raise the, perhaps pressing, question, How will increased host–host contact resulting from globalization affect viral evolution? Could new or renewed competition emerge among lineages of endemic viruses to drive diversification, evolution of increased pathogenicity, or even virus speciation? With these questions in mind, we emphasize that, in addition to zoonotic events, the ongoing surveillance of familiar, endemic viruses deserves global attention in effort to mitigate or prevent future pandemics.

Brief Methods

Genomes were retrieved for all viruses except IAV H3N2 and SARS-CoV-2 from NCBI virus (24). Members of related viral families were used to construct an outgroup when possible. IAV H3N2 (flu H3N2) segment HA was retrieved from the NCBI flu database (70). The SARS-CoV-2 tree and alignment analyzed in this work was subsampled from a larger alignment consisting of all high-quality genomes that were available as of January 8, 2021 in the GISAID database (25), as previously described (52). Subsampling was conducted to maximize the sequence diversity. Acknowledgments for the GISAID deposited sequences used in this study are displayed in *SI Appendix, Table S3*. Subalignments were considered for H3N2 and EVA, principally for the purpose of effective population size analysis. In all cases, sequences were harmonized to DNA (e.g., U was transformed to T to amend software compatibility) and aligned with MAFFT (26), using default settings. Sequences were clustered according to 100% identity with no coverage threshold using Cd-hit (71), and, otherwise, default settings for MERS and H3N2.

The longest sequence from each cluster was selected as a representative. Exterior ambiguous characters were removed, and sequences with more than 10 remaining ambiguous characters (“N”) were discarded. Outliers based on hamming distance to the nearest neighbor and consensus were identified and removed from the set. Sites corresponding to protein-coding ORFs (open reading frames) were then mapped to the alignments, and non-coding regions were discarded. Common gaps corresponding to multiples of three nucleotides were maintained as “true” insertions or deletions and mapped into the frame if necessary. Unique alignment rows were identified. Samples related to laboratory experiments, vaccine-related sequences, and patents were pruned based on an automated keyword search

Dates and locations of isolation are available for many isolates reported as calendar dates and city or country/administrative region of origin. These dates are referenced as calendar dates in the main text and as date indices (number of days before/after January 1, 1950) in *SI Appendix*. For the regional analysis, the latitude and longitude of each city of origin or a representative city for each country/administrative region of origin was identified from simplemaps (<https://simplemaps.com/data/world-cities>) (72). Ambiguity in metadata assignments was not problematic (*SI Appendix, Fig. S30*).

With the exception of SARS-CoV-2 and H3N2, tree topology was optimized using IQ-TREE (27) with the evolutionary model fixed to GTR+F+G4 and the minimum branch length decreased from the default 10e-6 to 10e-7 (options: -m GTR+F+G4 -st DNA -blmin 0.0000001). For SARS-CoV-2, the tree was drawn from the global topology previously described (52). The global H3N2 tree was approximated using FastTree (73) specifying GTR; a four-category gamma distribution; no support values; and

using a previously constructed maximum diversity subtree as a constraint (compiled at double precision, options: -nt -gtr -gamma -cat 4 -nosupport -constraints). Trees were rooted according to the position of an outgroup when possible, and by date or mid-point otherwise.

Viral lineages were both manually selected, based on available metadata, and algorithmically selected, into correlated clades we call GLs. GLs are defined as monophyletic clades with a strong correlation between the sequencing date and the distance to the tree root. Trees were used to construct date-constrained, genealogical trees using least-square dating (with software LSD2) (31). We considered the Shannon entropy of the clade distribution calculated over sliding windows based on the known or estimated date of isolation or based on distance to the tree root as an explicit measure of lineage turnover.

Fitch Traceback (74) was used to estimate ancestral states. Three classes of amino acid sites were identified on the basis of the nonsynonymous mutations within each site: 1) MD substitutions are “lineage defining”; 2) MS sites. 3) AS sites. We computed the site density of each class over a sliding window to assess signatures of positive selection. Selection pressures were also assessed through *dN/dS* analysis using PAML (44).

The effective population size *N_e* and the ratio of the census population size *N* over *N_e* was estimated as previously described (17) (<http://www.trevorbedford.com/pact>). Viral diversity (average time of any pair of leaves at a given time point to find their LCA) and average TMRCA over time were calculated with the PACT package as well. In order to demonstrate the potential equivalency between the impacts of selection strength and sampling density on effective population size, we additionally simulated an ensemble of trees.

Viral phylogenies were compared to trees simulated under a neutral model, shallowly sampled to reflect the fact that the MRCAs are predicted to have circulated much earlier than the invention of sequencing technologies. When necessary, maximally diverse subtrees of the respective viral phylogeny with the same number of leaves as in each simulated tree were generated. For all trees, real and simulated, the number of parallel branches from the root to the tree tips was computed. In order to probe the coalescence rate among, rather than within, extant lineages, the date associated with the deepest leaves was determined, and the period from the root to that point was further examined. The number of branches from the deepest leaves back toward the root was then fit to a single-parameter power-law expression [based on expectations of a power-law distribution of node descendants (75)] (*SI Appendix, Figs. S31–S35*). The derivative of each curve (*SI Appendix, Figs. S36–S40*) as well as the log-fold difference between the derivatives of the simulation and the real genealogical tree for each virus were computed (*SI Appendix, Fig. S41*).

Data Availability. The datasets generated and/or analyzed during the current study are available as Supplementary Data at Zenodo, <https://zenodo.org/record/5711959>, as well as through FTP, https://ftp.ncbi.nih.gov/pub/wolfi_suppl/virNiches/. Original virus sequences are publicly available for all viruses, except SARS-CoV-2, at NCBI virus (24). SARS-CoV-2 sequences are available at GISAID (25).

ACKNOWLEDGMENTS. We thank E.V.K. group members for helpful discussions. N.D.R., Y.I.W., P.M., and E.V.K. are supported by the Intramural Research Program of the NIH (National Library of Medicine).

Author affiliations: ^aNational Center for Biotechnology Information, National Library of Medicine, Bethesda, MD 20894; ^bBroad Institute of MIT and Harvard, Cambridge, MA 02142; ^cHHMI, Massachusetts Institute of Technology, Cambridge, MA 02139; ^dMcGovern Institute for Brain Research, Massachusetts Institute of Technology, Cambridge, MA 02139; ^eDepartment of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, MA 02139; and ^fDepartment of Biological Engineering, Massachusetts Institute of Technology, Cambridge, MA 02139

1. C. A. Suttle, Marine viruses—Major players in the global ecosystem. *Nat. Rev. Microbiol.* **5**, 801–812 (2007).
2. P. J. Chisholm, J. W. Busch, D. W. Crowder, Effects of life history and ecology on virus evolutionary potential. *Virus Res.* **265**, 1–9 (2019).
3. P. T. Johnson, J. C. de Roode, A. Fenton, Why infectious disease research needs community ecology. *Science* **349**, 1259504 (2015).
4. L. Yan, R. A. Neher, B. I. Shraiman, Phylogenetic theory of persistence, extinction and speciation of rapidly adapting pathogens. *eLife* **8**, e42405 (2019).
5. J. Marchi, M. Lässig, T. Mora, A. M. Walczak, Multi-lineage evolution in viral populations driven by host immune systems. *Pathogens* **8**, 115 (2019).
6. World Health Organization, Coronavirus (COVID-19) dashboard. <https://covid19.who.int/>. Accessed 20 May 2022.
7. W. T. Harvey *et al.*, COVID-19 Genomics UK (COG-UK) Consortium, SARS-CoV-2 variants, spike mutations and immune escape. *Nat. Rev. Microbiol.* **19**, 409–424 (2021).
8. N. Rochman, Y. Wolf, E. V. Koonin, Evolution of human respiratory virus epidemics [version 2; peer review: 2 approved]. *F1000 Res.* **10**, 447 (2021).
9. P. Simmonds, P. Aiewsakun, A. Katzourakis, Prisoners of war - Host adaptation and its constraints on virus evolution. *Nat. Rev. Microbiol.* **17**, 321–328 (2019).
10. J. E. Bryant, E. C. Holmes, A. D. Barrett, Out of Africa: A molecular perspective on the introduction of yellow fever virus into the Americas. *PLoS Pathog.* **3**, e75 (2007).
11. Y. Li, Z. Yang, Adaptive diversification between yellow fever virus West African and South American lineages: A genome-wide study. *Am. J. Trop. Med. Hyg.* **96**, 727–734 (2017).
12. L. Brouwer, G. Moreni, K. C. Wolthers, D. Pakr, World-wide prevalence and genotype distribution of enteroviruses. *Viruses* **13**, 434 (2021).
13. J. L. Geoghegan, E. C. Holmes, Evolutionary virology at 40. *Genetics* **210**, 1151–1162 (2018).
14. B. T. Grenfell *et al.*, Unifying the epidemiological and evolutionary dynamics of pathogens. *Science* **303**, 327–332 (2004).
15. H. Kim, R. G. Webster, R. J. Webby, Influenza virus: Dealing with a drifting and shifting pathogen. *Viral Immunol.* **31**, 174–183 (2018).
16. Y. I. Wolf, C. Viboud, E. C. Holmes, E. V. Koonin, D. J. Lipman, Long intervals of stasis punctuated by bursts of positive selection in the seasonal evolution of influenza A virus. *Biol. Direct* **1**, 34 (2006).
17. T. Bedford, S. Cobey, M. Pascual, Strength and tempo of selection revealed in viral gene genealogies. *BMC Evol. Biol.* **11**, 220 (2011).
18. K. Koelle, S. Cobey, B. Grenfell, M. Pascual, Epochal evolution shapes the phylogenetics of inter-pandemic influenza A (H3N2) in humans. *Science* **314**, 1898–1903 (2006).
19. P. A. Rota *et al.*, Measles. *Nat. Rev. Dis. Primers* **2**, 16049 (2016).
20. S. R. Hamblin, P. A. White, M. M. Tanaka, Viral niche construction alters hosts and ecosystems at multiple scales. *Trends Ecol. Evol.* **29**, 594–599 (2014).
21. E. M. Volz, K. Koelle, T. Bedford, Viral phylogenetics. *PLoS Comput. Biol.* **9**, e1002947 (2013).
22. C. J. Houldcroft, M. A. Beale, J. Breuer, Clinical and biological insights from viral genome sequencing. *Nat. Rev. Microbiol.* **15**, 183–192 (2017).
23. L. Liu, Fields Virology, 6th edition. *Clin. Infect. Dis.* **59**, 613 (2014).
24. E. L. Hatcher *et al.*, Virus variation resource - Improved response to emergent viral outbreaks. *Nucleic Acids Res.* **45**, D482–D490 (2017).
25. Y. Shu, J. McCauley, GISAID: Global initiative on sharing all influenza data - From vision to reality. *Euro Surveill.* **22**, 30494 (2017).
26. K. Katoh, K. Misawa, K. Kuma, T. Miyata, MAFFT: A novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* **30**, 3059–3066 (2002).
27. L.-T. Nguyen, H. A. Schmidt, A. von Haeseler, B. Q. Minh, IQ-TREE: A fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol. Evol.* **32**, 268–274 (2015).
28. S. O. Scholle, R. J. Ypma, A. L. Lloyd, K. Koelle, Viral substitution rate variation can arise from the interplay between within-host and epidemiological dynamics. *Am. Nat.* **182**, 494–513 (2013).
29. World Health Organization, *Ebola Virus Disease Fact-Sheets* (World Health Organization, 2021).
30. S. T. Jacob *et al.*, Ebola virus disease. *Nat. Rev. Dis. Primers* **6**, 13 (2020).
31. T.-H. To, M. Jung, S. Lycett, O. Gascuel, Fast dating using least-squares criteria and algorithms. *Syst. Biol.* **65**, 82–97 (2016).
32. A. B. Schneider *et al.*, Updated phylogeny of chikungunya virus suggests lineage-specific RNA architecture. *Viruses* **11**, 798 (2019).
33. C. Chancey, A. Grinev, E. Volkova, M. Rios, The global ecology and epidemiology of West Nile virus. *BioMed Res. Int.* **2015**, 376230 (2015).
34. F. Di Giallonardo *et al.*, Fluid spatial dynamics of West Nile virus in the United States: Rapid spread in a permissive host environment. *J. Virol.* **90**, 862–872 (2015).
35. W. Wang *et al.*, Assembly and infection efficacy of hepatitis B virus surface protein exchanges in 8 hepatitis D virus genotype isolates. *J. Hepatol.* **75**, 311–323 (2021).
36. M. Moriyama, W. J. Hugentobler, A. Iwasaki, Seasonality of respiratory viral infections. *Annu. Rev. Virol.* **7**, 83–101 (2020).
37. S. Lin, Y. J. Zhang, Advances in hepatitis E virus biology and pathogenesis. *Viruses* **13**, 267 (2021).
38. M. Kubinski *et al.*, Tick-borne encephalitis virus: A quest for better vaccines against a virus on the rise. *Vaccines (Basel)* **8**, 451 (2020).
39. P. T. Dolan, Z. J. Whitfield, R. Andino, Mechanisms and concepts in RNA virus population dynamics and evolution. *Annu. Rev. Virol.* **5**, 69–92 (2018).
40. M. Lynch *et al.*, Genetic drift, selection and the evolution of the mutation rate. *Nat. Rev. Genet.* **17**, 704–714 (2016).
41. M. Kimura, The neutral theory of molecular evolution: A review of recent evidence. *Jpn. J. Genet.* **66**, 367–386 (1991).
42. L. D. Hurst, The Ka/Ks ratio: Diagnosing the form of sequence evolution. *Trends Genet.* **18**, 486 (2002).
43. H. A. Hejase, N. Dukler, A. Siepel, From summary statistics to gene trees: Methods for inferring positive selection. *Trends Genet.* **36**, 243–258 (2020).
44. Z. Yang, PAML 4: Phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* **24**, 1586–1591 (2007).
45. J. J. Lin, M. J. Bhattacharjee, C. P. Yu, Y. Y. Tseng, W. H. Li, Many human RNA viruses show extraordinarily stringent selective constraints on protein evolution. *Proc. Natl. Acad. Sci. U.S.A.* **116**, 19009–19018 (2019).
46. K. Takeuchi, S. I. Kadota, M. Takeda, N. Miyajima, K. Nagata, Measles virus V protein blocks interferon (IFN)-alpha/beta but not IFN-gamma signaling by inhibiting STAT1 and STAT2 phosphorylation. *FEBS Lett.* **545**, 177–182 (2003).
47. J. Pappenburg *et al.*, Genetic diversity and molecular evolution of the major human metapneumovirus surface glycoproteins over a decade. *J. Clin. Virol.* **58**, 541–547 (2013).
48. P. P. Primadharisni, S. Nagashima, H. Okamoto, Mechanism of cross-species transmission, adaptive evolution and pathogenesis of hepatitis E virus. *Viruses* **13**, 909 (2021).
49. Z. Kyriakopoulou, V. Plaka, G. D. Amoutzias, P. Markoulatas, Recombination among human non-polio enteroviruses: Implications for epidemiology and evolution. *Virus Genes* **50**, 177–188 (2015).
50. E. C. Holmes, Patterns of intra- and interhost nonsynonymous variation reveal strong purifying selection in dengue virus. *J. Virol.* **77**, 11296–11298 (2003).
51. S. Kryazhimskiy, J. B. Plotkin, The population genetics of dN/dS. *PLoS Genet.* **4**, e1000304 (2008).
52. N. D. Rochman *et al.*, Ongoing global and regional adaptive evolution of SARS-CoV-2. *Proc. Natl. Acad. Sci. U.S.A.* **118**, e2104241118 (2021).
53. M. S. Oberste, K. Maher, D. R. Kilpatrick, M. A. Pallansch, Molecular evolution of the human enteroviruses: Correlation of serotype with VP1 sequence and application to picornavirus classification. *J. Virol.* **73**, 1941–1948 (1999).
54. G. W. Wertz, R. M. Moudy, Antigenic and genetic variation in human respiratory syncytial virus. *Pediatr. Infect. Dis. J.* **23**, S19–S24 (2004).
55. C. J. Neufeldt, M. Cortese, E. G. Acosta, R. Bartschlagler, Rewiring cellular networks by members of the Flaviviridae family. *Nat. Rev. Microbiol.* **16**, 125–142 (2018).
56. G. Shrivastava *et al.*, NS2A comprises a putative viroporin of Dengue virus 2. *Virulence* **8**, 1450–1456 (2017).
57. J. J. Fros, G. P. Pijlman, Alphavirus infection: Host cell shut-off and inhibition of antiviral responses. *Viruses* **8**, 166 (2016).
58. J. Wang, E. Santiago, A. Caballero, Prediction and estimation of effective population size. *Heredity* **117**, 193–206 (2016).
59. Y. Ishida, A. Rosales, The origins of the stochastic theory of population genetics: The Wright-Fisher model. *Stud. Hist. Philos. Biol. Biomed. Sci.* **79**, 101226 (2020).
60. P. A. P. Moran, Random processes in genetics. *Math. Proc. Camb. Philos. Soc.* **54**, 60–71 (1958).
61. K. Strimmer, O. G. Pybus, Exploring the demographic history of DNA sequences using the generalized skyline plot. *Mol. Biol. Evol.* **18**, 2298–2305 (2001).
62. O. G. Pybus, A. Rambaut, P. H. Harvey, An integrated framework for the inference of viral population history from reconstructed genealogies. *Genetics* **155**, 1429–1437 (2000).
63. J. T. McCrone, A. S. Lauring, Genetic bottlenecks in intraspecific virus transmission. *Curr. Opin. Virol.* **28**, 20–25 (2018).
64. S. Vinh, A. Fuehrer, A. von Haeseler, Random Tree-Puzzle leads to the Yule-Harding distribution. *Mol. Biol. Evol.* **28**, 873–877 (2011).
65. G. U. Yule II, A mathematical theory of evolution, based on the conclusions of Dr. J. C. Willis, F. R. S. *Philos. Trans. R. Soc. Lond., B. Contain. Pap. Biol. Character* **213**, 21–87 (1925).
66. G. Shah, J. L. Robinson, The particulars on parechovirus. *Can. J. Infect. Dis. Med. Microbiol.* **25**, 186–188 (2014).
67. J. O. Wertheim, S. L. Kosakovsky Pond, Purifying selection can obscure the ancient age of viral lineages. *Mol. Biol. Evol.* **28**, 3355–3365 (2011).
68. K. A. Tsetsarkin *et al.*, Multi-peaked adaptive landscape for chikungunya virus evolution predicts continued fitness optimization in *Aedes albopictus* mosquitoes. *Nat. Commun.* **5**, 4084 (2014).
69. H. Salje *et al.*, Dengue diversity across spatial and temporal scales: Local structure and the effect of host population size. *Science* **355**, 1302–1306 (2017).
70. Y. Bao *et al.*, The influenza virus resource at the National Center for Biotechnology Information. *J. Virol.* **82**, 596–601 (2008).
71. W. Li, A. Godzik, Cd-hit: A fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics* **22**, 1658–1659 (2006).
72. simplemaps, World cities database. <https://simplemaps.com/data/world-cities>. Accessed 20 May 2022.
73. M. N. Price, P. S. Dehal, A. P. Arkin, FastTree 2—Approximately maximum-likelihood trees for large alignments. *PLoS One* **5**, e9490 (2010).
74. W. M. Fitch, Toward defining the course of evolution: Minimum change for a specific tree topology. *Syst. Biol.* **20**, 406–416 (1971).
75. C. Xue, Z. Liu, N. Goldenfeld, Scale-invariant topology and bursty branching of evolutionary trees emerge from niche construction. *Proc. Natl. Acad. Sci. U.S.A.* **117**, 7879–7887 (2020).
76. Y. I. Wolf *et al.*, Origins and Evolution of the Global RNA Virome. *MBio* **9**, e02329-18 (2018).