**Title**:

**Molecular features similarities between SARS-CoV-2, SARS, MERS and key human genes could favour the viral infections and trigger collateral effects.**

**Authors**:**Lucas L. Maldonado[1] and Laura Kamenetzky[1]**

**Affiliations:**[1]IMPaM, CONICET, Facultad de Medicina, Universidad de Buenos Aires, Ciudad Autónoma de Buenos Aires, Argentina. [2]

**e-mails:**lucas.l.maldonado@gmail.com; lkamenetzky@fmed.uba.ar

## Abstract

In December 2019 rising pneumonia cases caused by a novel β-coronavirus (SARS-CoV-2) occurred in Wuhan, China, which has rapidly spread worldwide causing thousands of deaths. The WHO declared the SARS-CoV-2 outbreak as a public health emergency of international concern therefore several scientists are dedicated to the study of the new virus. Since human viruses have codon usage biases that match highly expressed proteins in the tissues they infect and depend on host cell machinery for replication and co-evolution, we selected the genes that are highly expressed in the tissue of human lungs to perform computational studies that permit to compare their molecular features with SARS, SARS-CoV-2 and MERS genes. In our studies, we analysed 91 molecular features for 339 viral genes and 463 human genes that consisted of 677873 codon positions. Hereby, we found that A/T bias in viral genes could propitiate the viral infection favoured by a host dependant specialization using the host cell machinery of only some genes. The envelope protein E, the membrane glycoprotein M and ORF7 could have been further benefited by a high rate of A/T in the third codon position. Thereby, the mistranslation or de-regulation of protein synthesis could produce collateral effects, as a consequence of viral occupancy of the host translation machinery due tomolecular similarities with viral genes. Furthermore, we provided a list of candidate human genes whose molecular features match those of SARS-CoV-2, SARSand MERS genes, which should be considered to be incorporated into genetic population studies to evaluate thesusceptibility to respiratory viral infections caused by these viruses.The results presented here, settle the basis for further research in the field of human genetics associated with the new viral infection, COVID-19, caused by SARS-CoV-2 and for the development of antiviral preventive methods.

33 **Keywords**

34 SARS-CoV-2, SARS, MERS, codon usage, viral genes, human infection, COVID-19

35 **1. Introduction**

36 Since its initial outbreak at Huanan Seafood Wholesale Market in Wuhan, China, in late

37 2019, COVID-19 has affected more than 4 million people and caused more than 300

38 thousand deaths all around the world. Thereafter, scientists are focused not only on

39 studying the biology and dissemination of COVID-19 to control the transmission and

40 design proper diagnostic tools and treatments, but also theyare racing to design a

41 vaccine that could prevent the infection caused by the coronavirus SARS-CoV-2.This

42 virus belongs to the Betacoronavirus(β-coronavirus) of the *Coronaviridae*family, which

43 is also composed of three more genera:

44 Alphacoronavirus(αCoV),Gammacoronavirus(γCoV) andDeltacoronavirus(δCoV)

45 (Chen et al., 2020a). Viruses from this family possess a single-stranded, positive-sense

46 RNA and thegenome ranges from 26 to 32 kb (Su et al., 2016).

47 Coronaviruses have been identified in several host species including humans, bats,

48 civets, mice, dogs, cats, cows and camels (Cavanagh, 2007; Clark, 1993; Wang et al.,

49 2006; Zhou et al., 2018). Since severe acute respiratory syndrome (SARS), caused by

50 the coronavirus SARS-CoV, emerged in southern China in 2002(Peiris et al., 2004),

51 several studies tracing the transmission and possible reservoirs for viruses have been

52 performed. In early 2007, it had already been warned that bats were a natural reservoir

53 for an increasing number of emerging zoonotic virusesas well as for a large number of

54 viruses that have a close genetic relationship with the coronavirusesthat causethe severe

55 acute respiratory syndrome. Furthermore, it was warned that these viruses possess more

56 risk than other pathogens for disease emergence in human and domestic mammals

57 because of their higher mutation rates(Wang et al., 2006).Moreover, legal and illegal

58 trading of wildlife animals propitiates the environment for cross-species virus

59 transmission contributing to the rapid spread of the viral infections around the

60 world(Wang et al., 2006; Wong et al., 2019) as already occurred with SARS and the

61 Middle East respiratory syndrome (MERS) (Zaki et al., 2012). Both viruses have likely

62 originated in bats and are genetically diverse coronaviruses (Cui et al., 2019). Currently,

63 the outbreak of an atypical pneumonia caused by the novel coronavirus SARS-CoV-2

64 appears to have also startedfrom a zoonotic and a cross-species virus transmission at a

2

65  market in Wuhan including bats and pangolins,where animals were kept together and

66  the meat was sold(Chan et al., 2020).

67  Due to viruses replicate exclusively inside of living cells and depend exclusively

68  on the protein synthesis and chaperones machinery of their host, we speculated that the

69  primary structure of viral genes might be determined by the same forces that shape the

70  codon usage in the hosts' genes. Thereby, viral molecular patterns and codon usage

71  preferences would be a reflection of the host machinery. Codon pair bias and

72  dinucleotide preferences of viruses have been suggestedas the main factors that reflect

73  the codon usage of their hosts. Indeed, virus attenuation by codon pair deoptimization is

74  used as an efficacious attenuation method of various small RNA viruses and has

75  resulted in the generation of superior experimental live virus vaccines (Coleman et al.,

76  2008; Mueller et al., 2010; Nouën et al., 2014; Shen et al., 2015; Wang et al., 2015;

77  Yang et al., 2013).

78  In order to contribute to solving the sanitary emergence, here we provide a thorough and

79  comprehensive analysis that could help to understand the viability of the virus as well as

80  the susceptibility of the human host to the viral infection based on the molecular

81  patterns of their genes. Therefore, the main goals of ourwork wereto study the

82  molecular and evolutionary aspects of the human coronaviruses SARS-CoV-2,SARS

83  and MERS andto determine the level of similarity of the codon usage and molecular

84  features betweenthe genes of human coronaviruses and the human genesin order to

85  identify the factors that are responsible for the codons selection in the viruses.Moreover,

86  we proposed to identify the essential viral genes for viral replication andhumangenes

87  whosetranslation machinery is involved in propitiating the system for viral replicationin

88  orderto determine whether the genetic population variability could be involved in

89  modelling the gene features andtherefore contributing to the human susceptibility to

90  viral infections.

91  **2. Methods**

92  Up to late April, a total of☐500 SARS-CoV-2 β-coronavirus genome became available.

93  The total available sequences of β-coronavirus were downloaded from the NCBI

94  (https://www.ncbi.nlm.nih.gov/labs/virus/vssi/#/) including the reference genomes of

95  MERS (NC_019843), SARS (NC_004718) and SARS-CoV-2 (NC_045512) and were

96  classified according to their host. Different SARS-CoV-2 isolates from different

97  countries were pre-analysed but only reference genomes were retaineddue to the low

98    variability of the data.The genomes qualitywas assessed and the genomes containing

99    more than 10 gaps were discarded. CDS of representative viruses fromthe previous

100    classification were selected and analysed. Since human viruses have codon usage biases

101    (CUB) that match highly expressed proteins in the tissues they infect (Miller et al.,

102    2017) we selected 463 highly expressed human genes in lungs tissues according to the

103    fold-change between the expression level in lung and the tissue with second-highest

104    expression level according to the"Human Protein Atlas"

105    (https://www.proteinatlas.org/humanproteome/tissue/lung). We considered valid CDS

106    when they started with an ATG codon, ended with an in-frame stop codon, and had no

107    undetermined nucleotides nor internal stop codons. The accession numbers of the

108    sequences that were used here can be found in Supplementary file 1.

109    The CUB analyses were performed with CodonW 1.4.4 (J Peden,

110    http://codonw.sourceforge.net/). The total GC content of the CDS as well as the GC

111    content of the first (P1), second (P2), and third (P3) codon positions were calculated

112    using custom PERL scripts. To correct the inequality composition at the third codon

113    position (Sueoka, 1988), the three stop codons (UAA, UAG, and UGA) were excluded

114    in the calculation of P3, and the two single codons for methionine (AUG) and

115    tryptophan (UGG) were excluded from P1, P2, and P3.

## 2.1. Codon usage indices

117    The following codon indices were calculated: relative synonymous codon usage

118    (RSCU) (Sharp and Li, 1987), the effective number of codons (ENc) (Wright, 1990),

119    codon adaptation index (CAI) (Lee et al., 2010; Sharp and Li, 1987), codon bias index

120    (CBI) (Bennetzen and Hall, 1982), the optimal frequency of codons (Fop) (Ikemura,

121    1981), General Average Hydropathicity (GRAVY) (Sharp and Li, 1987), aromaticity

122    (Aromo) (Lobry and Gautier, 1994) and GC-content at the first, second and third codon

123    positions (GC1, GC2 and GC3), frequency of either a G or C at the third codon position

124    of synonymous codons (GC3s), the average of GC1 and GC2 (GC12)and Translational

125    selection (TrS2).

126    ENc indicates the degree of codon bias for individual genes. Over a range of

127    values from 20 to 61, lower values indicate higher codon bias, while ENc equal to 61

128    means that all codons are used with equal probability (Novembre, 2002; Wright, 1990).

129    CAI values measure the extent of bias toward preferred codons in highly

130    expressed genes. CAI values range between 0 and 1.0, with higher CAI values

131    indicating higher expression and higher CUB (Lee et al., 2010; Sharp and Li, 1987)

132    under the assumption that translational selection would optimize gene sequences

133    according to their expression levels.

134        CBI is another measure of directional codon bias, based on the degree of

135    preferred codons used in a gene, like to the frequency of optimal codons. It measures

136    the extent to which a gene uses a subset of optimal codons. In genes with extreme codon

137    bias, CBI will be equal to 1, whereas in genes with random codon usage the CBI values

138    will be equal to 0 (Bennetzen and Hall, 1982).

139        Fop is a species-specific measure of bias towards particular codons that appear

140    to be translationally optimal in particular species. It can be calculated as the ratio

141    between the frequency of optimal codons and the total number of synonymous codons.

142    Its values range from 0 if a gene contains no optimal codons to 1 if a gene is entirely

143    composed of optimal codons (Ikemura, 1981). The determination of optimal codons was

144    carried out based on the axis 1 ordination, the top and bottom 5% of genes were

145    regarded as the high and low bias datasets, respectively. Codon usage in the two data

146    sets was compared using chi-square tests, with the sequential Bonferroni correction to

147    assess significance according to Peden(Peden, 1999). Optimal codons were defined as

148    those that are used at significantly higher frequencies (p-value < 0.01) in highly

149    expressed genes compared with the frequencies in genes expressed at low levels.

150        GRAVY values were calculated as a sum of the hydropathy values of all the

151    amino acids encoded by the codons in the gene product divided by the total number of

152    residues in the sequence of the protein. The more negative the GRAVY value, the more

153    hydrophilic the protein is, whereas while the more positive the GRAVY value, the more

154    hydrophobic the protein (Sharp and Li, 1987).

155        Aromo values denote the frequency of aromatic amino acids (Phe, Tyr, Trp)

156    encoded by the codons in the gene product. (Lobry and Gautier, 1994).

157        TrS2 estimates the codon-anticodon interaction efficiency revealing bias in

158    favour of optimal codon-anticodon energy and represents the translational efficiency of

159    a gene. TrS2 value > 0.5 shows bias in favour of translational selection according to

160    Gouy and Gautier (Gouy and Gautier, 1982; Uddin et al., 2017; Uddin and Chakraborty,

161    2018).

162    **2.2. Codon Pair Score and Codon Pair Bias**

163        The determination of codon pair biases in coding sequences was performed

164 using CPBias (https://rdrr.io/github/alex-sbu/CPBias/) developed in R. as described by

165 Coleman et al.(Coleman et al., 2008). The CPS is defined as the natural logarithm of the

166 ratio of the observed over the expected number of occurrences of a particular codon pair

167 in all protein-coding sequences of a species. The CPB was used as an index and also to

168 determine the bias in CPS among the virus and host genes. The expected number of

169 codon pair occurrences estimates the number of codon pairs to be present if there is no

170 association between the codons that form the codon pair. It is also calculated to be

171 independent of codon bias and amino acid frequency (Coleman et al., 2008). A negative

172 CPS value means that a particular codon pair is underrepresented, whereas a positive

173 CPS value indicates that a particular codon pair is overrepresented in the analysed

174 protein-coding sequences. Codon pairs that are equally under- or overrepresented have a

175 CPS equidistant from 0. We calculated CPS for each of the 3,721 possible codon pairs

176 (61 x 61 codons).

### 2.3. ENc-plots

178 The ENc-plot was used to analyse the influence of base the composition on the

179 codon usage in a genome (Hartl et al., 1994). The ENc values were plotted against

180 GC3s values and a standard curve was generated to show the functional relationship

181 between ENc and GC3s values under mutational bias rather than selection pressure. In

182 genes where codon choice is constrained only by a G+C mutational bias, the predicted

183 ENc values will lie on or close to the GC3s standard curve. However, the presence of

184 other factors, such as selection effects, causes the values to deviate considerably from

185 the expected GC3s curve. The values of ENc range from 20 (when only one codon is

186 used per amino acid) to 61 (when all codons are used with equal probability). The

187 predicted values of ENc were calculated according to Hartl*et al.*(Hartl et al., 1994).

### 2.4. Clustering analysis

189 A total of 91variables of codon usage of viruses and human genes, including the

190 gene composition, RSCU frequencies and the indices described in M&M section,were

191 integrated into an input matrix to feed the clustering algorithm. The analysed variables

192 can be found in supplementary file 1. Hierarchical clustering using Euclidean distances

193 was performed. clValid package clustering algorithm was used to choose and validate

194 the best clustering method. Flexible Procedures for Clustering (fpc:https://cran.r-

195 project.org/web/packages/fpc/index.html) was used for bootstrapping (n=100) and

196 evaluate the cluster stability. Clustering method was used to establish genes groups
197 from codon usage similarities among virus and human genes.

### 2.5. Principal component analysis (PCA)

199 Principal component analysis (PCA) was used to evaluate the codon usage
200 variation among genes as the multivariate statistical method. The axes represent and
201 allow to identify the most prominent factors contributing to the variation among the
202 genes. Since there are a total of 59 synonymous codons (including 61 sense codons,
203 minus the unique Met and Trpand stop codons), the degrees of freedom were reduced to
204 40 at removing variations caused by the unequal usage of amino acids during the
205 correspondence analysis of RSCU (Greenacre, 1984). The data were normalized
206 according to Sharp and Li (Sharp and Li, 1987) in order to define the relative
207 adaptiveness of each codon (Peden, 1999; Suzuki et al., 2005), codon usage indices
208 described above were also included as variables. PCA analyses were performed using
209 "factoextra R package" (https://cloud.r-
210 project.org/web/packages/factoextra/index.html).

### 2.6. Phylogenetic analyses

212 The DNA genome sequences of all the viruses were aligned usingClustalO
213 v1.2.4(Sievers et al., 2011). PartitionFinder 2 (Lanfear et al., 2016) was used to select
214 the best-fit partitioning schemes and models of evolution for the phylogenetic
215 analysisThe evolutionary model was set on generalised time-reversible substitution
216 model with gamma-distributed rate variation across the sites and a proportion of
217 invariable sites (GTR +G). The final phylogeny was calculated using fasttree(Price et
218 al., 2009). The bootstrap consensus trees inferred from 1000 replicates were retainedin
219 the bootstrap and the final trees were drawn usingFigtree
220 (https://github.com/rambaut/figtree/releases).

### 3. Results

### 3.1. Phylogeny

223 Up to late April, a total of □ 500 SARS-CoV-2 β-coronavirus genomes became
224 available and the number of available genomes incremented substantially.The total
225 available sequences of β-coronavirus were downloaded from the NCBI. In order to
226 evaluate the variability and to select the best candidates for codon usage and nucleotides

7

227     content analyses, we performed phylogenetic analyses by implementing the GTR + G

228     model according to partition finder results. Firstly, a phylogenetic tree was constructed

229     (data not shown) using the whole genome sequences of SARS-CoV-2 reported in

230     humans from Spain, USA, Italy, South America, China, Korea Japan and Australia,the

231     references genomes of MERS and SARSand the viruses genomes isolated from bats,

232     pangolins, civets, hedgehogs, *Bos taurus*, and canids (Supplementary file 1). Since all

233     SARS-CoV-2 genome sequences remained together in the same cluster we selected

234     representative viruses genomes randomly from each country and from each host and

235     constructed a final phylogenetic tree (Supplementary file 2). The topology of this tree

236     showed that SARS-CoV-2 samples diverge from a commonnode close to the bat virus

237     (Accession MN996532) and that all of them diverge from a common node very close to

238     pangolin viruses (Accessions MT040333, MT040334, MT040335, MT040336,

239     MT072864). From a distant node that contained the node of SARS-CoV-2 and also

240     clustereda set of bat viruses, SARS virus diverges and grouped within a node together

241     with Civets viruses, but also adjacent and very close to other bats viruses (Accession

242     KY417146, KT444582 and KY417150). MERS viruses grouped in a different node also

243     very close to other cluster containingbat viruses (Accession MF593268 and KC869678).

244     Furthermore, adjacent to this node we found hedgehogs viruses (Accession KC545383,

245     KC545386, MK679660, MK907286, MK907287 and NC_039207).

246        For further analyses of molecular features and codon usage properties of viral

247     genes, we selected the viruses according to the phylogenetic tree described above. First,

248     we identified the 3 nodes containing the human viruses: SARS-CoV-2, SARS and

249     MERS. Then we selected the closest viruses to each human virus within the nodes and

250     classified them according to their hosts (bats, civets, hedgehogs and pangolins). From

251     the human viruses only the references MERS (NC_019843), SARS (NC_004718) and

252     SARS-CoV-2 (NC_045512)were used. The total selected CDSs comprised104 viral

253     CDS.

254     **3.2. Viral gene codon usage patterns**

255        PCA of codon usage and molecular features of viral genesof SARS, MERS,

256     SARS-CoV-2 and related viral genes of the non-human host (bats, civets, hedgehogs

257     and pangolins) were performed in order to characterise the genes and distinguish

258     important gene features among the viral gene families and species. PCA showed that the

259     genes dispersed differentially according to the kind of gene rather than tothe host that

8

260    the viruses infect. The genesdistribution depending on the host was observedfor only

261    some genes(Figure 1 A and B).Most of the genes belonging to the same gene family

262    overlapped or positioned in a very short ratio from each other. The position for

263    nucleocapsid protein N is shared for all virusesexcept for SARS-CoV-2 that is the most

264    distant from the group followed by the ones of SARS and MERS. Regarding the

265    envelope protein E, the gene of civetsSARS and human SARSoccupied the same

266    positionwhile the gene of human SARS-CoV-2 and MERS distributed distantly from

267    the group and alsofrom each other. The genes that encode for the membrane

268    glycoprotein M also showed a distribution along positive axis 1. However,for thehuman

269    SARS-CoV-2 and hedgehog and bat MERSthese genes distributed away from the genes

270    of pangolin and bat SARS-CoV-2, human MERS and bat,civet and humanSARS toward

271    the inferior left quadrant.The main factors that contributedto the dispersion ofthese 3

272    genes (nucleocapsid protein N, envelope protein E and membrane glycoprotein M) were

273    CBI, Fop, TrS2, C3, C3s, CpG and GC.In addition, the membrane glycoprotein M

274    seems to be more influenced by the codon frequency of the codon TAC for Tyrosine

275    and CTG for Lysine. Whereas for hedgehogs MERS and human SARS-CoV-2 the

276    codons TAT for Tyrosine and TTA for Lysine as well as the GC bias, among others

277    contributed most. All of them are strongly influenced by the A/T composition in the

278    third codon position.On the other hand, the spike proteins S of all viruses distributed

279    toward positive values of PC2and were also highly influenced by the A/T compositionin

280    the third codon position. All these genes groupedvery closelyexcept for the gene of bats

281    and human MERS. Several ORF proteins occupied the same area overlapping or

282    positioning very close to each other.Other ORF genes such as ORF6, ORF8 and

283    ORF3of the human virus SARS-CoV-2 distributed distantly.

284         The CUBwas estimated based on ENc values. The values of ENc range from 20

285    (when only one codon is used per amino acid) to 61 (when all codons are used with

286    equal probability).Genes whose ENc values are lower than 50 are considered to have

287    skewed codon usage. All viruses presented a wide range of ENc values (26 <ENc< 58)

288    that varied mainly depending on the type of genes. The highest ENc median was

289    observed for the human MERS (ENc~50.40), followed by bats MERS (ENc ~49.64).

290    The hedgehogs MERS showed an ENc median of 47.70, being closer to the values that

291    SARS and SARS-CoV-2 presented. Civets SARS showed an ENc median of 47.76. For

292    bats SARS the ENc median was 47.80 and the ENc median for the human SARS was

293    47.66. The lowest ENc valueswere observed for the genes of SARS-CoV-2 viruses.

294  Bats SARS-CoV-2 showed a median ENc of 46.39. For pangolins SARS-CoV-2, the
295  ENc median was 44.80 and for the human SARS-CoV-2, the ENc median was 44.34
296  (Figure 1 C).

297  Codon bias of viral genes classified bythe type of gene or the gene
298  familyshowed that the envelope protein E presented ENc values that ranged from 42 to
299  61, being the genes of the human SARS-CoV-2 and bat SARS-CoV-2 the ones that
300  presented the lowest valuesof all the genes.Membrane glycoprotein M genes presented
301  ENc values that ranged from 43.32 to 61.00,andthe gene of hedgehogs MERS,the one
302  with the lowest value. The membrane glycoprotein M of Human MERS, Civets,bats and
303  human SARS showed virtually non-biased codon usage. The nucleocapsid protein N
304  showed ENc values that ranged from 49.08 to 55.30 being the human MERS gene,
305  followed by the bats MERSgene, the viral genes that presented the lowest values.

306  The genes that encode for the spike protein S presented ENc values that ranged
307  from 44.16 to 47.68. The gene of the humanSARS-CoV-2 showed the lowest ENc value
308  followed by the genes of bat and pangolin SARS-CoV-2. ORF genes presented ENc
309  values that ranged from 26.60 to 57.89, being the human SARS-CoV-2 the virus that
310  presented the lowest and the highest ENc value for ORF7 and ORF10 respectively.

311  **3.3. CPB analysis between viruses and lung tissue highly expressed genes**

312  Since human viruses have codon usage biases that match highly expressed
313  proteins in the tissues they infect (Miller et al., 2017) we selected the highly expressed
314  genes in lungs tissue to compare the gene molecular features, CUB and CPB of SARS,
315  SARS-CoV-2 and MERS viruses in relation with human genes. Since viruses replicate
316  exclusively inside living cells, many viruses are influenced by host codon pair
317  preferences,being thereflectionof the CPB or CPS of their hosts. Therefore, CPSof viral
318  genes was evaluated and compared with the CPS of human genesin order todetermine
319  whether viruses that infect humans have similar codon pair preferences to their
320  host.Viral genes showed lower CPS frequencies than human genes. The median CPS for
321  the viral genes was 0.053 for civet SARS, 0.051 for bat SARS, and 0.053 for human
322  SARS. The median CPS that presented the MERS genes was 0.048 for hedgehogMERS,
323  0.046 for bat MERS and 0.037 for human MERS. SARS-CoV-2 showed median CPS
324  values of 0.064 for bat SARS-CoV-2, 0.065 for pangolins SARS-CoV-2 and 0.061 for
325  human SARS-CoV-2. The median CPS for highly expressed human genes in lungs was
326  0.11 (Figure 2A). We also calculated the median CPS of viral genes for each gene

327   family. The genes that encode for ORF7b and spike protein S were the genes with the

328   highest median of CPB values (0.081 and 0.061 respectively), followed by ORF3a and

329   nucleocapsid protein N (0.058 and 0.057) indicating that particular codon pairs are

330   overrepresented in these genes. Envelope protein E and the membrane glycoprotein M

331   showed values of 0.048 and 0.044 respectively (Figure 2B).

332        Furthermore, different CPB among the viral genes families of different

333   virusesspecies that infect different hostswere also observed (Figure 2C). Bat SARS

334   showed the highest CPB for ORF7b (~0.1) and ORF7a (~0.048). The human SARS-

335   CoV-2 showed the highest CPB values for ORF1a (0.089), ORF1ab (0.051),

336   nucleocapsid protein N (0.048) and spike protein S (0.042). Envelope protein E and

337   membrane glycoprotein M presented CPB values of 0.0005 and 0.02 respectively. In

338   pangolin SARS-CoV-2 the highest CPB value was for the spike protein S (0.060)

339   followed by ORF3a (0.037) and ORF8 (0.033) while inbat SARS-CoV-2 the highest

340   CPB              value              was              for              the

341   spike protein S (0.076) followed by ORF1ab (0.050).The highest CPB values for human

342   MERS genes were observedfor a non-structural protein (0.032) followed by the

343   spike protein S (0.021) and the nucleocapsid protein N (0.008). In bat MERS the highest

344   CPB value was for spike protein S(0.062) followed by ORF1a (0.022)and by the

345   membrane glycoprotein M (0.022).In hedgehog MERS the highest CPB value was for

346   ORF3a (0.049), followed by ORF1a/b (0.042)and the membrane glycoprotein M

347   (0.029).In orderto evaluate the fitness and specialization of the viruses in theirhosts, we

348   compared the CPS of the viral genes derived from the different hosts against the CPS of

349   highly expressed genes in human lungs tissue by performing CPS correlation analysis.

350   The 3721 codon pairs of all viruses were compared with the 3721 codon pairs of human

351   genes. Correlation analyses showed low R values with a not clear dependence on the

352   human host codon pairs (Supplementaryfile3).

353   **3.4. Viral genes clustering analysis**

354        Further analysis using hierarchical clustering of viral genes provided additional

355   qualitative information about how similar certain genes are in terms of molecular and

356   codon usage parameters(Supplementaryfile4). All the genes encoding for viral

357   nucleocapsid protein N of all the virusesgrouped together demonstrating a high level of

358   conservation. The human SARS-CoV-2 gene presented a very low CPB value in

359   comparison with the genes of the other viruses.

11

360     Envelope protein E genes grouped in different clusters. Thethree SARS viruses

361     (civets, bats and human hosts) grouped with the bat and pangolin SARS-CoV-2 and

362     showed high and similar CPB values.Whereas the human SARS-CoV-2 grouped better

363     with the human and bat MERS, both showinghigh CPB values. Hedgehog envelope

364     protein E gene grouped distantly with ORF and hypothetical proteins.These genes

365     presented a high preference for the codon GCG for Alanine. Furthermore, the codons

366     GGT and CCT were preferable for Glycine and Proline respectively. The

367     hydrophobicity, CAI and Fop values are higher for SARS and SARS-CoV-2 genes.

368     The membrane glycoprotein M genes also appeared in different clusters. The

369     humanSARS-CoV-2 gene grouped with the human and bat MERS gene.Only the human

370     SARS-CoV-2 gene showed a preference for the codon GCA for Alanine and AGG for

371     Arginine. A preference for the codons CCA, ACT and GTT for Proline, Threonine and

372     Tyrosinewas also observed. However, this was the gene that presented the lowest CPB

373     value. The pangolin and bat SARS-CoV-2 gene grouped with the threeSARS genes in a

374     different clusterand showed a higher CPB. The hedgehogs MERS grouped distantly.The

375     codon CAC was observed more frequently for Histidine in the three SARS viruses.

376     Pangolin SARS-CoV-2 showed a preference for GAC for Aspartate and CCA for

377     Proline whereas the bat SARS-CoV-2 showed a preference for GGA for Glycine and

378     GTA for Valine.

379     The genes encoding for the spike protein S appeared in different clusters. For the

380     human SARS-CoV-2,this gene grouped alone with the genes that encode for

381     hypothetical and ORF proteins and showed high bias only for the codon TTG for Lysine

382     and GAA for Glutamate. The genes for the rest of the viruses clustered together with

383     other ORF genesand with the membrane glycoprotein M of hedgehogs, although

384     distantly.This analysis also showed that CPB is highly related to the dinucleotide bias.

385     In this cluster the genes that encode for the ORF genes showed low values of CPB,

386     being the lowest of all the clusters. Conversely, the genes that encode for the spike

387     protein S presented high CPB values.

388     **3.5. Viral and human genes clustering analysis and CPS correlation**
389     **analysis**

390     Since virus fitness specialization could be dependent on the translation

391     machinery for only some particular genes, we performed clustering analysis based on

392     the codon usage and molecular features including human genes and the genes of

12

393    SARS-CoV-2, SARS and MERS(Supplementary file 1). From a total of 463 human

394    genes, 70 genes (15.1%)grouped in clusters together with viral genes. These genes were

395    selectedand extracted to make an illustrating heatmap including both, human and all

396    viral genes (Figure 3).Furthermore, in order to figure out whether the viral genes that

397    composed particular clusters present codon pair frequencies that correlate with the

398    human genes of the same clusters, we evaluated the CPS of both, viral and human genes

399    for each block and the correlation between each other.

400        For SARS-CoV-2, 8 out of 12 genesgrouped with 40 human genes distributed in

401    4 clusters. The first cluster comprised3 viral genes: the nucleocapsid protein N, ORF1a

402    and ORF1ab together with 19 human genes: COL6A5 (3), OVCH1 (2), DNAH12 (2),

403    ROS1 (3), SCN7A, RP1, ABCA13 (4) and LRRK2 (3) and the CPS correlation was

404    R☐0.11 (p-val☐$7.1 \times 10^{-8}$). The CPB values for the human and viral genes were -0.04

405    and 0.15 respectively. The second cluster contained 2 viral genes: the spike protein S

406    and ORF7b together with 9 human genes: AGTR2, TMEM212, CALCRL, FMO2 (2),

407    FAM216B, CXCL10, MMP13 and SLC6A14. For this group, the CPS correlation

408    between viral and human genes was R☐0.21 (p-val☐0.0027) and the CPB values for the

409    human and viral genes were 0.05 and 0.12 respectively. The third cluster grouped only 1

410    viral gene, the envelope protein E, with 10 human genes: IFNG, ST8SIA6 (2),

411    SLC39A8 (2), CDC20B (3) and DCSTAMP (2), and the CPS correlation between viral

412    and human genes wereR☐0.29 (p-val☐0.0052) and the CPB values for the human and

413    the viral genes were 0.25 and 0.16 respectively. In the fourth cluster 2 human genes,

414    TSPAN19 (2)grouped with the viral genes ORF6 and ORF7a and showed a CPS

415    correlation of R☐0.5 (p-val☐0,0016). The CPB values for the human and viral genes

416    were 0.13 and 0.14 respectively.

417        For SARS virus (Supplementary file 5), 12 out of 14 genes grouped with 39

418    human genes distributed in 5 clusters. The first cluster contained 11 human genes:

419    AGTR2, TMEM212, CALCRL, FMO2 (2), FAM216B, CXCL10, MMP13, SLC6A14,

420    and BCL2A1 (2) and 7 viral genes: envelope protein E, hypothetical proteins (5) and

421    membrane glycoprotein M.The CPS correlation for human and viral genes was

422    R☐0.28 (p-val☐$2 \times 10^{-16}$). The second cluster was composed of 7 human genes:

423    C7orf77, SNTN (2), MUC1 (2) and WIF1 (2).The CPS correlation between human and

424    viral genes was R☐0.77 (p-val☐0.00078). The third cluster consisted of 19 human

425    genes:

426    COL6A5 (3), OVCH1 (2), DNAH12 (2), ROS1 (3), SCN7A, RP1, ABCA13 (4) and

13

427  LRRK2 (3) and 2 viral genes: ORF1a and ORF1ab. The CPS correlation between

428  human and viral genes was R☐0.31 (p-val☐2.2x10$^{-16}$). The fourth cluster was

429  composed of 2 TSPAN19 human genes and one viral hypothetical protein. CPS

430  correlation for this cluster was statistically not significant.

431  For MERS virus (Supplementary file 6), the 10 genes grouped into 6 clusters

432  together with 65 human genes. The first cluster comprised 19 human genes: COL6A5

433  (3), OVCH1 (2), DNAH12 (2), ROS1 (3), SCN7A, RP1, ABCA13 (4) and LRRK2 (3)

434  and 2 viral genes: ORF1a and ORF1b. CPS correlation for human and viral genes was

435  R☐0.29 (p-val☐0.021). The second cluster was composed of 8 human genes: AGTR2,

436  TMEM212, CALCRL, FMO2 (2), CXCL10, MMP13 and SLC6A14 and 1 viral gene

437  that encodes for the spike protein S.The third cluster contained 6 human genes:

438  CLEC12A (3), FAM216B, DNAAF6 and PIH1D3 and 3 viral genes: non-structural

439  protein (2) and the nucleocapsid protein N. The fourth cluster consisted of 7 human

440  genes:                                                                                   C7orf77,

441  SNTN(2), MUC1 (2) and WIF1 (2) and 1 viral gene encoding for a non-structural

442  protein. The fifth cluster was composed of 9 human genes: ST8SIA6 (2), SLC39A8 (2),

443  CDC20B (3) and DCSTAMP (2) and 2 viral genes: non-structural protein and

444  membrane glycoprotein M. The sixth cluster comprised 16 human genes: CLEC6A,

445  IL1RL1 (2), VNN2 (4), RTKN2 (2), SDR16C5 (3), IL18R1 (2), ACADL, and

446  DNAH12 and the viral gene that encodes for envelope protein E.

447  Furthermore, we observed that 28 out of 70 human genes comprised a core of

448  genes that appeared into the clusters together with the viral genes for the three human

449  viruses (Figure 4 and Table 1). Between MERS and SARS-CoV-2 only 9 human

450  geneswere shared, 7 genes were shared between onlyMERS and SARS and 2 genes

451  were shared between only SARS-CoV-2 and SARS. The 28 human genes that appeared

452  in the clusters together with the genes of the 3 human viruses were retrieved against

453  PheGenI     (https://www.ncbi.nlm.nih.gov/gap/phegeni/#pgGAP)     and     DisGeNET

454  (https://www.disgenet.org/) in order to classify the human diseases associated with the

455  malfunction of the identified genes as an approach for possible human diseases or

456  collateral effects caused by the viral infections (Supplementary file 1). According to

457  PheGenI and DisGeNET,14 out of the 28 genes were associated with 27 diseases. All

458  the diseases appeared in nearly equal proportions. A slightly higher frequency was

459  observed for Nervous System Diseases (12 genes), Neoplasms (12 genes), Respiratory

460  Tract Diseases (11 genes), Pathological Conditions Signs and Symptoms(11 genes),

461     Neonatal Diseases and Abnormalities (11 genes) and Cardiovascular Diseases (10

462     genes) among others (Figure 4 B).

463     **4. Discussion**

464         In humans, coronaviruses cause mainly respiratory tract infections.Previously,

465     six coronaviruseswere identified as human-susceptible viruses, being theβ-

466     coronaviruses, SARS-CoV and MERS-CoV, responsible for severe and potentially fatal

467     respiratory tract infections (Yin and Wunderink, 2018). In December 2019 rising

468     pneumonia cases caused by a novel β-coronavirus (SARS-CoV-2), occurred in Wuhan,

469     China. The disease was officially named coronavirus disease 2019 (COVID-19). It was

470     found that the genome sequence of SARS-CoV-2 is 96.2% identical to the bat

471     coronavirus RaTG13, whereas it shares 79.5% of identity to SARS-CoV. It has been

472     proposed that SARS-CoV-2 may have originated from bats or unknown intermediate

473     hosts that could involve pangolins, crossing the species barrier into humans. (Guo et al.,

474     2020). Rather, bats are the natural reservoir of a wide variety of coronaviruses,

475     including SARS-CoV-like and MERS-CoV-like viruses (Banerjee et al., 2019;

476     Hampton, 2005; Li et al., 2005).

477         Up to late April, a total of ⬜ 500 SARS-CoV-2 β-coronavirus genomes became

478     available and this number is continuously increasing at an unprecedented rate. In our

479     studies, we analysed the coronaviruses from different host species and retrieved the

480     phylogeny in order to select the best candidate genomes for further analysis of codon

481     usage and molecular relationship with the human host. As previously reported, we

482     found that for SARS-CoV-2, SARS and MERS, bats seem to be a common natural host

483     or reservoir.

484         Because viruses do not have tRNAs and rely on host cell machinery for

485     replication, co-evolution between RNA viruses and their hosts' codon usages have been

486     observed (Franzo et al., 2017; Rahman et al., 2018; Simón et al., 2017). Furthermore,

487     the adaptation of viruses that replicates in multiple hosts should involve a trade-off

488     between precise and functional matching to fit the diverse tRNA pools of multiple hosts

489     (Tian et al., 2018). Conversely, single host viruses are expected to have specialised to

490     match the host tRNA repertoire.

491         Since human viruses have CUB that matches highly expressed proteins in the

492     tissues they infect (Miller et al., 2017) we selected the highly expressed genes in lung

493     tissues to compare gene molecular features of SARS-CoV-2,SARS and MERS in

494   relation with human genes. In general, in our studies for all the analysed viruseswe

495   found that the total gene repertoire had a similar ENc average that differs only □1 unit

496   with respect to their non-human host they come from, reflecting the molecular features

497   of their original host. Furthermore, as demonstrated in our clustering analysis, codon

498   pair usage seems to be dependent on the dinucleotide bias and the human CPB was

499   higher for human genes than for viruses genes as previously reported (Kames et al.,

500   2020; Kunec and Osterrieder, 2016).

501       Moreover, our analyses allowed us to distinguish not only the main factors that

502   contribute to the distribution of the genes along the axes in PCA, but also to determine

503   some particular different features among human and non-human viruses in specific

504   genes that could be important for explaining the virus infection evolution. In contrast to

505   SARS-CoV-2 of bats and pangolins, human SARS-CoV-2 exhibited a differential

506   distribution in particular genes that depended mostly on the A/T content in the third

507   codon position, which is in accordance with human SARS-CoV-2 gene composition

508   reports (Alnazawi et al., 2017; Alonso and Diambra, 2020; Kames et al., 2020; Tort et

509   al., 2020).

510       Important viral genes such as the membrane glycoprotein M, that is involved in

511   the membrane transport of nutrients, the bud release, the formation of the envelope, the

512   virus assembly and in the biosynthesis of new virus particles (Guo et al., 2020),

513   distributed differentially from the non-human viruses indicating that it is highly

514   influenced by A/T content. Surprisingly, this gene was positionednear of the hedgehog's

515   MERS gene, suggesting similar molecular patterns between two distant viruses.

516       The envelope protein E, that functions as an ion channel and regulates virion

517   assembly and the immune system of the host (Guo et al., 2020; Yin and Wunderink,

518   2018), showed the same tendency toward A/T ending codons for the human viruses

519   SARS-CoV-2 and MERS. Human genes using A/T ending codons is also a common

520   feature in several human genes since they present a wide rate of GC content (27-97%)

521   (Bernardi, 2015; D'Onofrio et al., 1991). Both, membrane glycoprotein M and envelope

522   protein E genes have a higher CUB in comparison with human MERS and SARS which

523   is not in accordance with the trade-off theory that postulates that cross-species virus

524   transmission demands relaxing the codon usage pattern (Kunec and Osterrieder, 2016).

525   However, this phenomenon could be explainedby a selection pressure in favour of the

526   virus replication in the new host or due to the recent cross-species virus transmission as

527   we know it occurred. If this is the case, the analysis of newisolates from infectedhumans

16

528  should tend to show incremented ENc values for the envelope protein E and membrane

529  glycoprotein M genes.

530        Nevertheless, only the envelope protein E clustered together with human genes,

531  demonstrating similar molecular patterns that could mean an advantage for virus

532  replication in humans facilitating the virion assembly and the regulation of the immune

533  system of the host. Furthermore, positive CPB and an incremented CPS correlation for

534  the cluster that grouped the envelope protein E with human genes supports the

535  hypothesis of a facilitated translation depending on codon usage and codon pairs.

536  Similar patterns were observed for ORF6 and ORF8genes, which are involved in the

537  viral pathogenesis, apoptosis induction and inflammatory responses in the host (Chen et

538  al., 2020b; Diemer et al., 2010). These genes grouped with human genes in different

539  clusters and showed also and incremented CPS correlation.

540        Studies in different viruses species have reported high conservation of the genes

541  that encode for the nucleocapsid protein N among virus families (HORNE, 2013; Kunec

542  and Osterrieder, 2016; Masters, 2019; Nathan and Scobell, 2012; Parker and Masters,

543  1990). Therefore, it is expected that codon usage and molecular patterns present similar

544  features as observed in our studies. Nevertheless, SARS-CoV-2 nucleocapsid gene

545  tends to distribute slightly far from the rest of viruses' capsids genes and toward a

546  higher A/T content in PCA. Both, the human SARS and SARS-CoV-2 nucleocapsid

547  protein N present high CPB suggesting a specialization acquirement in the human host.

548        Two viral genes that also present high CPB are ORF1a/b, that encodes for the

549  replicase complex (polyproteins pp1a and pp1ab) and the Spike protein S that

550  participates in the early viral infection by attaching to the host receptor ACE2 and

551  mediating the internalization of the virus (Guo et al., 2020). In our studies, ORF1a/b

552  grouped with the gene that encodes for the nucleocapsid protein N, indicating that their

553  molecular features are highly conserved and are also presentin several human genes.

554  This result is in concordance with previous works that proposed these genes as

555  candidates for deoptimization for the design of attenuated vaccines due to their high

556  positive CPB values (Kames et al., 2020). Instead, the gene that encodes for the spike

557  protein S, grouped with ORF7 (involved in viral pathogenesis and apoptosis induction)

558  that also presents high and similar positive CPB values. For all of them, a higher rate of

559  A/T composition in the third codon position was observed. Changes in the third position

560  produce synonymous substitutions that could have conducted to a codon optimization in

561  human cells using the host machinery that translates only genes whose molecular

562  features match the viral needs. Some viral genes seem to have been favoured for an

563  increased viral replication in humans and optimized by using or mimicking some

564  particular molecular patterns of human genes. But only some genes, such as the

565  envelope E, the ORF 6 and 8, could be the key for an exacerbated viral pathogenesis.

566  Furthermore, because of these molecular and codon usage similarities between some

567  highly expressed human genes and viral genes that occupy the same clusters, the

568  translation machinery of the host could propitiate the translation of viral genes to the

569  detriment of human gene expression in lung tissues.Indeed, mistranslation or de-

570  regulation of protein synthesis has been reported as a consequence of tRNA miss-

571  modification and imbalanced tRNA expression, causing diseases(Lant et al., 2019).

572  Recent studies have also proposed that an unbalance in the tRNAs pools of the infected

573  cells could occur and would explain the collateral effects observed in some viral

574  infections (Alonso and Diambra, 2020). Since COVID-19 outbreak, several studies

575  associated with different pathologies have been performed in orderto find out how

576  damaging this new virus is for the human being. Hereby, in our studies we provided a

577  list of human genes that could be particularly affected as a consequence oftheir

578  molecular similarities with viral genes, not only belonging to SARS-CoV-2 but also to

579  SARS andMERS. The malfunction of these genes has been associated with different

580  human pathologies and is in continuous increase. Patients infected with COVID-19

581  typically present fever and respiratory symptoms.Nevertheless, it has been reported an

582  increased risk for complications of hypertension, congestive heart failure, and

583  atherosclerosis conducting to anincreased presence of cardiovascular comorbidities

584  (Clerkin et al., 2020; Li et al., 2020; Zheng et al., 2020).Also, some patients have

585  experimented gastrointestinal manifestations (Wong et al., 2020),neurologic

586  complications (Bridwell et al., 2020; Dugue et al., 2020), and complications associated

587  with the endocrine and urogenital systems,among others (Wang et al., 2020; Wu et al.,

588  2020). Diseases and collateral effects caused by COVID-19 infections could be a

589  consequence of the malfunction of the genes listed in our work. Therefore,they should

590  be considered to be incorporated into susceptibility population studies for respiratory

591  viral infections.Hereby, these results lay the groundwork for further research in the field

592  of human genetics associated with the new viral infection, COVID-19, caused by

593  SARS-CoV-2 and for the development of antiviral preventive methods.

594  **5. Conclusions**

595    In our study, we described the main factors that shape CUB in SARS-CoV-2,
596    SARS and MERS in comparison with highly expressed genes in human lung tissue and
597    revealed matching features with human genes that could have favoured the virus for an
598    incremented pathogenesis. Furthermore, we provided a list of candidate human genes
599    that could be involved in the viral infection and had not been described yet which could
600    be the key for explaining collateral effects and the human susceptibility to viral
601    infectionsandshould be considered to be incorporated into genetic population studies.
602

## 6. Declarations

### 6.1. Competing interests

605    The authors declare that the research was conducted in the absence of any
606    commercial or financial relationships that could be construed as a potential conflict of
607    interest.

### 6.2. Funding

### 6.3. Authors' contributions

612    L. K wrote and revised the manuscript L.M. designed the study, performed the
613    bioinformatics analysis, wrote and revised the manuscript.

### Legends to figures

615    **Figure 1:** viral genes distribution in PCA plot in the first 2 axes and ENc-GC3s
616    plot of SARS-CoV-2 (NC_045512), SARS (NC_004718) and MERS (NC_038294),
617    and the related virus of non-human hosts: CivetsSARS (AY686864), bats MERS
618    (KC869678), bats SARS (KY417150), bats SARS-CoV-2 (MN996532), pangolins
619    SARS-CoV-2 (MT040336) and hedgehogs MERS (NC_039207). **A)**. Distribution of
620    viral genes in PC1 and PC2. **B)**. Main factors represented by vectors that contribute to
621    the distribution of viral genes in PC1 and PC2. **C).** Distribution of the effective number
622    of codons (ENc) in relation to the GC3s of viral genes. The standard curve of ENcis
623    indicated in solid line.
624    **Figure 2:**codon pair score of viral genes of SARS-CoV-2 (NC_045512), SARS

625    (NC_004718) and MERS (NC_038294), and the related virus of non-human hosts:

626    civets SARS (AY686864), bats MERS (KC869678), bats SARS (KY417150), bats

627    SARS-CoV-2 (MN996532), pangolins SARS-CoV-2 (MT040336) and hedgehogs

628    MERS (NC_039207) and human genes. **A).** codon pair frequencies for each virus. **B).**

629    codon pair frequencies for each gene and classified by type of viral gene. **C).** Codon

630    pair bias for each viral gene vs protein length.

631            **Figure 3:** Heatmap of clusters(1 to 4) using a hierarchical method of viral genes

632    forSARS-CoV-2 (NC_045512) of the human host and human genes based on the

633    molecular features. CPB correlation is included in the left for each cluster relating the

634    CPB of human genes (horizontal axis) and CPB of the viral genes (vertical axis).

635            **Figure 4: A).**Venn diagram representing the number of human genes that

636    clustered together with viral genes for SARS-CoV-2 (NC_045512), SARS

637    (NC_004718) and MERS (NC_038294) based on the molecular features.**B).**Diseases

638    frequencies associated to human genes grouped with viral genes of SARS-CoV-2,

639    SARS and MERS in the clustering analysis.

640    **Tables**

Table 1: Human genes shared among the clusters of the three coronaviruses SARS-CoV-2, SARS and MERS

| Accession Id | Gene name | Proteins |
|---|---|---|
| NP_000677 | AGTR2 | type-2 angiotensin II receptor |
| NP_001157908 | TMEM212 | transmembrane protein 212 |
| NP_001258680 | CALCRL | calcitonin gene-related peptide type 1 receptor precursor |
| NP_001265227 | COL6A5 | collagen alpha-5(VI) chain isoform 1 precursor |
| NP_001288276 | FMO2 | dimethylaniline monooxygenase |
| NP_001305861 | FAM216B | protein FAM216B |
| NP_001340108 | OVCH1 | ovochymase-1 precursor |
| NP_001352957 | DNAH12 | dynein heavy chain 12, axonemal isoform 4 |
| NP_001365831 | ROS1 | proto-oncogenetyrosine-proteinkinase ROS isoform 3 precursor |
| NP_001451 | FMO2 | dimethylaniline monooxygenase |
| NP_001556 | CXCL10 | C-X-C motif chemokine 10 precursor |
| NP_002418 | MMP13 | collagenase 3 preproprotein |
| NP_002935 | ROS1 | proto-oncogenetyrosine-proteinkinase ROS isoform 1 precursor |
| NP_002967 | SCN7A | sodium channel protein type 7 subunit alpha |
| NP_006260 | RP1 | oxygen-regulated protein 1 isoform 1 |
| NP_009162 | SLC6A14 | sodium- and chloride-dependent neutral and basic amino acid transporter B(0+) |
| NP_689914 | ABCA13 | ATP-binding cassette sub-family A member 13 |
| NP_694996 | COL6A5 | collagen alpha-5(VI) chain isoform 2 precursor |
| NP_940980 | LRRK2 | leucine-rich repeat serine/threonine-protein kinase 2 |

20

| XP_005268686 | LRRK2 | leucine-rich repeat serine/threonine-protein kinase 2 isoform X1 |
| XP_011510923 | COL6A5 | collagen alpha-5(VI) chain isoform X1 |
| XP_011513433 | ABCA13 | ATP-binding cassette sub-family A member 13 isoform X2 |
| XP_011513434 | ABCA13 | ATP-binding cassette sub-family A member 13 isoform X3 |
| XP_011513438 | ABCA13 | ATP-binding cassette sub-family A member 13 isoform X6 |
| XP_011534357 | ROS1 | proto-oncogene tyrosine-protein kinase ROS isoform X10 |
| XP_011536184 | LRRK2 | leucine-rich repeat serine/threonine-protein kinase 2 isoform X8 |
| XP_024100505 | DNAH12 | dynein heavy chain 12, axonemal isoform X2 |
| XP_024304736 | OVCH1 | ovochymase-1 isoform X1 |

641

## Legends to supplementary files

**Supplementary file 2:** Phylogenetic tree using 118 virus genomes including the references SARS-CoV-2, SARS, and MERS and related viruses belonging to non-human host as described in Supplementary file 1. The clusters were virus grouped with SARS-CoV-2, SARS, and MERSare highlighted. The accession id is followed by the host the sample was isolated from or the county in the case of different isolation of SARS-CoV-2.

**Supplementary file 3:** Codon pair bias plots and correlation of viral genes against highly expressed human genes in lungs tissues according to the fold-change between the expression level in lung and the tissue with second-highest expression level according to "Human Protein Atlas" (https://www.proteinatlas.org/humanproteome/tissue/lung). P1: KT444582, P2: KY417146, P3: KY417150, P4: MG772934, P5: MN996532, P6: C869678, P7: MF593268, P8: FJ938064, P9: FJ938066, P10: KX432213, P11: Y572035, P12: AY686864, P13: MK679660, P14: NC_039207, P15: NC_038294, P16: T040336, P17: MT040333, P18: MT126808, P19: NC_045512, P20: MT066156, P21: T263074, P22: NC_004718, P23: MT198652, P24: MT233519, P25: MT233523, P26: T118835, P27: MT233526, P28: MT259235, P29: MT263435.

**Supplementary file 4:** Heatmaps clustering using a hierarchical method of viral genes based on the measures of molecular and codon usage patterns. The viruses whose genes were included here were SARS-Cov-2, SARS, and MERS and the closest viruses to the common nodes in the phylogenetic tree of supplementary file 2. Civets SARS (AY686864), bats MERS (KC869678), bats SARS (KY417150), bats SARS-CoV-2 (MN996532), pangolins SARS-CoV-2 (MT040336), human SARS (NC_004718), human MERS (NC_038294), hedgehogs MERS (NC_039207) and the human SARS-

21

667     CoV-2 (NC_045512).

668         **Supplementary file 5:** Heatmap of clusters (1 to 4) using a hierarchical method

669     of viral genes forSARS (NC_004718) of the human host and human genes based on the

670     molecular features. CPB correlation is included in the left for each cluster relating the

671     CPB of human genes (horizontal axis) and CPB of the viral genes (vertical axis).

672         **Supplementary file 6:** Heatmap of clusters (1 to 6) using a hierarchical method

673     of viral genes forSARS (NC_038294) of the human host and human genes based on the

674     molecular features. CPB correlation is included in the left for each cluster relating the

675     CPB of human genes (horizontal axis) and CPB of the viral genes (vertical axis).

676     **7.  References**

677     Alnazawi, M., Altaher, A., and Kandeel, M. (2017). Comparative genomic analysis
678         MERS CoV isolated from humans and camels with special reference to virus
679         encoded helicase. *Biol. Pharm. Bull.* 40, 1289–1298. doi:10.1248/bpb.b17-00241.

680     Alonso, A. M., and Diambra, L. (2020). SARS-CoV-2 codon usage bias downregulates
681         host expressed genes with similar codon usage. *bioRxiv*, 2020.05.05.079087.
682         doi:10.1101/2020.05.05.079087.

683     Banerjee, A., Kulcsar, K., Misra, V., Frieman, M., and Mossman, K. (2019). Bats and
684         coronaviruses. *Viruses* 11. doi:10.3390/v11010041.

685     Bennetzen, J. L., and Hall, B. D. (1982). Codon selection in yeast. *J. Biol. Chem.* 257,
686         3026–3031.    Available    at:    http://www.jbc.org/content/257/6/3026.full.pdf
687         [Accessed September 11, 2017].

688     Bernardi, G. (2015). Chromosome Architecture and Genome Organization. *PLoS One*
689         10, e0143739. doi:10.1371/journal.pone.0143739.

690     Bridwell, R., Long, B., and Gottlieb, M. (2020). Neurologic complications of COVID-
691         19. *Am. J. Emerg. Med.* doi:10.1016/j.ajem.2020.05.024.

692     Cavanagh, D. (2007). Coronavirus avian infectious bronchitis virus. *Vet. Res.* 38, 281–
693         297. doi:10.1051/vetres:2006055.

694     Chan, J. F. W., Kok, K. H., Zhu, Z., Chu, H., To, K. K. W., Yuan, S., et al. (2020).
695         Genomic  characterization  of  the  2019  novel  human-pathogenic  coronavirus

696     isolated from a patient with atypical pneumonia after visiting Wuhan. *Emerg.*
697     *Microbes Infect.* 9, 221–236. doi:10.1080/22221751.2020.1719902.

698     Chen, Y., Liu, Q., and Guo, D. (2020a). Emerging coronaviruses: Genome structure,
699     replication, and pathogenesis. *J. Med. Virol.* 92, 418–423. doi:10.1002/jmv.25681.

700     Chen, Y., Liu, Q., and Guo, D. (2020b). Emerging coronaviruses: Genome structure,
701     replication, and pathogenesis. *J. Med. Virol.* 92, 418–423. doi:10.1002/jmv.25681.

702     Clark, M. A. (1993). Bovine coronavirus. *Br. Vet. J.* 149, 51–70. doi:10.1016/S0007-
703     1935(05)80210-6.

704     Clerkin, K. J., Fried, J. A., Raikhelkar, J., Sayer, G., Griffin, J. M., Masoumi, A., et al.
705     (2020). Coronavirus Disease 2019 (COVID-19) and Cardiovascular Disease.
706     *Circulation*. doi:10.1161/CIRCULATIONAHA.120.046941.

707     Coleman, J. R., Papamichail, D., Skiena, S., Futcher, B., Wimmer, E., and Mueller, S.
708     (2008). Virus Attenuation by Genome-Scale Changes in Codon Pair Bias. *Science*
709     *(80-. ).* 320, 1784–1787. doi:10.1126/science.1155761.

710     Cui, J., Li, F., and Shi, Z. L. (2019). Origin and evolution of pathogenic coronaviruses.
711     *Nat. Rev. Microbiol.* 17, 181–192. doi:10.1038/s41579-018-0118-9.

712     D'Onofrio, G., Mouchiroud, D., Aïssani, B., Gautier, C., and Bernardi, G. (1991).
713     Correlations between the compositional properties of human genes, codon usage,
714     and amino acid composition of proteins. *J. Mol. Evol.* 32, 504–510.
715     doi:10.1007/BF02102652.

716     Diemer, C., Schneider, M., Schätzl, H. M., and Gilch, S. (2010). "Modulation of host
717     cell death by SARS coronavirus proteins," in *Molecular Biology of the SARS-*
718     *Coronavirus* (Springer Berlin Heidelberg), 231–245. doi:10.1007/978-3-642-
719     03683-5_14.

720     Dugue, R., Cay-Martínez, K. C., Thakur, K., Garcia, J. A., Chauhan, L. V., Williams, S.
721     H., et al. (2020). Neurologic manifestations in an infant with COVID-19.
722     *Neurology*, 10.1212/WNL.0000000000009653.
723     doi:10.1212/wnl.000000000009653.

724 Franzo, G., Tucciarone, C. M., Cecchinato, M., and Drigo, M. (2017). Canine
725   parvovirus type 2 (CPV-2) and Feline panleukopenia virus (FPV) codon bias
726   analysis reveals a progressive adaptation to the new niche after the host jump. *Mol.*
727   *Phylogenet. Evol.* 114, 82–92. doi:10.1016/j.ympev.2017.05.019.

728 Gouy, M., and Gautier, C. (1982). Codon usage in bacteria: Correlation with gene
729   expressivity. *Nucleic Acids Res.* 10, 7055–7074. doi:10.1093/nar/10.22.7055.

730 Greenacre, M. J. (1984). *Theory and applications of correspondence analysis.*
731   Academic Press Available at:
732   https://books.google.com.ar/books/about/Theory_and_Applications_of_Correspon
733   denc.html?id=LsPaAAAAMAAJ&redir_esc=y [Accessed May 27, 2018].

734 Guo, Y. R., Cao, Q. D., Hong, Z. S., Tan, Y. Y., Chen, S. D., Jin, H. J., et al. (2020).
735   The origin, transmission and clinical therapies on coronavirus disease 2019
736   (COVID-19) outbreak- A n update on the status. *Mil. Med. Res.* 7, 1–10.
737   doi:10.1186/s40779-020-00240-0.

738 Hampton, T. (2005). Bats may be SARS reservoir. *J. Am. Med. Assoc.* 294, 2291.
739   doi:10.1001/jama.294.18.2291.

740 Hartl, D. L., Moriyama, E. N., and Sawyer, S. A. (1994). Selection intensity for codon
741   bias. *Genetics* 138, 227–234. doi:10.3168/jds.S0022-0302(75)84789-8.

742 HORNE, R. W. (2013). "The Structure of Viruses," in *Scientific American* (Elsevier),
743   153–178. doi:10.1090/gsm/146/03.

744 Ikemura, T. (1981). Correlation between the abundance of Escherichia coli transfer
745   RNAs and the occurrence of the respective codons in its protein genes: A proposal
746   for a synonymous codon choice that is optimal for the E. coli translational system.
747   *J. Mol. Biol.* 151, 389–409. doi:10.1016/0022-2836(81)90003-6.

748 Kames, J., Holcomb, D. D., Kimchi, O., DiCuccio, M., Hamasaki-Katagiri, N., Wang,
749   T., et al. (2020). Sequence analysis of SARS-CoV-2 genome reveals features
750   important for vaccine design. *bioRxiv*, 2020.03.30.016832.
751   doi:10.1101/2020.03.30.016832.

752    Kunec, D., and Osterrieder, N. (2016). Codon Pair Bias Is a Direct Consequence of
753        Dinucleotide Bias. *Cell Rep.* 14, 55–67. doi:10.1016/j.celrep.2015.12.011.

754    Lanfear, R., Frandsen, P. B., Wright, A. M., Senfeld, T., and Calcott, B. (2016).
755        PartitionFinder 2: New Methods for Selecting Partitioned Models of Evolution for
756        Molecular and Morphological Phylogenetic Analyses. *Mol. Biol. Evol.* 34,
757        msw260. doi:10.1093/molbev/msw260.

758    Lant, J. T., Berg, M. D., Heinemann, I. U., Brandl, C. J., and O'Donoghue, P. (2019).
759        Pathways to disease from natural variations in human cytoplasmic tRNAs. *J. Biol.*
760        *Chem.* 294, 5294–5308. doi:10.1074/jbc.REV118.002982.

761    Lee, S., Weon, S., Lee, S., and Kang, C. (2010). Relative codon adaptation index, a
762        sensitive measure of codon usage bias. *Evol. Bioinforma.* 2010, 47–55.
763        doi:10.4137/EBO.S4608.

764    Li, B., Yang, J., Zhao, F., Zhi, L., Wang, X., Liu, L., et al. (2020). Prevalence and
765        impact of cardiovascular metabolic diseases on COVID-19 in China. *Clin. Res.*
766        *Cardiol.* 109, 531–538. doi:10.1007/s00392-020-01626-9.

767    Li, W., Shi, Z., Yu, M., Ren, W., Smith, C., Epstein, J. H., et al. (2005). Bats are natural
768        reservoirs of SARS-like coronaviruses. *Science (80-. ).* 310, 676–679.
769        doi:10.1126/science.1118391.

770    Lobry, J. R., and Gautier, C. (1994). Hydrophobicity, expressivity and aromaticity are
771        the major trends of amino-acid usage in 999 escherichia coli chromosome-encoded
772        genes. *Nucleic Acids Res.* 22, 3174–3180. doi:10.1093/nar/22.15.3174.

773    Masters, P. S. (2019). Coronavirus genomic RNA packaging. *Virology* 537, 198–207.
774        doi:10.1016/j.virol.2019.08.031.

775    Miller, J. B., Hippen, A. A., Wright, S. M., Morris, C., and Ridge, P. G. (2017). Human
776        viruses have codon usage biases that match highly expressed proteins in the tissues
777        they infect. *Res. Artic. Biomed. Genet. Genomics* 2, 1–5.
778        doi:10.15761/BGG.1000134.

779    Mueller, S., Coleman, J. R., Papamichail, D., Ward, C. B., Nimnual, A., Futcher, B., et

25

780    al. (2010). Live attenuated influenza virus vaccines by computer-aided rational

781    design. *Nat. Biotechnol.* 28, 723–726. doi:10.1038/nbt.1636.

782    Nathan, A. J., and Scobell, A. (2012). How China sees America. *Foreign Aff.* 91, 287.

783    doi:10.1017/CBO9781107415324.004.

784    Nouën, C. Le, Brock, L. G., Luongo, C., McCarty, T., Yang, L., Mehedi, M., et al.

785    (2014). Attenuation of human respiratory syncytial virus by genome-scale codon-

786    pair deoptimization. *Proc. Natl. Acad. Sci. U. S. A.* 111, 13169–13174.

787    doi:10.1073/pnas.1411290111.

788    Novembre, J. A. (2002). Accounting for Background Nucleotide Composition When

789    Measuring Codon Usage Bias. *Mol. Biol. Evol.* 19, 1390–1394.

790    doi:10.1093/oxfordjournals.molbev.a004201.

791    Parker, M. M., and Masters, P. S. (1990). Sequence comparison of the N genes of five

792    strains of the coronavirus mouse hepatitis virus suggests a three domain structure

793    for the nucleocapsid protein. *Virology* 179, 463–468. doi:10.1016/0042-

794    6822(90)90316-J.

795    Peden, J. F. (1999). Analysis of codon usage. *Biosystems.* 5.

796    doi:10.1016/j.biosystems.2011.06.005.

797    Peiris, J. S. M., Guan, Y., and Yuen, K. Y. (2004). Severe acute respiratory syndrome.

798    *Nat. Med.* 10, S88–S97. doi:10.1038/nm1143.

799    Price, M. N., Dehal, P. S., and Arkin, A. P. (2009). Fasttree: Computing large minimum

800    evolution trees with profiles instead of a distance matrix. *Mol. Biol. Evol.* 26,

801    1641–1650. doi:10.1093/molbev/msp077.

802    Rahman, S. U., Yao, X., Li, X., Chen, D., and Tao, S. (2018). Analysis of codon usage

803    bias of Crimean-Congo hemorrhagic fever virus and its adaptation to hosts. *Infect.*

804    *Genet. Evol.* 58, 1–16. doi:10.1016/j.meegid.2017.11.027.

805    Sharp, P. M., and Li, W. H. (1987). The codon adaptation index-a measure of

806    directional synonymous codon usage bias, and its potential applications. *Nucleic*

807    *Acids Res.* 15, 1281–1295. doi:10.1093/nar/15.3.1281.

808    Shen, S. H., Stauft, C. B., Gorbatsevych, O., Song, Y., Ward, C. B., Yurovsky, A., et al.
809        (2015). Large-scale recoding of an arbovirus genome to rebalance its insect versus
810        mammalian preference. *Proc. Natl. Acad. Sci. U. S. A.* 112, 4749–4754.
811        doi:10.1073/pnas.1502864112.

812    Sievers, F., Wilm, A., Dineen, D., Gibson, T. J., Karplus, K., Li, W., et al. (2011). Fast,
813        scalable generation of high□quality protein multiple sequence alignments using
814        Clustal Omega. *Mol. Syst. Biol.* 7, 539. doi:10.1038/msb.2011.75.

815    Simón, D., Fajardo, A., Sóñora, M., Delfraro, A., and Musto, H. (2017). Host influence
816        in the genomic composition of flaviviruses: A multivariate approach. *Biochem.*
817        *Biophys. Res. Commun.* 492, 572–578. doi:10.1016/j.bbrc.2017.06.088.

818    Su, S., Wong, G., Shi, W., Liu, J., Lai, A. C. K., Zhou, J., et al. (2016). Epidemiology,
819        Genetic Recombination, and Pathogenesis of Coronaviruses. *Trends Microbiol.* 24,
820        490–502. doi:10.1016/j.tim.2016.03.003.

821    Sueoka, N. (1988). Directional mutation pressure and neutral molecular evolution. *Proc.*
822        *Natl. Acad. Sci. U. S. A.* 85, 2653–2657. doi:10.1073/pnas.85.8.2653.

823    Suzuki, H., Saito, R., and Tomita, M. (2005). A problem in multivariate analysis of
824        codon usage data and a possible solution. *FEBS Lett.* 579, 6499–6504.
825        doi:10.1016/j.febslet.2005.10.032.

826    Tian, L., Shen, X., Murphy, R. W., and Shen, Y. (2018). The adaptation of codon usage
827        of +ssRNA viruses to their hosts. *Infect. Genet. Evol.* 63, 175–179.
828        doi:10.1016/j.meegid.2018.05.034.

829    Tort, F. L., Castells, M., and Cristina, J. (2020). A comprehensive analysis of genome
830        composition and codon usage patterns of emerging coronaviruses. *Virus Res.* 283.
831        doi:10.1016/j.virusres.2020.197976.

832    Uddin, A., and Chakraborty, S. (2018). Codon Usage Pattern of Genes Involved in
833        Central Nervous System. *Mol. Neurobiol.* doi:10.1007/s12035-018-1173-y.

834    Uddin, A., Choudhury, M. N., and Chakraborty, S. (2017). Factors influencing codon
835        usage of mitochondrial ND1 gene in pisces, aves and mammals. *Mitochondrion* 37,

836        17–26. doi:10.1016/j.mito.2017.06.004.

837   Wang, B., Yang, C., Tekes, G., Mueller, S., Paul, A., Whelan, S. P. J., et al. (2015).
838        Recoding of the vesicular stomatitis virus L gene by computer-aided design
839        provides a live, attenuated vaccine candidate. *MBio* 6. doi:10.1128/mBio.00237-
840        15.

841   Wang, L. F., Shi, Z., Zhang, S., Field, H., Daszak, P., and Eaton, B. T. (2006). Review
842        of    bats    and    SARS.    *Emerg.    Infect.    Dis.*    12,    1834–1840.
843        doi:10.3201/eid1212.060401.

844   Wang, S., Zhou, X., Zhang, T., and Wang, Z. (2020). The need for urogenital tract
845        monitoring in COVID-19. *Nat. Rev. Urol.* doi:10.1038/s41585-020-0319-7.

846   Wong, A., Li, X., Lau, S., and Woo, P. (2019). Global Epidemiology of Bat
847        Coronaviruses. *Viruses* 11, 174. doi:10.3390/v11020174.

848   Wong, S. H., Lui, R. N., and Sung, J. J. (2020). Covid 19 and the digestive system. *J.*
849        *Gastroenterol. Hepatol.* 35, 744–748. doi:10.1111/jgh.15047.

850   Wright, F. (1990). The "effective number of codons" used in a gene. *Gene* 87, 23–29.
851        doi:10.1016/0378-1119(90)90491-9.

852   Wu, Z., Zhang, Z., and Wu, S. (2020). Focus on the Crosstalk Between COVID-19 and
853        Urogenital Systems. *J. Urol.* doi:10.1097/ju.0000000000001068.

854   Yang, C., Skiena, S., Futcher, B., Mueller, S., and Wimmer, E. (2013). Deliberate
855        reduction of hemagglutinin and neuraminidase expression of influenza virus leads
856        to an ultraprotective live vaccine in mice. *Proc. Natl. Acad. Sci. U. S. A.* 110,
857        9481–9486. doi:10.1073/pnas.1307473110.

858   Yin, Y., and Wunderink, R. G. (2018). MERS, SARS and other coronaviruses as causes
859        of pneumonia. *Respirology* 23, 130–137. doi:10.1111/resp.13196.

860   Zaki, A. M., van Boheemen, S., Bestebroer, T. M., Osterhaus, A. D. M. E., and
861        Fouchier, R. A. M. (2012). Isolation of a Novel Coronavirus from a Man with
862        Pneumonia    in    Saudi    Arabia.    *N.    Engl.    J.    Med.*    367,    1814–1820.
863        doi:10.1056/NEJMoa1211721.

864   Zheng, Y. Y., Ma, Y. T., Zhang, J. Y., and Xie, X. (2020). COVID-19 and the
865       cardiovascular system. *Nat. Rev. Cardiol.* 17, 259–260. doi:10.1038/s41569-020-
866       0360-5.

867   Zhou, P., Fan, H., Lan, T., Yang, X. Lou, Shi, W. F., Zhang, W., et al. (2018). Fatal
868       swine acute diarrhoea syndrome caused by an HKU2-related coronavirus of bat
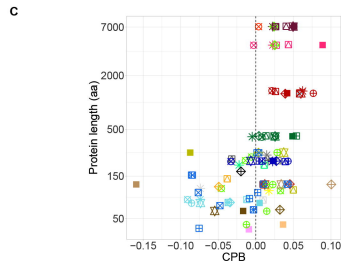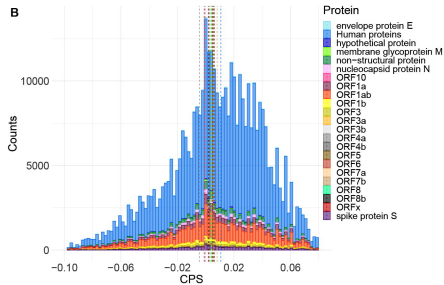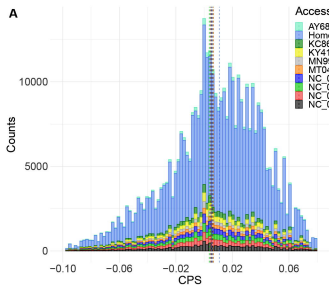869       origin. *Nature* 556, 255–259. doi:10.1038/s41586-018-0010-9.

**A**

**B**

**C**

Proteins

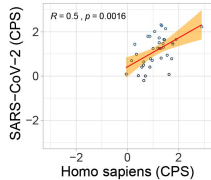| | |
|---|---|
| envelope protein E | ORF3b |
| hypothetical protein | ORF4a |
| membrane glycoprotein M | ORF4b |
| non-structural protein | ORF5 |
| nucleocapsid protein N | ORF6 |
| ORF10 | ORF7a |
| ORF1a | ORF7b |
| ORF1ab | ORF8 |
| ORF1b | ORF8b |
| ORF3 | ORFx |
| ORF3a | spike protein S |

Accession | Host | Virus

AY686864|Civets|SARS
KC869678|Bats|MERS
KY417150|Bats|SARS
MN996532|Bats|SARS-CoV-2
MT040336|Pangolins|SARs-CoV-2
NC_004718|Human|SARS
NC_038294|Human|MERS
NC_039207|Hedgehogs|MERS
NC_045512|Human|SARS-CoV-2

A



B

- Musculoskeletal Diseases
- Neoplasms
- Nervous System Diseases
- Nutritional and Metabolic Diseases
- Occupational Diseases
- Otorhinolaryngologic Diseases
- Pathological Conditions, Signs and Symptoms
- Respiratory Tract Diseases
- Skin and Connective Tissue Diseases
- Stomatognathic Diseases
- Wounds and Injuries
- Behavior and Behavior Mechanisms
- Cardiovascular Diseases
- Chemically-Induced Disorders
- Digestive System Diseases
- Endocrine System Diseases
- Eye Diseases
- Female Urogenital Diseases and Pregnancy Complications
- Hemic and Lymphatic Diseases
- Immune System Diseases
- Infections
- Male Urogenital Diseases
- Mental Disorders