# Analysis of partial sequences of the RNA-dependent RNA polymerase gene as a tool for genus and subgenus classification of coronaviruses

David A Wilkinson, Léa Joffrin, Camille Lebarbenchon, Patrick Mavingui

MICROBIOLOGY
SOCIETY

OPEN
MICROBIOLOGY

# Analysis of partial sequences of the RNA-dependent RNA polymerase gene as a tool for genus and subgenus classification of coronaviruses

David A. Wilkinson*, Léa Joffrin†, Camille Lebarbenchon and Patrick Mavingui

## Abstract

The recent reclassification of the *Riboviria*, and the introduction of multiple new taxonomic categories including both sub-families and subgenera for coronaviruses (family *Coronaviridae,* subfamily *Orthocoronavirinae*), represents a major shift in how official classifications are used to designate specific viral lineages. While the newly defined subgenera provide much-needed standardization for commonly cited viruses of public health importance, no method has been proposed for the assignment of subgenus based on partial sequence data, or for sequences that are divergent from the designated holotype reference genomes. Here, we describe the genetic variation of a 387 nt region of the coronavirus RNA-dependent RNA polymerase (RdRp), which is one of the most used partial sequence loci for both detection and classification of coronaviruses in molecular epidemiology. We infer Bayesian phylogenies from more than 7000 publicly available coronavirus sequences and examine clade groupings relative to all subgenus holotype sequences. Our phylogenetic analyses are largely coherent with whole-genome analyses based on designated holotype members for each subgenus. Distance measures between sequences form discrete clusters between taxa, offering logical threshold boundaries that can attribute subgenus or indicate sequences that are likely to belong to unclassified subgenera both accurately and robustly. We thus propose that partial RdRp sequence data of coronaviruses are sufficient for the attribution of subgenus-level taxonomic classifications and we supply the R package, MyCoV, which provides a method for attributing subgenus and assessing the reliability of the attribution.

## INTRODUCTION

Coronaviruses are widely studied for their impact on human and animal health [1], as well as their broad diversity and host/reservoir associations. In recent years, the emergence of betacoronaviruses in human populations has resulted in widespread morbidity and mortality. The severe acute respiratory syndrome (SARS) coronavirus was responsible for 8096 cases and 774 deaths during the 2002–2003 outbreak [World Health Organization (WHO), http://www.who.int/csr/sars/en]. Since 2013, the Middle East respiratory syndrome (MERS) coronavirus has infected 2506 people and led to 862 deaths (WHO data). At the time of writing, the pandemic caused by severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) [also defined as

nCoV-2019, causing the disease coronavirus disease 2019 (COVID-19)] is an ongoing public health emergency of international concern, having infected more than 16000000 people worldwide, resulting in more than 600000 deaths. Thanks to molecular epidemiology studies, we know that SARS, MERS and SARS-CoV-2 had their origins in wild animal reservoir species before spilling over into humans. Indeed, numerous molecular studies have identified a wealth of coronavirus diversity harboured by equally diverse animal hosts [1–9], and phylogenetic analysis of sequence data from these studies is helping to increase our understanding of many aspects of disease ecology and evolution [10–12]. This includes the role of reservoir hosts in disease maintenance and transmission [3, 4, 6, 13–15], the evolutionary origins of human-infecting coronaviruses

[3, 16], the importance of bats as reservoirs of novel coronaviruses [9, 17], the role of intermediate hosts in human disease emergence [3, 14, 18], and the understanding of risk that might be related to coronavirus diversity and distributions [19, 20].

The precise taxonomic classification of all organisms undergoes constant drift as new discoveries are made that inform their evolutionary histories, and as methods for studying evolutionary histories change and improve [21, 22]. However, in 2018, the International Committee for the Taxonomy of Viruses (ICTV, https://talk.ictvonline.org/) introduced a shift in the taxonomic designations of all RNA viruses, introducing the realm *Riboviria*, grouping 'all RNA viruses that use cognate RNA-dependent RNA polymerases (RdRps) for replication' [23]. In addition to this basal classification, many new taxonomic classifications were defined, or existing taxa reclassified. This included the separation of the family *Coronaviridae* into two subfamilies – the amphibian-infecting *Letovirinae* and *Orthocoronavirinae* encompassing the genera *Alphacoronavirus*, *Betacoronavirus*, *Gammacoronavirus* and *Deltacoronavirus*, which are classically recognized to infect mammals and birds and have significance in human and livestock diseases. The subgenus level of classification for members of the *Orthocoronavirinae* was also introduced, providing specific taxa for commonly cited groups of similar viruses such as *Betacoronavirus* lineages β-A, β-B, β-C and β-D [24], which were designated *Embecovirus*, *Sarbecovirus*, *Merbecovirus* and *Nobecovirus*, respectively [23]. The nomenclature for the designated subgenera was assigned with respect to known host species for each subgenus (e.g. rhinacoviruses for alphacoronviruses known to be hosted by bats of the family *Rhinolophidae*) or based on commonly used disease terminology (e.g. merbecoviruses for betacoronaviruses related to MERS coronavirus). Whole-genome data from holotype specimens selected to represent an exhaustive spectrum of coronavirus diversity were used to test the phylogenetic repartition and support for each of these taxa [25]. Due to the limited diversity of holotype specimens classified into these subgenera, there is currently no method for attributing subgenus to isolates with divergent sequences, and no proposed method for partial sequence data. However, sequence data from the RdRp region of the polymerase gene are one of the most commonly used tools for the purposes of coronavirus detection, identification and classification in molecular epidemiology [26].

Here, we examine the phylogenetic relationship from all identifiable public partial RdRp sequences of coronaviruses using Bayesian inference in BEAST 2 and examine the clade associations of all defined subgenus holotypes. We use this analysis to explore the range of logical similarity thresholds for the designation of subgenus-level classifications to partial RdRp sequence and predict subgenus classifications for all reference sequences. We cross-validate a sequence identity-based classification method against phylogenetically inferred classifications showing that alignment identity is >99% specific for the assignment of subgenus-level classifications to partial RdRp sequences. We compiled a database of our assigned classifications and developed the R package 'MyCoV' for assignment of user-generated sequences to these taxa.

## METHODS

### Sequence data, curation and alignment

Sequence data were obtained from the National Center for Biotechnology Information (NCBI) nucleotide database on the 5 July 2019, using the search term 'coronavir*'. This resulted in the identification of 30249 sequences. A preliminary set of representative partial RdRp sequences was compiled with reference to recent publications describing coronavirus diversity across the subfamily *Orthocoronavirinae* [27], in order to include starting reference sequences with the largest possible diversity of coronaviruses. This preliminary list was then used to identify partial RdRp sequences from retrieved NCBI records by annotating regions that had at least 70% identity to any reference sequence in the Geneious software package (version 9.4.1). Annotated regions and 200 bp of flanking sequence data were then extracted. Data containing incomplete sequences in the form of strings of Ns or significant numbers of ambiguities (>5) were removed. Open reading frames with a minimum length of 300 bp were identified and extracted from the remaining sequences. In the case where the correct reading frame was ambiguous, pairwise alignment to reference sequence data was used to determine reading frame. The remaining sequences were then aligned in-frame using the multiple alignment program for amino acid or nucleotide sequences (MAFFT) [28], and the resulting alignment was further curated by visual inspection. Retained sequences were then trimmed to include only the most frequently sequenced partial region of RdRp and so that each sequence contained a minimum of 300 gap-free bases. The final alignment was 387 bp in length with 7544 individual sequences, of which 3155 were unique. The relevant 387 bp region is positioned within the second half of the RNA-binding non-structural protein 12 (RdRp) region of ORF1ab, corresponding to nucleotide positions 15287:15673 in *Merbecovirus* (*Betacoronavirus*) holotype reference sequence JX869059.2, positions 14456:14842 in *Pedacovirus* (*Alphacoronavirus*) holotype AF353511.1, positions 12725:13108 in *Andecovirus* (*Deltacoronavirus*) holotype reference sequence JQ065048.1 and positions 14221:14607 in *Cegacovirus* (*Gammacoronavirus*) holotype EU111742.1.

### Genetic analyses

Phylogenies were inferred from all unique sequences using the BEAST 2 software [29]. Parameters were estimated for a GTR substitution model with four gamma categories and an estimated proportion of invariant sites. The Yule population model was used, and a log-normal distribution of priors was specified for birth rate and proportion of invariant site priors. Convergence of estimated parameters was assessed
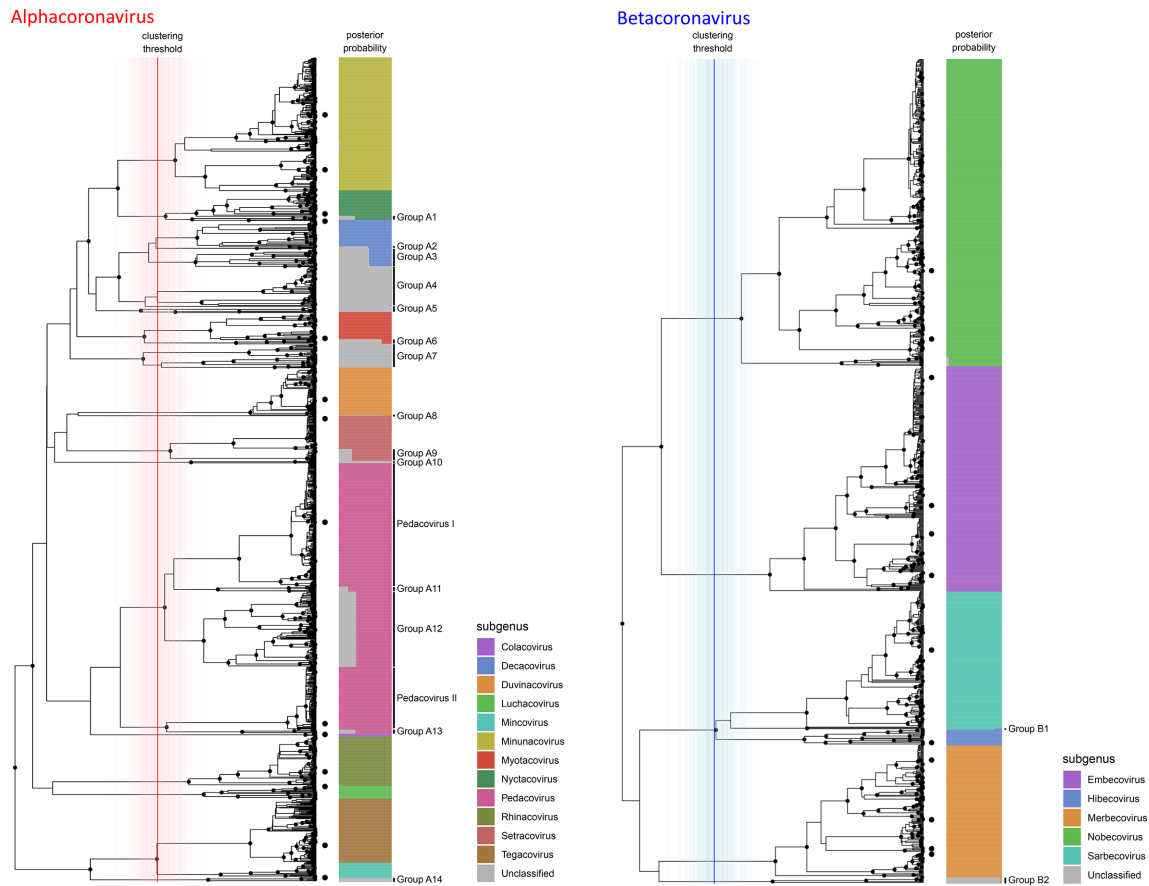
**Fig. 1.** Phylogenetic subgenus classifications for partial RdRp sequences of alphacoronaviruses (left) and betacoronaviruses (right). Depicted trees are subtrees of the consensus phylogeny presented in Fig. 4. Dots on leaf tips indicate sequences belonging to holotype reference sequences for each subgenus. Coloured bars show the proportion of trees from the Bayesian analysis where the corresponding leaf was assigned to each subgenus, which is indicated by colour according to the legend. Vertical lines show the distribution of cluster-defining height thresholds that were identified to assign subgenus classifications, with the median of all clustering thresholds displayed in bold; lines are coloured by genus as in Fig. 4. Monophyletic groups where all members have posterior probabilities of being assigned to a known subgenus of <90 % are highlighted and assigned sequential IDs.

in Tracer v1.7.1 [30]. Three independent MCMC chains were run until effective sample sizes were above 200 for all estimated parameters after removing MCMC chains from the start of the analysis prior to convergence (burn-in). Analyses were run until convergence criteria could be fulfilled whilst providing equal chain lengths after burn-in for all three repeats, meaning that the number of trees in the posterior distributions was the same for each independent repeat.

Genetic distance measures were calculated using the ape package [31] in RStudio as the proportion of variant sites in pairwise comparisons after removing regions containing gaps in either compared sequence.

### Taxonomic classification

Sequences originating from known references were used to identify common ancestral nodes for the *Orthocoronaviridae* genera within each phylogenetic tree. Genus-level subtrees

were then extracted and treated independently for subgenus-level analyses.

Sequences originating from defining subgenus holotype samples were identified in the genus-level topologies. Clustering thresholds were defined as the highest node positions at which clusters of leaves could be defined without combining holotype specimens from different subgenera into the same clade. Clusters defined at these thresholds that contained no holotype specimens were designated as 'unclassified'. Clustering thresholds were calculated, and subgenera were assigned to all sequences across a random subsample of 453 trees, 151 from each independent repeat of the phylogenetic analyses. The proportion of trees in which each sequence was assigned to a given subgenus was used as the 'posterior probability' of that sequence belonging to that subgenus. Sequences with <90% majority posterior probabilities were designated as 'unclassified'. Potential positioning of new subgenus level clades (as indicated by 'Group *X*' in Figs 1 and 2) was inferred
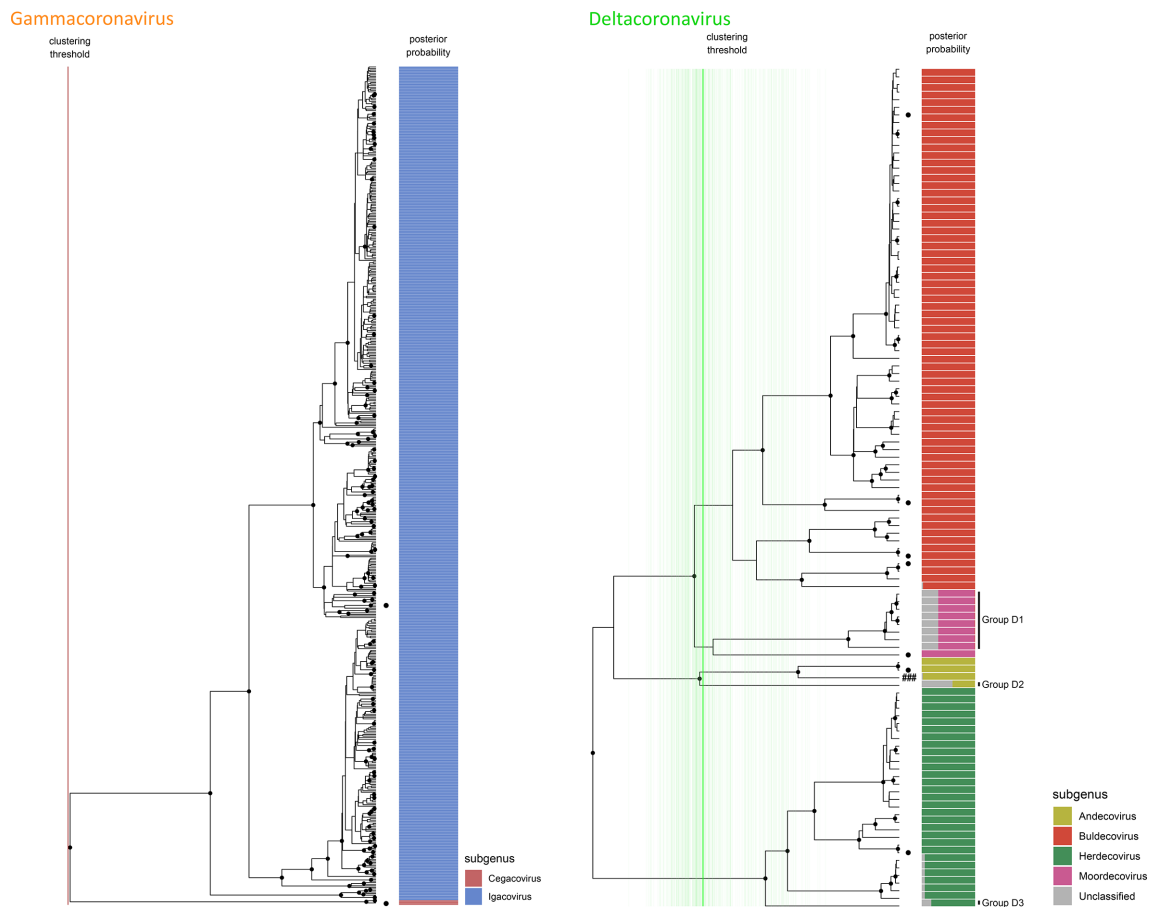
**Fig. 2.** Phylogenetic subgenus classifications for partial RdRp sequences of gammacoronaviruses (left) and deltacoronaviruses (right). Depicted trees are subtrees of the consensus phylogeny presented in Fig. 4. Dots on leaf tips indicate sequences belonging to holotype reference sequences for each subgenus. Coloured bars show the proportion of trees from the Bayesian analysis where the corresponding leaf was assigned to each subgenus, which is indicated by colour according to the legend. Vertical lines show the distribution of cluster-defining height thresholds that were identified to assign subgenus classifications, with the median of all clustering thresholds displayed in bold; lines are coloured by genus as in Fig. 4. In the case of gamma coronaviruses, both subgenera are consistently separated at the root of the tree, thus all cluster defining heights equate to the root. Monophyletic groups where all members have posterior probabilities of being assigned to a known subgenus of <90% are highlighted and assigned sequential IDs.

using the maximum clade-credibility consensus tree from all BEAST analyses, identifying monophyletic clades where all descendants were not classified into defined subgenera.

## Cross-validation

The assignment of sequences to the relevant subgenus using best hit and pairwise identity data from BLASTN [32] was tested by iteratively removing each sequence from the test database and reassigning its classification. Sequences that could not be reassigned to the same subgenus by this method were reclassified as 'atypical' members of their respective subgenera.

## Recombination analysis

To explore the extent to which ancestral recombination events confound our ability to classify coronaviruses at the subgenus level from a single gene locus, sequence accession numbers

were identified from the pool of 7544 sequences analysed above that possessed data for both RdRp and structural spike gene regions. This search resulted in the identification of 2649 sequences. Due to sequence heterogeneity in the spike region, sequences were grouped and compared by genus based on the predicted classification obtained using the RdRp locus. The full-length spike gene was extracted, aligned in-frame using MAFFT, and maximum-likelihood phylogenies were generated for each genus using MEGA X removing all gap-containing sites. Tanglegram representations were used to compare tree topologies obtained using data from the two different loci. Differences in tree topology that resulted in paraphyly of RdRp-predicted subgenus classifications within the spike gene tree were deemed to be the result of recombination. We did not use this analysis to examine within-subgenus recombination events, but focused on those cases where recombination would result in different subgenus

classifications for a given virus, depending on which genomic locus was used for its classification.

## R package for assignment of user-generated sequences

The purpose of the R package MyCoV is to allow users to classify coronavirus sequence data that include the relevant portion of the RdRp gene to the taxonomic level of subgenus, and to assess to what extent the classification is optimally based on the criteria presented herein.

In order to achieve this, the 3155 unique partial sequences from the phylogenetic analyses were used to establish a reference BLAST database. Metadata pertaining to host organism, country of origin and date of collection were mined from NCBI and standardized by taxonomic grouping of the host and geographical region of origin to generate corresponding metadata for all 7544 NCBI reference sequences from which the unique sequence list was established.

MyCoV was written as a basic wrapper script for BLASTN, which queries sequences of interest against the established database, and summarizes subgenus classification, subgenus posterior support of the most similar sequence in the phylogenetic analysis, pairwise distances to the most similar sequence in the database and their metadata using R packages 'ggplot2', 'formattable' (available at https://github.com/renkun-ken/formattable) and 'ggtree' [33].

As the MyCoV database was established prior to the recent emergence of the SARS-CoV-2, we used genomic sequence data from this virus as a test case for the utility of the MyCoV package. Outputs from this analysis are shown in Fig. 3.

MyCoV is available at https://github.com/dw974/MyCoV.

## RESULTS

Our three independent, randomly seeded phylogenetic analyses converged on similar estimates for all parameters in BEAST 2. The resulting predictions of tree topology had well-supported major nodes with narrow posterior distributions around most node heights (Fig. 4a). The four known genera associated with these sequences fell into four well-supported clades, divided close to the root of the tree. Genetic distance measures between all members of the four genera had logical thresholds for the distinction between genera except in the case of some betacoronaviruses, which had major clade divisions close to the root of the tree (Fig. 4b) and therefore had genetic distances between members of the same genus that overlapped with distances between members of the alpha- and betacoronaviruses. In practice, this is likely to mean that identity-based phylogenetic topologies based on this partial region of RdRp may incorrectly infer paraphyly between members of the alpha- and betacoronaviruses.

At the subgenus level, separation of the inferred tree topologies into monophyletic clades based on the positions of

reference holotype sequences produced logical and well-supported groupings that covered the majority of coronavirus diversity explored to date by RdRp sequencing (Figs 1 and 2). In total, 88% of unique sequences fell into clade groups containing subgenus holotypes with subgenus assignment posterior probabilities of >90%. The remaining 12% of unique sequences fell into 19 separate monophyletic groups, of which 14 were alphacoronaviruses (Fig. 1), 2 were betacoronaviruses (Fig. 1) and 3 were deltacoronaviruses (Fig. 2). When host and geographical origins of isolates falling within unclassified clades were examined, the majority were associated with regional radiations for which little or no genomic or phenotypic data are available. For example, unclassified deltacoronaviruses were all from bird species in Oceania, and many unclassified alpha- and betacoronaviruses originated in bat species that are exclusively found in Central and South America (Fig. S1, available in the online version of this article).

The pedacoviruses, for which multiple genome holotypes were supplied for the description of subgenus, were split into multiple clades by the imposition of a common height threshold for cluster definition using the presented methodology. The two holotype-containing clades corresponded to a single group of porcine epidemic diarrhoea virus (PEDV) – related viruses distributed globally but entirely from pigs (Pedacovirus I in Fig. 1) and a monophyletic group of viral sequences obtained from Asian *Scotophilus* bats. The monophyletic group that contained both *Pedacovirus* holotypes also enclosed other major viral clades (clades A11, A12 and A13 in Fig. 1), which were mainly associated with other bat species of the family *Vespertillionidae* (Fig. S1).

Cross-validation of sub-genus assignments by best hit using BLASTN was successful in more than 99.9% of cases, with a handful of lone sequences that branched at basal positions of each phylogenetic clade group being assigned to different subgenera.

For sequence members of each genus, genetic distance measurements between and within sequences attributed to each subgenus displayed logical and discrete threshold boundaries for the distinction of individual subgenus members. The one exception, again, was members of the subgenus *Pedacovirus*, which displayed overlapping within-taxon distances with between-taxon distances for other *Alphacoronavirus* subgenera (Fig. 5). The distinct *Pedacovirus* clades displayed in Fig. 1 were thus treated as separate subgenera for distance threshold calculation. Optimal thresholds were identified as the midpoint of a fitted binomial probability distribution for intra- and inter- subgenus pairwise distances. The optimal identity thresholds for distinguishing same vs different subgenera were as follows: (i) 77.6% identity, resulting in 99.7% precision and 95.3% accuracy of classification for subgenera of the alphacoronaviruses; (ii) 71.7% identity, resulting in 99.9% precision and 99.6% accuracy of classification for subgenera of the betacoronaviruses; (iii) 74.9% identity, resulting in 98.8% precision and 99.2% accuracy of classification for subgenera of the deltacoronaviruses; and
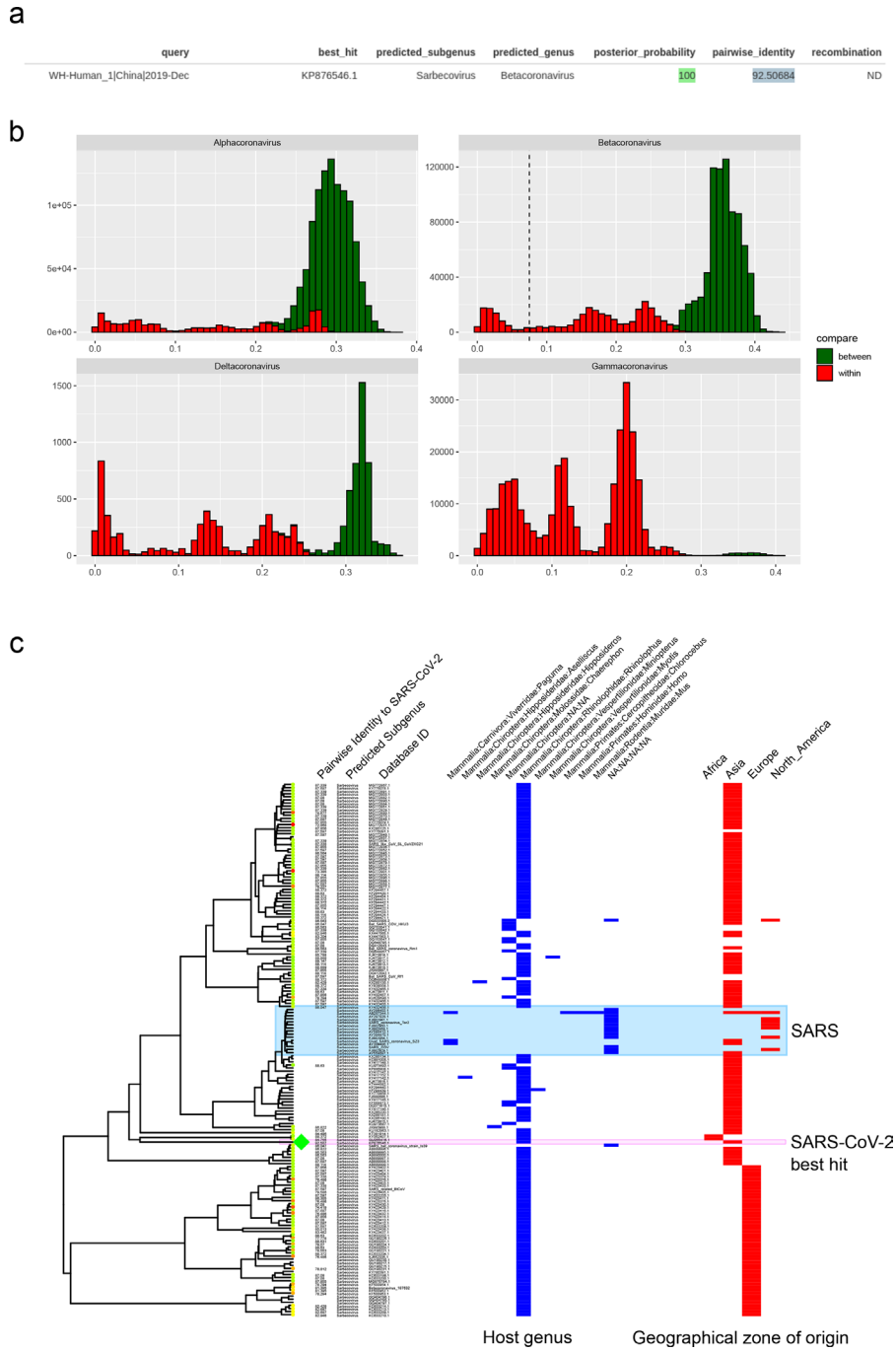
**Fig. 3.** MyCoV output plots for the analysis of the 2019nCoV. (a) Tabular output of best blastn hit of the query sequence to the reference database. The predicted subgenus and genus of the best-matching hit are displayed, as well as the posterior support for the assignment to the predicted subgenus (see the Methods section). Pairwise identity between the two sequences is shown and is calculated relative to the maximum possible alignment length against the reference sequences (387 bp). (b) For each queried sequence, pairwise identity values are mapped to all observations from pairwise comparisons between sequences in the database. The vertical dashed line represents the pairwise dissimilarity of the queried sequence. (c) Phylogenetic positioning and metadata from the analysis of the reference sequences are displayed. Reference sequences with blast hits matching the queried sequence are highlighted on the leaves, and tips are coloured from red to green with increasing pairwise identity. The hit with the best score is highlighted by a large green diamond on the tip. Pairwise identity scores are displayed for all leaves, as well as predicted subgenus. Host genus associations (blue) and geographical region of origin (red) from available metadata are indicated by binary heatmaps. Note that multiple metadata observations are possible for each leaf, as leaves are displayed for unique sequences only. The ID next to each leaf is that of the representative sequence for that leaf, and other IDs are left off for clarity.
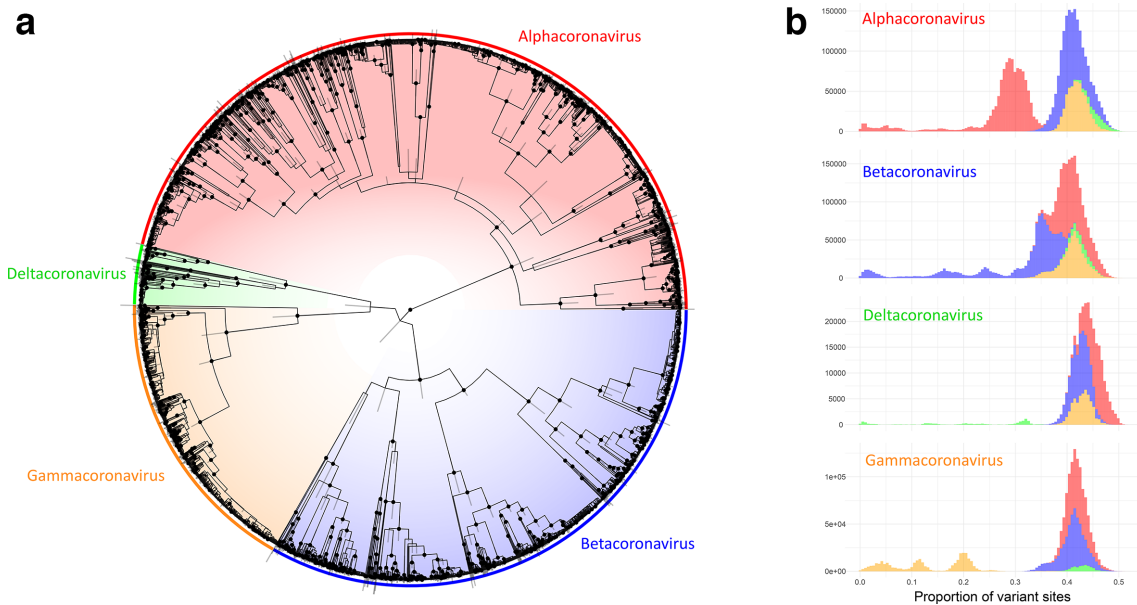
**Fig. 4.** Comparative analysis of 3155 partial RdRp sequences belonging to members of the *Orthocoronavirinae*. (a) Consensus phylogeny from three independent beastBEAST analyses. Nodes with posterior support >90% are highlighted with dots and bars display the 95% highest posterior distributionHighest Posterior Distribution of the heights of each node. Colours indicate genus-level classification for sequences, clades and pairwise comparisons throughout. (b) Histograms of genetic distances, measured as the proportion of variant sites, between sequences belonging to each genus and grouped by the genus of the queried sequence.
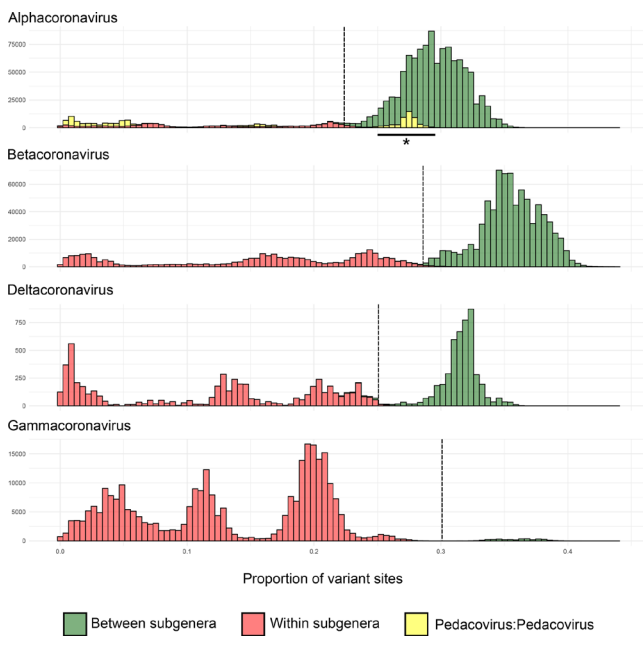


**Fig. 5.** Histograms of genetic distances for intra- and inter-subgenus comparisons. Vertical dashed lines represent the optimal genetic distance cut-offs for the subgenus threshold, calculated as the midpoint of the fitted binomial probability distribution. *Pedacovirus*: *Pedacovirus* distance comparisons that fall outside of the optimal distance threshold are highlighted by an asterisk.

(iv) 69.9% identity, resulting in 100% precision and 100% accuracy of classification for subgenera of the gammacoronaviruses (Figs 1, 2 and 5).

Comparing tree topologies inferred from both the partial RdRp region and the spike gene allowed us to estimate whether the predicted subgenus classifications were robust to the choice of genomic locus used for their classification. Out of the 2649 sequences available that had both RdRp and spike data available, only members of two subgenera from the genus *Alphacoronavirus* showed evidence of recombination. Assuming that the differences in tree topology are explained by the most parsimonious route, we predict only 2 independent inter-subgenus recombination events resulting in 10 of 2649 isolates that would have been assigned to different subgenera if using the spike gene compared to when using the partial RdRp region. These likely recombination events occurred between members of the tegacoviruses and pedacoviruses, in viruses associated with pig hosts as well as members of the minunacoviruses and myotacoviruses, in viruses associated with fruit bat hosts (Fig. S1). Tanglegram tree topologies assessing tree paraphyly for alphacoronaviruses, betacoronaviruses, gammacoronaviruses and deltacoronaviruses are displayed in Figs S2–S5, respectively.

Our R package, MyCoV, successfully identified SARS-CoV-2 as a member of the subgenus *Sarbecovirus*, with the closest match being to reference sequence KP876545.11 (*Rhinolophus* bat coronavirus BtCov/3990), which showed 92.5% pairwise identity to SARS-CoV-2 in the RdRp region.

This sequence had been assigned 100% posterior support for being attributed to the subgenus *Sarbecovirus* (Fig. 3a). Distributions of pairwise identities within members of the subgenera of the betacoronaviruses fell between 71 and 100%, whereas pairwise distances between *Betacoronavirus* subgenera were less than or equal to 71%. Thus, the output of MyCoV allows us to state with certainty that SARS-CoV-2 belongs to this subgenus (Fig. 3b). Positioning of the closest match in the phylogenetic tree shows that the SARS-CoV-2 forms a distinct lineage from SARS coronavirus, and that its closest match belonged to a *Rhinolophus* bat from China (Fig. 3c). Interestingly, this sequence came from an abandoned mine in 2013, suggesting that SARS-CoV-2 predecessors circulated in bat communities for a number of years prior to the 2019 emergence in human populations. The provided visualization of host and geographical origins for these partial reference sequences allows for a rapid assessment of the distribution of similar viruses; for example, it highlights the fact that SARS-related and SARS-CoV-2-related viruses have also been identified in bats in Africa (specifically *Rhinolophus* bats in Kenya), and that they are not just restricted to Asian bat hosts.

## DISCUSSION

The recent reclassification of the *Riboviria* is a logical progression in viral taxonomy, as the unique mechanism of replication of all negative-sense, single-stranded, RNA viruses results in the conservation of many viral characteristics, including relative sequence conservation of regions of the cognate RNA-dependent RNA polymerase. Consequently, such genomic loci lend themselves to the design of primers for virus detection in diagnostic and molecular epidemiology, and to the phylogenetic inference of evolutionary histories. Furthermore, establishing the classification level of subgenus has provided a useful tool for researchers, attributing standardized terminology for many commonly referenced viral lineages that, in general, demonstrate a level of specificity in their host associations and epidemiological characteristics (Fig. S1).

Our analyses have shown that the phylogenetic interpretation of short sequences of the RdRp locus of members of the *Orthocoronaviridae* is largely coherent with genome-scale analyses based on designated holotype members for each subgenus. The vast majority of known RdRp sequences (88%) can be classified into the defined subgenera, and their classification cross-validated based on simple distance thresholds established from a 387 bp fragment of RdRp. Globally, these distance measures form discrete clusters between taxa, offering logical threshold boundaries that can attribute subgenus or indicate sequences that are sufficiently distinct in their sequences so that no subgenus can be assigned both accurately and robustly without the need for complex phylogenetic inference. The provided R package, MyCoV, provides a method for achieving this and for the assessment of the reliability of the attribution.

An alternative strategy for coronavirus classification from partial sequence data may be using the spike protein-encoding S-gene, which our analyses demonstrate would provide identical subgenus classifications in the vast majority of cases. However, the use of this region is more common in epidemic outbreak scenarios and thus there are many S gene sequences in public databases that are either identical or extremely closely related. Performing comparative sequence searches by querying the NCBI nucleotide database with the two search terms '((coronavir* spike) OR (coronavir* S gene)) AND 'viruses'[porgn:__txid10239]' and '((coronavir* RdRp) OR (coronavir* polymerase)) AND 'viruses'[porgn:__txid10239]' shows that there are approximately three times more sequences from the S gene, but that these sequences originate from approximately three times fewer viral taxa. Also, unlike the highly conserved RdRp region, the spike region is subject to large-scale structural variation, meaning that global alignment of spike sequences is challenging (here, we restricted our analyses of spike sequences to within-genus comparisons). We therefore favour the use of the RdRp, which provides a more exhaustive representation of known coronavirus diversity whilst allowing global sequence comparisons.

Of course, this form of interpretation is subject to the same caveats as any other that is based on partial sequence data from a short, single genomic locus; indeed, the effects of potential recombination events [17, 34–36] cannot be captured, and some uncertainties will exist in the presented phylogenetic trajectories that may be resolvable by the addition of longer sequence data. For these reasons, we do not suggest the definition of new subgenera for unclassified clade groups presented in Figs 1 and 2. The limits of the phylogenetic resolving power of this partial region of RdRp are most clear for members of the genus *Alphacoronavirus*, where there is an elevated level of mid-distance genetic diversity and a large number of unclassified genetic clade groups associated with regional, likely host-specific radiations. And thus, precise taxonomic delineation of emerging alphacoronaviruses will require more information than is offered by this RdRp locus. Conversely, the clear genetic distinction and corresponding epidemiological associations that exist between clade groups of the pedacoviruses raise the question as to whether the definition of this subgenus should be revisited.

References

1. Lau SKP, Chan JFW. Coronaviruses: emerging and re-emerging pathogens in humans and animals. *Virol J* 2015;12:209.

2. Chan JFW, Lau SKP, To KKW, Cheng VCC, Woo PCY *et al*. Middle East respiratory syndrome coronavirus: another zoonotic beta-coronavirus causing SARS-like disease. *Clin Microbiol Rev* 2015;28:465–522.

3. Cui J, Li F, Shi Z-L. Origin and evolution of pathogenic coronaviruses. *Nat Rev Microbiol* 2019;17:181–192.

4. Woo PCY, Lau SKP, Wernery U, Wong EYM, Tsang AKL *et al*. Novel betacoronavirus in dromedaries of the middle East, 2013. *Emerg Infect Dis* 2014;20:560–572.

5. Mihindukulasuriya KA, Wu G, St Leger J, Nordhausen RW, Wang D. Identification of a novel coronavirus from a Beluga whale by using a panviral microarray. *J Virol* 2008;82:5084–5088.

6. Razanajatovo NH, Nomenjanahary LA, Wilkinson DA, Razafimanahaka JH, Goodman SM *et al*. Detection of new genetic variants of Betacoronaviruses in endemic Frugivorous bats of Madagascar. *Virol J* 2015;12:42.

7. Zhang T, Wu Q, Zhang Z. Probable Pangolin origin of SARS-CoV-2 associated with the COVID-19 outbreak. *Curr Biol* 2020.

8. Lam TT-Y, Jia N, Zhang Y-W, Shum MH-H, Jiang J-F *et al*. Identifying SARS-CoV-2-related coronaviruses in Malayan pangolins. *Nature* 2020;583:282-285.

9. Li W, Shi Z, Yu M, Ren W, Smith C *et al*. Bats are natural reservoirs of SARS-like coronaviruses. *Science* 2005;310:676–679.

10. Drexler JF, Corman VM, Drosten C. Ecology, evolution and classification of bat coronaviruses in the aftermath of SARS. *Antiviral Res* 2014;101:45–56.

11. Corman VM, Baldwin HJ, Tateno AF, Zerbinati RM, Annan A *et al*. Evidence for an ancestral association of human coronavirus 229E with bats. *J Virol* 2015;89:11858–11870.

12. Anthony SJ, Gilardi K, Menachery VD, Goldstein T, Ssebide B *et al*. Further evidence for bats as the evolutionary source of middle East respiratory syndrome coronavirus. *mBio* 2017;8:e00373-17.

13. Corman VM, Muth D, Niemeyer D, Drosten C. *Hosts and Sources of Endemic Human Coronaviruses. In: Advances in Virus Research*. pp. 163–188.

14. Menachery VD, Graham RL, Baric RS. Jumping species-a mechanism for coronavirus persistence and survival. *Curr Opin Virol* 2017;23:1–7.

15. Song H-D, Tu C-C, Zhang G-W, Wang S-Y, Zheng K *et al*. Cross-host evolution of severe acute respiratory syndrome coronavirus in palm civet and human. *Proc Natl Acad Sci U S A* 2005;102:2430–2435.

16. Forni D, Cagliani R, Clerici M, Sironi M. Molecular evolution of human coronavirus genomes. *Trends Microbiol* 2017;25:35–48.

17. Tao Y, Shi M, Chommanard C, Queen K, Zhang J *et al*. Surveillance of bat coronaviruses in Kenya identifies relatives of human coronaviruses NL63 and 229E and their recombination history. *J Virol* 2017;91.

18. Joffrin L, Dietrich M, Mavingui P, Lebarbenchon C. Bat pathogens hit the road: but which one? *PLoS Pathog* 2018;14:e1007134.

19. Anthony SJ, Johnson CK, Greig DJ, Kramer S, Che X *et al*. Global patterns in coronavirus diversity. *Virus Evol* 2017;3:vex012.

20. Han BA, Kramer AM, Drake JM. Global patterns of zoonotic disease in mammals. *Trends Parasitol* 2016;32:565–577.

21. King AMQ, Lefkowitz EJ, Mushegian AR, Adams MJ, Dutilh BE *et al*. Changes to taxonomy and the International Code of virus classification and nomenclature ratified by the International Committee on taxonomy of viruses (2018). *Arch Virol* 2018;163:2601–.

22. Simmonds P, Adams MJ, Benkő M, Breitbart M, Brister JR *et al*. Consensus statement: virus taxonomy in the age of metagenomics. *Nat Rev Microbiol* 2017;15:161-168.

23. Walker PJ, Siddell SG, Lefkowitz EJ, Mushegian AR, Dempsey DM *et al*. Changes to virus taxonomy and the International Code of virus classification and nomenclature ratified by the International Committee on taxonomy of viruses (2019). *Arch Virol* 2019;164:2417–2429.

24. Woo PCY, Huang Y, Lau SKP, Yuen K-Y. Coronavirus genomics and bioinformatics analysis. *Viruses* 2010;2:1804-20.

25. Ziebuhr J, Baric RS, Baker S, de Groot RJ, Drosten C *et al*. ICTV Report 2017;013S.

26. Lehmann KC, Gulyaeva A, Zevenhoven-Dobbe JC, Janssen GMC, Ruben M *et al*. Discovery of an essential nucleotidylating activity associated with a newly delineated conserved domain in the RNA polymerase-containing protein of all nidoviruses. *Nucleic Acids Res* 2015;43:8416–.

27. Joffrin L, Goodman SM, Wilkinson DA, Ramasindrazana B, Lagadec E *et al*. Bat coronavirus phylogeography in the Western Indian Ocean. *Sci Rep* 2020;10:742866.

28. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol Biol Evol* 2013;30:772–.

29. Bouckaert R, Heled J, Kühnert D, Vaughan T, Wu C-H *et al*. Beast 2: a software platform for Bayesian evolutionary analysis. *PLoS Comput Biol* 2014;10:e1003537.

30. Rambaut A, Drummond AJ, Xie D, Baele G, Suchard MA. Posterior summarization in Bayesian phylogenetics using tracer 1.7. *Syst Biol* 2018;67:901–904.

31. Paradis E, Schliep K. Ape 5.0: an environment for modern phylogenetics and evolutionary analyses in R. *Bioinformatics* 2019;35:526–528.

32. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J *et al*. BLAST+: architecture and applications. *BMC Bioinformatics* 2009;10:421.

33. Yu G, Smith DK, Zhu H, Guan Y, Lam TT. ggtree: an r package for visualization and annotation of phylogenetic trees with their covariates and other associated data. *Methods Ecol Evol* 2017;8:28–36.

34. Li X, Giorgi EE, Marichann MH, Foley B, Xiao C *et al*. Emergence of SARS-CoV-2 through recombination and strong purifying selection. *bioRxiv* 2020

35. Jackwood MW, Boynton TO, Hilt DA, McKinley ET, Kissinger JC *et al*. Emergence of a group 3 coronavirus through recombination. *Virology* 2010;398:98–.

36. Woo PCY, Lau SKP, Yip CCY, Huang Y, Tsoi H-W *et al*. Comparative analysis of 22 coronavirus HKU1 genomes reveals a novel genotype and evidence of natural recombination in coronavirus HKU1. *J Virol* 2006;80:7136–7145.