# Predicting Hosts Based on Early SARS-CoV-2 Samples and Analyzing Later World-wide Pandemic in 2020

Qian Guo[1#], Mo Li[1#], Chunhui Wang[1#], Jinyuan Guo[1#], Xiaoqing Jiang[1#], Jie Tan[1],

Shufang Wu[1], Peihong Wang[1], Tingting Xiao[2], Man Zhou[1], Zhencheng Fang[1],

Yonghong Xiao[2*] & Huaiqiu Zhu[1*]

[1] *State Key Laboratory for Turbulence and Complex Systems, Department of*

*Biomedical Engineering, College of Engineering, and Center for Quantitative*

*Biology, and School of life Sciences, Peking University, Beijing 100871, China.*

[2] *State Key Laboratory for Diagnosis and Treatment of Infectious Diseases, National*

*Clinical Research Center for Infectious Diseases, Collaborative Innovation Center for*

*Diagnosis and Treatment of Infectious Diseases, The First Affiliated Hospital, College*

*of Medicine, Zhejiang University, Hangzhou 310006, China.*

[#] These authors contributed equally.

[*] Corresponding authors: hqzhu@pku.edu.cn (Zhu H) & xiao-yonghong@163.com (Xiao Y)

**Running title:** *Guo Q et al / A host prediction algorithm, DeepHoF*

## Abstract

The SARS-CoV-2 pandemic has raised the concern for identifying hosts of the virus since the early-stage outbreak. To address this problem, we proposed a deep learning method, DeepHoF, based on extracting the viral genomic features automatically, to predict host likelihood scores on five host types, including plant, germ, invertebrate, non-human vertebrate and human, for novel viruses. DeepHoF made up for the lack of an accurate tool applicable to any novel virus and overcame the limitation of the sequence similarity-based methods, reaching a satisfactory AUC of 0.987 on the five-classification. Additionally, to fill the gap in the efficient inference of host species for SARS-CoV-2 using existed tools, we conducted a deep analysis on the host likelihood

29   profile calculated by DeepHoF. Using the isolates sequenced in the earliest stage of

30   COVID-19, we inferred minks, bats, dogs and cats were potential hosts of SARS-CoV-

31   2, while minks might be one of the most noteworthy hosts. Several genes of SARS-

32   CoV-2 demonstrated their significance in determining the host range. Furthermore, the

33   large-scale genome analysis, based on DeepHoF's computation for the later world-wide

34   pandemic in 2020, disclosed the uniformity of host range among SARS-CoV-2 samples

35   and the strong association of SARS-CoV-2 between humans and minks.

36   **KEYWORDS:** Host prediction; Deep learning, Mink; SARS-CoV-2; Early stage of

37   pandemic

38

## Introduction

40   The global COVID-19 pandemic caused by severe acute respiratory syndrome

41   coronavirus 2 (SARS-CoV-2) has raised the long-lasting quest for hosts of the virus

42   since the pandemic outbreak, meanwhile the majority view is that the virus probably

43   originated from bats [1]. So far there have been many discussions for the potential hosts

44   despite an initial pointer to *Manis javanica* (pangolins) [2, 3], most of the suppositions

45   were based on the increasing cases of animal infection, such as dogs, cats, tigers, lions,

46   and minks [4, 5], *etc.* Several studies performed experiments to investigate the

47   susceptibility of a limited number of model animals [6-8]. At the same time, some

48   studies attempted to reveal the range of hosts based on analysis of molecular sequence

49   or structural information [9, 10]. For instance, Damas *et al.,* [10] conducted a

50   computational analysis based on host receptor similarity using the angiotensin-

51   converting enzyme 2 (ACE2) protein and evaluated the infection risks for a broad range

52   of animals. As the pandemic spreads, minks, which were even not referred to as high

53   infection animal in above peer-review articles, have been frequently reported massively

54   infected with COVID-19 over the world [5], and were the only known animal reported

55   to transmit SARS-CoV-2 to humans [11, 12]. It is worth mentioning that, in January,

56   2020, we have reported in the form of a preprint archive with predicting minks as a

57   potential host based on the six earliest sequenced SARS-CoV-2 isolates [13]. However,

58    the later complication of pandemic prompts peoples again to have a full review of the

59    issue of host determination for SARS-CoV-2. This raises a new challenge, which is how

60    to implement and improve the capability of computational methods to predict the hosts

61    of a novel virus like SARS-CoV-2, especially when we have relatively small amounts

62    of samples of sequencing viral data at the early stage of the pandemic outbreak. It is

63    certainly constructive for similar pandemic caused by novel viruses in the future.

64         Generally, the host range of viruses is dependent on molecular interactions between

65    viruses and host cells including receptor recognition, adaptions to the host cellular

66    machinery and evading innate immune recognition [14]. Of these, receptor recognition

67    that facilitates the attachment of viruses to the host cells is the most primary step. Thus,

68    the glycoproteins that viruses use to recognize the host receptor as well as the whole

69    genome sequences are widely used in identifying the potential hosts of viruses [1]. To

70    detect the potential host and pathogenicity of novel viruses, the conventional

71    computational methods are almost based on similarity of either virus genome

72    composition or host receptor. Limitations of the both strategies lie in that they assume

73    phylogeny may reflect host association. However, this assumption is untenable from

74    the perspective of epidemiology and evolution. On the one hand, viruses occasionally

75    shift between distantly related host species. On the other hand, owing to the long-term

76    adaptation to the hosts, the viral genomic characteristics acquired from hosts can be

77    quite incompatible with the virus phylogenetic groups [15]. The specificity of

78    recognition between viruses and host species also involves structural information in

79    some key domains of both viral proteins and host receptor proteins, such as the receptor-

80    binding domain, that sequence similarity is insufficient to explain. For example, the

81    civet-specific K479 and S487 residues of SARS-CoV spike glycoprotein can efficiently

82    bind to civet ACE2 but have much less affinity to human ACE2 [16, 17]. This is also

83    the reason that the similarity-based method of host ACE2 proteins sequences fails to

84    predict minks as host of high and very high risk for SARS-CoV-2 infection [10].

85         Until now, several published tools aimed to identify the hosts of viruses exceeded

86    the limitation of sequence-similarity-based strategies by machine learning methods

87  with viral sequences or their genomic traits related to virus-host interactions, such as

88  ViralHostPredictor [15], HostPhinder [18], WIsH [19], Host Taxon Predictor [20], and

89  VIDHOP [21]. While these tools performed well under some conditions, they are

90  actually not considered feasible to be applied to a novel virus without the knowledge

91  of host range, like SARS-CoV-2. HostPhinder and WIsH predict hosts for only

92  bacteriophages and they are inappropriate for non-phage viruses. Host Taxon Predictor

93  focuses on distinguish bacteriophages and eukaryotic viruses. ViralHostPredictor

94  predicts hosts and the existence and identity of arthropod vectors for human-infecting

95  RNA viruses by Gradient boosting machines with the features of selected evolutionary

96  genomic traits and phylogenetic information. It also illustrated the better ability of

97  machine learning methods to predict virus hosts compared to the way of sequence

98  similarity comparison. However, ViralHostPredictor cannot determine whether human

99  is the host of a novel virus. With the utilization of evolutionary signatures,

100 ViralHostPredictor lacks power to predict incidental hosts which do not maintain long-

101 term circulation of new viruses. Moreover, the predictive abilities of the methods above

102 rely on the handcrafted features like codon pair scores, $k$-mer frequencies and amino

103 acid biases, which might neglect other important information encoded in the virus

104 genomes. VIDHOP, a deep-learning-based tool, is designed to predict potential hosts of

105 viruses, but its application was limited into three viral species: influenza A, rabies

106 lyssavirus and rotavirus A.

107      To address the challenge of predicting probable hosts of a novel virus like SARS-

108 CoV-2, we proposed the host prediction algorithm DeepHoF (**Deep** learning-based

109 **Ho**st **F**inder) in the current study. Developed based on BiPath Convolutional Neural

110 Network (BiPathCNN), DeepHoF automatically extracts the genomic features from the

111 input viral sequences. The model finally outputs five host likelihood scores and their $p$-

112 values on five host types, including plant, germ, invertebrate, non-human vertebrate

113 (refers to other vertebrates except humans) and human, where all the living organism

114 hosts are covered. DeepHoF was designed as a five-class classifier containing five

115 independent nodes in the output layer with sigmoid activation and binary cross-entropy

4

116   loss function for each node, corresponding to five independent binary classifications on

117   the five host types individually. DeepHoF made up for the lack of efficient method

118   applicable for any novel virus and significantly outperformed the Basic Local

119   Alignment Search Tool (BLAST)-based strategy with the evidently high AUC of 0.987

120   on the classification of five host types. In January 2020, we have reported the host

121   prediction for six earliest sequenced SARS-CoV-2 isolates employing our algorithm

122   [13]. In this study, we furthered the work using the 17 earliest sampled SARS-CoV-2

123   isolates, which provides essential information in the early epidemic of the virus.

124   DeepHoF evaluated the host likelihood scores on humans and non-human vertebrates

125   for the earliest samples and characterized the isolates with their host likelihood score

126   profiles. As there existed a blank in the inference of host species for SARS-CoV-2 using

127   the tools which were state of the art, we conducted a deep analysis on the host likelihood

128   score profile predicted by DeepHoF to find the detailed hosts, including both reservoirs

129   and susceptible hosts which are not discriminated in this study. We inferred minks, bats,

130   dogs and cats were the probable hosts, while minks maybe one of the most noteworthy

131   hosts. The inference was supported by the infection facts or animal experiments in the

132   later pandemic. Based on our model, several genes of SARS-CoV-2 were further

133   investigated and demonstrated their significance in determining the host likelihood

134   scores on human or the host range for SARS-CoV-2, respectively. With a large-scale

135   genome analysis based on DeepHoF's computation for the later world-wide pandemic,

136   the uniformity of host inference among a large number of SARS-CoV-2 samples was

137   verified, and the association of SARS-CoV-2 between humans and minks was disclosed.

138   Supported by the satisfactory performance on five host type classification and the

139   successful application in SARS-CoV-2, DeepHoF has the capability to provide reliable

140   host information of novel virus, and is expected to narrow the time lag between novel

141   virus discovery and prevention at the early-stage of epidemic prevention.

142

143   **Results**

144   **Performance of the DeepHoF algorithm**

5

145   The DeepHoF algorithm is designed as a five-class classifier using the deep learning

146   method of BiPathCNN (see Methods). Herein five likelihood scores on five host types,

147   including plants, germs, invertebrates, non-human vertebrates, and humans, are

148   calculated by DeepHoF. The host likelihood score profile consisting of five predicted

149   scores, is then analysed in depth to find the specific hosts of a novel virus such as SARS-

150   CoV-2 in this study. As mentioned above, the existed bioinformatics tools [15, 18-21]

151   were not designed to perform the prediction of the host likelihood scores on the five

152   host types for any given virus, and thus cannot be compared with DeepHoF directly.

153   And therefore, we compared the performance of DeepHoF model with BLAST (details

154   of finding host using BLAST are described in Supplementary Methods), adopting six

155   classification metrics: true-positive rate (TPR), false-positive rate (FPR), area under the

156   curve (AUC), precision, accuracy and F1-score. To assess the performance of predicting

157   novel viruses, we used training and test datasets divided in chronological order [22]

158   (Methods). There is no overlap of virus species in training and test sets. With an evident

159   higher AUC of 0.987, DeepHoF can significantly outperform BLAST (with the average

160   AUC of 0.833) as shown in **Figure 1A** and **Table 1** (a detailed comparison on each host

161   type is illustrated in Supplemental Figure S1 and Table S1).

162   In addition, we compared the utility of DeepHoF and a phylogenetic tree to

163   discriminate the human-infecting and non-human-infecting coronaviruses using their

164   whole genome sequences. As shown in Figure 1B (the left), DeepHoF could identify

165   evidently higher probabilities of human-infecting coronaviruses to infect humans (two-

166   sided unpaired Welch Two Sample $t$-test, $p$-value $= 1.732 \times 10^{-10}$). However, the

167   phylogenetic analysis result was not satisfactory owing to the weak homology among

168   the human-infecting coronaviruses, which were scattered around the phylogenetic tree

169   of coronaviruses (Figure 1C). The comparison was similar for the inferences using their

170   spike glycoprotein coding genes (S genes) as shown in Figure 1B (the right), and D

171   (two-sided unpaired Welch Two Sample $t$-test, $p$-value$=3.657 \times 10^{-5}$). This result is

172   nontrivial because S genes are essential in coronavirus-host interaction [23]. Clearly,

6

173 DeepHoF can overcome the limitation of sequence similarity-based method and shows

174 superior predictive ability especially for novel viruses.

**Host prediction of SARS-CoV-2**

176 The accurate prediction of hosts of earliest detected isolates can undoubtedly assist the

177 public health system to take more appropriate preventive measures at the early stage of

178 the pandemic outbreak. In view of this, we focused on the prediction with SARS-CoV-

179 2 isolates sequenced in the earliest stage of COVID-19 detection, which is closer to the

180 most recent common ancestor of SARS-CoV-2. Previous to this paper, we have reported

181 the prediction for the six earliest sequenced SARS-CoV-2 isolates using our algorithm

182 on 21 January, 2020 [13]. In this study, we further strengthened the prediction of hosts

183 of SARS-CoV-2 with all 17 earliest detected isolates (including the six earliest ones)

184 sequenced in December, 2019. Herein we take NC_045512 (complete genome of

185 SARS-CoV-2 isolate, Wuhan-Hu-1, collected on 31 December 2019 in Wuhan, China,

186 and used as the representative genome of SARS-CoV-2 in most studies) as an example

187 to illustrate the workflow of DeepHoF on SARS-CoV-2 isolates (**Figure 2**).

188 For all the 17 SARS-CoV-2 isolates listed in **Figure 3A**, the host likelihood scores

189 on non-human vertebrates and humans were assigned $p$-values less than 0.05 (0.002

190 and 0.027 respectively), illustrating a high possibility of non-human vertebrates and

191 humans (Methods) to be the hosts of SARS-CoV-2. Besides, compared to other

192 coronaviruses released on RefSeq [24], the high similarity of human and non-human

193 vertebrate host likelihood scores among SARS-CoV-2, SARS-CoV and MERS-CoV

194 (Figure 3B), would raise an alarm when the infection capabilities of SARS-CoV-2 was

195 uncertain in the early stage of pandemic.

196 To describe the contribution of each gene in the determination of the host likelihood

197 scores of SARS-CoV-2 isolates (use NC_045512 as a representation), we used each

198 gene sequence of SARS-CoV-2 as the input of DeepHoF and predicted the host

199 likelihood scores for each gene. We found that the S gene, ORF1ab and ORF7b indeed

200 acquired high likelihood scores on human host type and thus playing important roles in

201 determining human as the host (Figure 3C). The fact that several domains on S gene

202 and ORF1ab are essential for the coronavirus-host fusion process, host survival or viral

203 replication [25-27] suggests the rationality of our findings. It is noteworthy that the

204 linear correlation between the lengths and the host likelihood scores for genes is not

205 tenable (Supplemental Figure S2). This shows that the importance of ORF1ab is not

206 due to the remarkable length of the gene. Additionally, our prediction proposes the

207 necessity of further experimental research on the function of ORF7b in SARS-CoV-2.

208 Furthermore, we explored how each gene functioning on coronavirus life circle [25-28]

209 contributed to the human host likelihood scores of SARS-CoV-2, SAR-CoV and

210 MERS-CoV using the earliest sequenced samples, including 12 SARS-CoV isolates, 9

211 MERS-CoV isolates and 17 SARS-CoV-2 isolates released in NCBI in 2003, 2012 and

212 2019, respectively (Supplemental Table S2). The contributions of these genes were

213 represented by their host likelihood scores on human. We found that ORF1ab was

214 relatively important in the prediction for all these viruses, which was possibly due to its

215 functions in viral replication and host survival [27]. The structural genes (S, M, N, and

216 E genes) in these three viruses contributed differently on the human host type,

217 illustrating these genes functioned inconsistently in these viruses. Specifically, S gene,

218 participating in virus-host fusion process, contributed more in SARS-CoV-2 and SARS-

219 CoV, while N gene, eliciting the strong specific antibody responses, played the most

220 important role in MERS-CoV. Two equivalent genes, ORF9b, attaching membrane in

221 virion assembly of SARS-CoV, and ORF8b, related to or immune evasion of MERS-

222 CoV, made high contributions on human host likelihood scores for the two viruses.

223 Moreover, two group-specific genes, ORF7b with unclear function in SARS-CoV, and

224 ORF3 associated with virial replication and pathogenesis in MERS-CoV contributed

225 significantly in the two viruses (Figure 3C, Supplemental Figure S3). These

226 discrepancies might indicate the different significance of these genes among the three

227 coronaviruses in the interaction with human and give hints to the target of drug design.

228 It is disappointed that host determination for SARS-CoV-2 is extremely difficult due to

229 the limited knowledge of the virus world. Therefore, the sequences and host

230 information of viruses contained in the public database should be valued and fully

231 utilized. To fill the gap in the efficient inference of host species for SARS-CoV-2 using

232 the tools which were state of the art, we deeply analyzed the host likelihood profiles of

233 viruses output by DeepHoF to seek specific vertebrate hosts of the early-stage SARS-

234 CoV-2 isolates. In this study, we proposed that viruses with the same host species

235 possessed the host likelihood score profiles close in the five-dimensional space. Based

236 on this assumption, we compared the host likelihood score profile of SARS-CoV-2 with

237 those of the non-human vertebrate viruses released in GenBank [29] before the

238 pandemic outbreak of SARS-CoV-2 (Methods). We found that minks (*Mustela*

239 *lutreola/Neovison vison*) were the most probable host, followed by Chinese rufous

240 horseshoe bats (*Rhinolophus sinicus*), dogs (*Canis lupus familiaris*), Pomona roundleaf

241 bats (*Hipposideros Pomona*) and cat family (*Felidae*) (**Table 2**, Supplemental Table S3).

242 In contrast, minks, Chinese rufous horseshoe bats, dogs and cat family were

243 respectively classified into very low, low or medium groups by Damas *et al.,* [10], who

244 divided 410 vertebrate species into five categories from very high to very low

245 depending on the susceptibility to SARS-CoV-2 based on the analysis of sequence

246 similarity of ACE2 and protein structure of ACE2/SARS-CoV-2 S-binding interface

247 from the vertebrates. In the later world-wide pandemic, it should be pointed out that all

248 the probable hosts we predicted were proved by animal experiments or the infection

249 events [5], which illustrated the usefulness of such analysis for the host inference of

250 SARS-CoV-2. Remarkably, SARS-CoV-2 has been reported largely to infect farmed

251 minks in Netherlands, Denmark, Spain, the United States, Sweden, Italy, Greece,

252 France, Lithuania, Canada, and Poland from April to Febrary, 2021. As of Febrary, 2021,

253 SARS-CoV-2 had been reported to sweep 69 and 207 mink farms in Netherlands and

254 Denmark, respectively, which accelerated the cull of minks and killed the fur industry

255 in the two countries. On 9 October, 2020, at least 10,000 minks were reported dead at

256 Utah and Wisconsin mink farms in the USA, and they were believed infected by SARS-

257 CoV-2 [5] (Table 2).

258 When evaluating the contributions of 11 genes of SARS-CoV-2 in determining

259 mink as the most probable host, we found ORF1ab and ORF8 contributed the most

260   (Supplemental Table S4), which suggesting that genes show different contributions

261   when determining different hosts. The rationality of this result is supported by the roles

262   of ORF1ab in viral replication and host survival [27], and the roles of ORF8 related to

263   immune evasion [30]. However, the interaction between the two genes and the mink

264   cell should merit the further attention and investigation.

265      Additionally, novel coronaviruses, which possess high sequence similarity with

266   SARS-CoV-2, were found on pangolin [2, 3] in China. Even though these pangolin-

267   associated coronaviruses were assigned similar host likelihood score profiles with

268   early-stage SARS-CoV-2 isolates, our analysis demonstrated that the similarity of

269   profiles between SARS-CoV-2 and pangolin-associated coronaviruses was lower than

270   those between SARS-CoV-2 and certain viruses of mink and Chinese rufous horseshoe

271   bat.

**272   Association of SARS-CoV-2 between humans and minks**

273   In April 2020, farmed minks in Netherlands were noticed to be infected by SARS-CoV-

274   2 because of the abnormal mortality [4]. Even though all the mink farms in Netherlands

275   have been screened mandatorily since 28 May 2020, the transmission of coronavirus

276   among the mink population did not seem to cease. Thus, a million farmed minks were

277   culled in Netherlands, and followed by a plan to cull 2.5 million farmed minks in

278   Denmark.

279      Characterizing SARS-CoV-2 by their host likelihood score profiles, we found the

280   isolates detected on humans and minks in Netherlands distributed in a consistent mode,

281   where both groups were divided into a major cluster and a divergence (Figure 3D, 1,746

282   SARS-CoV-2 samples collected from humans in Netherlands as of September 15 and

283   153 SARS-CoV-2 samples collected from farmed minks in Netherlands as of October

284   15 were used respectively, Methods). For SARS-CoV-2, as the host likelihood score on

285   susceptible hosts such as human and mink can also indicate the likelihood to infect

286   these animals, the mode of host likelihood score profile can reflect its property of viral

287   infection. Consequently, the consistency mentioned above hinted the close infection-

288     related behaviors of SARS-CoV-2 on humans and minks in Netherlands and thus

289     illustrated the association of SARS-CoV-2 isolates collected from the two populations.

290     Furthermore, nine of 14 high-frequency variants in human-derived SARS-CoV-2

291     genomes sequenced in Netherlands were absent in the genomes detected in other

292     countries. Herein we used NC_045512 as the reference for variant calling, regarded the

293     variants with ≥5% frequency as high-frequency ones and filtered out the synonymous

294     single nucleotide polymorphisms (SNPs) (Supplemental Table S5). Among these

295     unique high-frequency variants in Dutch human-derived SARS-CoV-2, two were found

296     in Dutch mink-derived SARS-CoV-2, thus proved the circulation of SARS-CoV-2

297     between humans and minks in Netherlands. It was remarkable that our findings could

298     be supported by the conclusions from a research team in Netherland, who utilized more

299     detailed information about patients and related mink farms [12]. In the 2020 world-

300     wide pandemic, minks are the only animal that has been reported to transmit SARS-

301     CoV-2 to humans [11, 12]. We further compared the high-frequency variants of SARS-

302     CoV-2 isolates in humans and minks in Netherlands. Except for four common variants,

303     SARS-CoV-2 isolates derived from minks still had 23 unique high-frequency variants

304     and six were found on S protein that is related to virus-host fusion process. This result

305     indicated that the virus might have gained higher diversity after the intra-species

306     circulation among mink herd and inter-species circulation between minks and human.

307     As the mink infections are expanding worldwide, the association and circulation of

308     SARS-CoV-2 between humans and minks in Netherlands notifies us of the importance

309     to take precautions of the bidirectional transmission in other regions.

310     **Retrospective analysis of the world-wide pandemic**

311     To verify the stability and uniformity of the host inference among SARS-CoV-2

312     samples, retrospective analysis of more isolates in the lasting pandemic was required.

313     As the surge in variants of SARS-CoV-2 complicated the host prediction of the novel

314     virus, we utilized 102,804 SARS-CoV-2 genomes released on GISAID EpiCoV

315     Database (https://www.gisaid.org/) [31] as of 15 September 2020, before the rapid

316     accumulation of mutations in SARS-CoV-2. We picked out 53,759 genomes which met

317 the quality standard given by Chinese Academy of Sciences [32] and trimmed their

318 varied-length 5′- and 3′-untranslated regions (UTR) based on the annotation of

319 NC_045512 (Methods). We calculated the host likelihood score profiles of the 53,759

320 isolates (Supplemental Table S5) and conducted principal component analysis (PCA)

321 on the profiles. As shown in **Figure 4A**, we found a clear cluster of all SARS-CoV-2

322 isolates with 17 earliest ones locating in the center. The kernel density estimation curves

323 displayed on the first two principal components were approximately normally

324 distributed. As the profiles of the 53,759 isolates are under the normal distribution

325 mentioned above, the host range of SARS-CoV-2 isolates keep consistent throughout

326 the pandemic and it is therefore reasonable that the validity of the host inference using

327 the earliest 17 isolates would be efficient in the later pandemic.

328 However, when the SARS-CoV-2 isolates were divided chronologically using 15

329 April 2020 as the split date, which divided 53,759 isolates into two parts more evenly

330 than other dates, we found that the two subsets have divergent distributions in each of

331 the two dimensions of PCA (two-sided two-sample Kolmogorov-Smirnov test, $p$-value

332 $= 0$, $n_{isolates} = 26,167$ before 15 April 2020 and 27,592 after 15 April 2020) (Figure 4B).

333 The approximately normal distribution of SARS-CoV-2 genomes and their time-

334 dependent feature indicate the overall consistency and a certain extent of divergence in

335 the host likelihood score profiles of SARS-CoV-2 isolates.

336 To explain the divergence among host likelihood score profiles, we identified all

337 variants in 53,759 genomes (Supplemental Table S5). The 13 high-frequency variants

338 were located on S gene, N gene, ORF1ab, ORF8 and ORF3a, some of which are related

339 to virus-host fusion process [22, 33]. Furthermore, we annotated our PCA result with

340 the GISAID nomenclature system [31] which divides all SARS-CoV-2 genomes into

341 six major clades based on marker variants that appeared over time. Most of the marker

342 variants were recognized as high-frequency variants in the variant calling. As we can

343 see in Figure 4C, SARS-CoV-2 isolates fell into several clear fusiform clusters

344 according to their clades. This indicated that those marker variants might explain the

345 divergence among host likelihood score profiles. When we manually mutated the 17

12

346 earliest sequenced genomes with those marker variants, we found the variants marking

347 each clade drove the earliest sequenced SARS-CoV-2 to the corresponding cluster of

348 the clade (Supplemental Figure S4), which verified our previous speculation and

349 demonstrated the efficacy of DeepHoF to identify the important variants emerging in

350 the virus's evolution. However, as the consistency of the distribution of host likelihood

351 score profiles were not disturbed, it hinted that these mutations did not change the host

352 range of SARS-CoV-2.

353 Furthermore, to explore the trend of host likelihood of the SARS-CoV-2 over time,

354 we finally examined the relationships between sampling time and the host likelihood

355 scores on non-human vertebrates and humans (Figure 4D). We found that both scores

356 gradually descended. As the host likelihood scores on susceptible hosts also indicate

357 the likelihood to be infected by SARS-CoV-2 from a computational point of view, the

358 trends might indicate the gradually descending infectiousness to human and other

359 vertebrates from the outbreak to 15 September 2020. Those trends may not be so

360 pronounced, but they should arouse our attention.

361

## Discussion

363 In summary, we proposed a deep learning method, DeepHoF, based on extracting the

364 viral genomic features, to calculate the host likelihood scores on five host types.

365 DeepHoF made up for the vacancy of a universal tool feasible to any novel virus. For

366 the identification of five host types, our model can significantly outperform BLAST

367 and well discriminate the human-infecting and non-human-infecting viruses like

368 coronaviruses. Overcoming the limitation of sequence similarity-based methods to

369 disclose the host information of novel viruses, DeepHoF demonstrated the practicality

370 to SARS-CoV-2 in the 2020 pandemic. Using 17 SARS-CoV-2 isolates sequenced in

371 the earliest stage of COVID-19 detection, DeepHoF evaluated the host likelihood

372 scores on humans and non-human vertebrates for SARS-CoV-2. Filling the gap in

373 predicting the host species for any novel virus that remained unsolved using the tools

374 which were state of the art, we further analyzed the host likelihood score profile to

13

375    further infer the specific hosts of SARS-CoV-2. The hosts determined by DeepHoF can

376    be either reservoirs or susceptible middle hosts, which are not discriminated in this

377    study. We found minks, bats, dogs and cats could be potential hosts of SARS-CoV-2,

378    while minks might be one of the most noteworthy animal hosts. Due to mutations, the

379    host likelihood score profiles of the isolates in the long period of the later pandemic had

380    slightly varied, but followed normal distribution where those of the early 17 isolates

381    locate in the center. As a consequence, the host range inferred with the profiles of the

382    isolates during the pandemic was consistent with the inference using the early samples.

383    Additionally, based on the model, we further found three genes (S gene, ORF7b and

384    ORF1ab) and two genes (ORF1ab and ORF8) were significant in determining the host

385    likelihood score on human and the host range for SARS-CoV-2, respectively. The genes

386    involving virus-host fusion process (S gene), viral replication (ORF1ab) and host

387    survival (ORF1ab) played a significant role in determining human as the host, while

388    the genes related to viral replication (ORF1ab), host survival (ORF1ab) and immune

389    evasion (ORF8) were significant to determine the host range for SARS-CoV-2. For the

390    prevention and control of a novel epidemic disease such as COVID-19, the prediction

391    of probable hosts is essential at the early stage of the epidemic outbreak. In view of this,

392    our study is expected to play a potentially effective role in support of those efforts.

393       Furthermore, according to the analysis results of host likelihood score profiles of

394    humans and minks in Netherlands, we found a strong association of SARS-CoV-2

395    isolates collected from the two populations and disclosed the contribution of mink on

396    higher divergence in SARS-CoV-2. The phenomenon coincided with the analysis result

397    of variant calling and could be explained by characteristics of minks in virus circulation.

398    As reported by previous studies about avian-derived influenza A virus, minks serve as

399    a significant node in the viral transmission network, connecting animals from different

400    families and acting as domesticators for viral adaptation to mammals [34]. As the only

401    one animal that has been reported to transmit SARS-CoV-2 to humans, the role of minks

402    in the evolution of SARS-CoV-2 should be studied in depth. Therefore, with a large-

403    scale genome analysis based on DeepHoF's computation for the later world-wide

14

404  pandemic, it should not be slighted for the relationship of SARS-CoV-2 between

405  humans and minks.

406  Although we have applied DeepHoF to SARS-CoV-2 in the current study, the

407  application of DeepHoF is not limited to this virus. DeepHoF is also feasible to

408  determine the host ranges for many other novel viruses, such as the small circular rep-

409  encoding ssDNA viruses newly discovered on wild animals and domestic animals or in

410  the environment. However, limitations of DeepHoF lie in that it does not consider the

411  host sequence information, which can be improved in the future. DeepHoF also does

412  not discriminate between reservoir hosts, vector hosts and other susceptible hosts.

413  Meanwhile, the present study is expected to be further confirmed with both the ongoing

414  events of pandemic and additional experimental findings, and the interpretation of our

415  analysis should be still kept a certain caution.

416  Represented by SARS-CoV-2, more complex and larger numbers of viral genome

417  data will be produced in similar epidemics in the future. In addition, the metagenome

418  and the metavirome can also be used in the prevention and control of the epidemic. The

419  United States Agency for International Development launched the Global Virus

420  Program in 2018 to reduce possible epidemiological threats by studying metaviromic

421  samples from more than 35 countries around the world [35]. It is estimated that there

422  are about 1.67 million novel viruses in mammals, birds and other important hosts of

423  zoonotic viruses. Among them, 631,000-827,000 have the potential to cause zoonotic

424  diseases [35]. However, only 263 viruses from 25 virus families have been confirmed

425  to infect humans [36]. Newly emerged infectious viruses keep threatening our health

426  and well-being. Under the circumstances, using computational methods to discover

427  pathogenetic viruses and acquire knowledge, including the host range, about novel

428  viruses can provide timely response in the prevention of epidemics and pandemics. In

429  the future, the detection of novel viruses will rely more heavily on high-throughput

430  sequencing technologies such as metagenomics and metaviromics. Thus, more robust

431  tools designed for metagenomes and metaviromes are required.

432

## Materials and Methods

### Datasets construction for training and test

We downloaded 63,049 whole viral genomes from GenBank by 9 July, 2019, and tagged them with five host labels (plant, germ, invertebrate, non-human vertebrate and human), which were integrated from the host metadata provided by GenBank (Supplemental Table S6). The five host types covered all the living organism hosts. For viruses infecting multiple host types, multiple labels were given. Following the data collection procedure, short fragments were generated randomly from those tagged whole genomes because of the computational cost in long sequence processing. The training set was constructed with short fragments from 55,283 genomes released before 1 January, 2018, and the test set was constructed with the rest (the Accession list and the host information of the genomes used for training and test are in Supplemental Table S7). There is non-overlap of virus species in the training and test sets.

### Mathematical representation of viral whole genomes

Due to the long-term adaptation to natural reservoirs, viruses share some evolutionary signatures in nucleotide sequences, such as codon pair, dinucleotide, codon, and amino acid biases, with their natural reservoirs [15]. Besides, viral proteins, especially the receptors that are effectively attached to the host cell membrane, are crucial factors for viruses to invade and infect the host cells [37]. In brief, the genome compositions of viruses can inform host-virus correlation.

Herein, we represent a given viral sequence with a base one-hot matrix (BOH) and a codon one-hot matrix (COH), digitizing the genetic information of the virus on nucleotide and codon level respectively. To start with, bases and codons are encoded with one-hot format to work with deep learning algorithms. In the coding of BOH, each consecutive base of a query sequence linked by its complementary strand is encoded by one-hot. For COH, we do not extract ORFs since coding sequences make up most of the viral genome. Instead, we directly concatenate the six phases of the input sequence (Supplemental Figure S5), and then each consecutive codon of the joined sequences is encoded by one-hot. Consequently, for an input sequence of length L, it

16

462    will be transformed to a BOH matrix, with the size of 2L×4, and a COH matrix, with

463    the size of 2L×64.

**BiPathCNN Model descriptions**

465    In building the framework of DeepHoF, we firstly utilize a BiPathCNN [38], containing

466    two CNN paths, digging information from the BOH matrix and COH matrix

467    respectively. The information is naturally corresponding to the viral genomic features

468    for the viruses which infect the same kind of hosts. After independent convolution and

469    pooling operations at the beginning, the two paths are combined by a concatenation

470    layer. Following a normalization layer, five prediction scores will be provided by five

471    sub-paths, containing five independent nodes, corresponding to five independent binary

472    classifications on plant, germ, invertebrate, non-human vertebrate and human

473    individually, in the output layer with sigmoid activation and binary cross-entropy loss

474    function for each node. The architecture of DeepHoF is shown in Supplemental Figure

475    S6 and the details of each layer in BiPathCNN are described in Supplementary

476    Information.

**Implementation of DeepHoF**

478    In the practical application, viral nucleotide sequence is the only input required by

479    DeepHoF. For a viral whole genome sequence (or a partial genome sequence), a cut

480    window moves along the long sequence without overlapping to separate it into suitable

481    fragments for the pre-trained BiPathCNN model. DeepHoF firstly predicts the host

482    infection scores for each fragment. Then it calculates the final score by weighting and

483    summing the predicted scores of each fragment. For example, a 2,000 bp query

484    sequence is separated into three consecutive fragments, corresponding to the first 800

485    bp, the middle 800 bp and the last 400 bp of the query sequence. Then DeepHoF

486    predicts the three fragments independently and calculates the weighted average of the

487    three predicted score vectors with the weights of 800/2,000, 800/2,000, and 400/2,000

488    respectively. For each input sequence, DeepHoF outputs five scores on five host types,

489    respectively. Besides, DeepHoF provides the *p*-values of each score, statistically

490    measuring of how distinct the scores are compared with those of non-infectious viruses

491 [22]. For example, if an input virus has a score of 0.4 on human, we compare 0.4 with

492 the scores of non-human viruses in our dataset and provide the *p*-value as a judgment

493 basis. If the *p*-value is less than 0.05, we conclude that human is the probable host of

494 the input virus with a significantly higher score on human host type than non-human

495 viruses.

496 As the host likelihood score profile of a virus, consisting of the five predicted scores

497 given by DeepHoF, can be regarded as a host-related feature vector extracted by

498 DeepHoF, we utilize it to characterize the virus. It is logistical to regard the viruses with

499 the same host species possess the similar host likelihood score profiles. Based on this

500 assumption, the potential host species of a virus can be inferred by the analysis of the

501 profiles. To quantitatively compare host likelihood score profiles between viruses, we

502 calculated the Euclidean distance between the profiles. In the case of SARS-CoV-2, we

503 searched the detailed vertebrate host of the earliest detected isolates, which are closer

504 to the most recent common ancestor of SARS-CoV-2. To start with, we added the host

505 annotations provided by Virus-Host DB [39] to the vertebrate viruses included in

506 GenBank. Here, the average of host likelihood score profiles of 17 earliest sequenced

507 isolates was used as the representation of SARS-CoV-2. We calculated the Euclidean

508 distance between the profile of SARS-CoV-2 and that of each non-human vertebrate

509 virus (discovered before the outbreak of SARS-CoV-2). We regarded the vertebrate

510 infected by a virus possessing profile close to that of SARS-CoV-2 was the probable

511 host of SARS-CoV-2.

512 **Data filtering and trimming for SARS-CoV-2 genome sequences**

513 There were 102,804 SARS-CoV-2 genomes released on GISAID EpiCoV Database as

514 of 15th September 2020. We downloaded all the sequences and filtered them with the

515 quality standard given by the Chinese Academy of Sciences [32]. Because the UTRs

516 were not taken as seriously as the protein-coding regions and the lengths of sequenced

517 UTRs varied a lot in different SARS-CoV-2 genomes, we trimmed the 5′- and 3′- UTR

518 according to the annotation of NC_045512 to get rid of noises. Thus, we finally got

519 53,759 clean sequences.

**Phylogenetic analysis and single nucleotide polymorphisms analysis**

In this study, we applied Clustal Omega software [40] (version 1.2.4) for multiple sequence alignment and RAxML software [41] (version 8.2.12) for phylogenetic tree building using maximum likelihood methods with 1000 bootstrap replicates. Snippy [42] (version 4.4.3) was utilized for variant calling, using NC_045512 as the reference genome. In this study, we filtered out the synonymous SNPs and regarded the variants with $\geq$ 5% frequency as high-frequency ones. Commands of the three tools are included in Supplementary Information.

## Authors' contributions

HQZ and YHX co-supervised the study. QG, ML, CHW, JYG, XQJ and HQZ developed the DeepHoF model, conducted the analyses and wrote the manuscript. HQZ and ZCF helped with designing the model. MZ calculated performance metrics of DeepHoF and BLAST. PHW helped with phylogeny analysis and SNP analysis. JT, SFW and TTX made plots and table for the results. All authors read and approved the final manuscript.

## Competing interests

The authors have declared no competing interests.

## Acknowledgements

## References

19

548 [1] Zhou P, Yang X-L, Wang X-G, Hu B, Zhang L, Zhang W, et al. A pneumonia
549 outbreak associated with a new coronavirus of probable bat origin. nature
550 2020;579:270-3.

551 [2] Lam TT-Y, Jia N, Zhang Y-W, Shum MH-H, Jiang J-F, Zhu H-C, et al. Identifying
552 SARS-CoV-2-related coronaviruses in Malayan pangolins. Nature 2020;583:282-5.

553 [3] Xiao K, Zhai J, Feng Y, Zhou N, Zhang X, Zou J-J, et al. Isolation of SARS-CoV-
554 2-related coronavirus from Malayan pangolins. Nature 2020;583:286-9.

555 [4] Oreshkova N, Molenaar RJ, Vreman S, Harders F, Munnink BBO, Hakze-van Der
556 Honing RW, et al. SARS-CoV-2 infection in farmed minks, the Netherlands, April and
557 May 2020. Eurosurveillance 2020;25:2001005.

558 [5] OIE. COVID-19 Portal: Events in Animals. https://www.oie.int/en/scientific-
559 expertise/specific-information-and-recommendations/questions-and-answers-on-
560 2019novel-coronavirus/events-in-animals/ (Oct 25 2020, date last accessed).

561 [6] Shi J, Wen Z, Zhong G, Yang H, Wang C, Huang B, et al. Susceptibility of ferrets,
562 cats, dogs, and other domesticated animals to SARS–coronavirus 2. Science
563 2020;368:1016-20.

564 [7] Sia SF, Yan L-M, Chin AW, Fung K, Choy K-T, Wong AY, et al. Pathogenesis and
565 transmission of SARS-CoV-2 in golden hamsters. Nature 2020;583:834-8.

566 [8] Munster VJ, Feldmann F, Williamson BN, Van Doremalen N, Pérez-Pérez L, Schulz
567 J, et al. Respiratory disease in rhesus macaques inoculated with SARS-CoV-2. Nature
568 2020;585:268-72.

569 [9] Santini JM, Edwards SJ. Host range of SARS-CoV-2 and implications for public
570 health. The Lancet Microbe 2020;1:e141-e2.

571 [10] Damas J, Hughes GM, Keough KC, Painter CA, Persky NS, Corbo M, et al. Broad
572 host range of SARS-CoV-2 predicted by comparative and structural analysis of ACE2
573 in vertebrates. Proceedings of the National Academy of Sciences 2020;117:22311-22.

574 [11] Mallapaty S. What's the risk that animals will spread the coronavirus. Nature 2020.

575 [12] Munnink BBO, Sikkema RS, Nieuwenhuijse DF, Molenaar RJ, Munger E,

576 Molenkamp R, et al. Transmission of SARS-CoV-2 on mink farms between humans

577 and mink and back to humans. Science 2021;371:172-7.

578 [13] Guo Q, Li M, Wang C, Wang P, Fang Z, tan J, et al. Host and infectivity prediction

579 of Wuhan 2019 novel coronavirus using deep learning algorithm. bioRxiv

580 2020:2020.01.21.914044.

581 [14] Rothenburg S, Brennan G. Species-specific host–virus interactions: implications

582 for viral host range and virulence. Trends in microbiology 2020;28:46-56.

583 [15] Babayan SA, Orton RJ, Streicker DG. Predicting reservoir hosts and arthropod

584 vectors from evolutionary signatures in RNA virus genomes. Science 2018;362:577-80.

585 [16] Lu G, Wang Q, Gao GF. Bat-to-human: spike features determining 'host jump'of

586 coronaviruses SARS-CoV, MERS-CoV, and beyond. Trends in microbiology

587 2015;23:468-78.

588 [17] Li W, Zhang C, Sui J, Kuhn JH, Moore MJ, Luo S, et al. Receptor and viral

589 determinants of SARS-coronavirus adaptation to human ACE2. The EMBO journal

590 2005;24:1634-43.

591 [18] Villarroel J, Kleinheinz KA, Jurtz VI, Zschach H, Lund O, Nielsen M, et al.

592 HostPhinder: a phage host prediction tool. Viruses 2016;8:116.

593 [19] Galiez C, Siebert M, Enault F, Vincent J, Söding J. WIsH: who is the host?

594 Predicting prokaryotic hosts from metagenomic phage contigs. Bioinformatics

595 2017;33:3113-4.

596 [20] Gałan W, Bąk M, Jakubowska M. Host taxon predictor-a tool for predicting taxon

597 of the host of a newly discovered virus. Scientific reports 2019;9:1-13.

598 [21] Mock F, Viehweger A, Barth E, Marz M. VIDHOP, viral host prediction with deep

599 learning. Bioinformatics 2020.

600 [22] Ren J, Ahlgren NA, Lu YY, Fuhrman JA, Sun F. VirFinder: a novel k-mer based

601 tool for identifying viral sequences from assembled metagenomic data. Microbiome

602 2017;5:1-20.

603    [23] Belouzard S, Millet JK, Licitra BN, Whittaker GR. Mechanisms of coronavirus

604    cell entry mediated by the viral spike protein. Viruses 2012;4:1011-33.

605    [24] Brister JR, Ako-Adjei D, Bao Y, Blinkova O. NCBI viral genomes resource.

606    Nucleic acids research 2015;43:D571-D7.

607    [25] Li Y-H, Hu C-Y, Wu N-P, Yao H-P, Li L-J. Molecular characteristics, functions,

608    and related pathogenicity of MERS-CoV proteins. Engineering 2019;5:940-7.

609    [26] Cheng VC, Lau SK, Woo PC, Yuen KY. Severe acute respiratory syndrome

610    coronavirus as an agent of emerging and reemerging infection. Clinical microbiology

611    reviews 2007;20:660-94.

612    [27] Hu B, Guo H, Zhou P, Shi Z-L. Characteristics of SARS-CoV-2 and COVID-19.

613    Nature Reviews Microbiology 2020:1-14.

614    [28] Wong L-YR, Ye Z-W, Lui P-Y, Zheng X, Yuan S, Zhu L, et al. Middle East

615    Respiratory Syndrome Coronavirus ORF8b Accessory Protein Suppresses Type I IFN

616    Expression by Impeding HSP70-Dependent Activation of IRF3 Kinase IKKε. The

617    Journal of Immunology 2020;205:1564-79.

618    [29] Sayers EW, Cavanaugh M, Clark K, Pruitt KD, Schoch CL, Sherry ST, et al.

619    GenBank. Nucleic acids research 2021;49:D92-D6.

620    [30] Young BE, Fong S-W, Chan Y-H, Mak T-M, Ang LW, Anderson DE, et al. Effects

621    of a major deletion in the SARS-CoV-2 genome on the severity of infection and the

622    inflammatory response: an observational cohort study. The Lancet 2020;396:603-11.

623    [31] Shu Y, McCauley J. GISAID: Global initiative on sharing all influenza data–from

624    vision to reality. Eurosurveillance 2017;22:30494.

625    [32] Zhao W-M, Song S-H, Chen M-L, Zou D, Ma L-N, Ma Y-K, et al. The 2019 novel

626    coronavirus resource. Yi chuan = Hereditas 2020;42:212-21.

627    [33] Lan J, Ge J, Yu J, Shan S, Zhou H, Fan S, et al. Structure of the SARS-CoV-2 spike

628    receptor-binding domain bound to the ACE2 receptor. Nature 2020;581:215-20.

629    [34] Xue R, Tian Y, Hou T, Bao D, Chen H, Teng Q, et al. H9N2 influenza virus isolated

630    from minks has enhanced virulence in mice. Transboundary and emerging diseases

631    2018;65:904-10.

632 [35] Carroll D, Daszak P, Wolfe ND, Gao GF, Morel CM, Morzaria S, et al. The global

633 virome project. Science 2018;359:872-4.

634 [36] King AM, Lefkowitz E, Adams MJ, Carstens EB. Virus taxonomy: ninth report of

635 the International Committee on Taxonomy of Viruses. Elsevier, 2011.

636 [37] Dimitrov DS. Virus entry: molecular mechanisms and biomedical applications.

637 Nature Reviews Microbiology 2004;2:109-22.

638 [38] Fang Z, Tan J, Wu S, Li M, Xu C, Xie Z, et al. PPR-Meta: a tool for identifying

639 phages and plasmids from metagenomic fragments using deep learning. Gigascience

640 2019;8:giz066.

641 [39] Mihara T, Nishimura Y, Shimizu Y, Nishiyama H, Yoshikawa G, Uehara H, et al.

642 Linking virus genomes with host taxonomy. Viruses 2016;8:66.

643 [40] Sievers F, Wilm A, Dineen D, Gibson TJ, Karplus K, Li W, et al. Fast, scalable

644 generation of high-quality protein multiple sequence alignments using Clustal Omega.

645 Molecular systems biology 2011;7:539.

646 [41] Stamatakis A. RAxML-VI-HPC: maximum likelihood-based phylogenetic

647 analyses with thousands of taxa and mixed models. Bioinformatics 2006;22:2688-90.

648 [42] Seemann T. Snippy: rapid bacterial SNP calling and core genome alignments.

649 https://github.com/tseemann/snippy.git (Oct 25 2020, date last accessed).

650 [43] Letunic I, Bork P. Interactive Tree Of Life (iTOL) v4: recent updates and new

651 developments. Nucleic acids research 2019;47:W256-W9.

652

653 **Figure legends**

654 **Figure 1    DeepHoF outperforms BLAST and well learns the information of virus**

655 **hosts**

656 **A**. Average ROC curves and AUC values of DeepHoF and BLAST. DeepHoF performs

657 better than BLAST on average AUC of five host types. **B**. Comparison of host

658 likelihood scores predicted by DeepHoF between human-infecting and non-human-

659 infecting coronaviruses on human. The former performed higher probabilities than the

660 latter (two-sided unpaired Welch Two Sample $t$-test, $t_{(43.843)} = 8.265$ and $t_{(38.016)} = 4.674$,

661     *p*-values = $1.732 \times 10^{-10}$ and $3.657 \times 10^{-5}$. *** *p*-value $<$ 0.0001, *t*-values and degrees of

662     freedom were presented as $t_{(df)}$). **C.** Phylogenetic analyses of whole genomes of

663     coronaviruses. **D**. Phylogenetic analyses of S genes of coronaviruses. Maximum-

664     likelihood phylogenic trees were built by RAxML [41] with 1,000 bootstrap replicates

665     and visualized with iTOL [43]. The whole genomes and the S genes of the human-

666     infecting coronaviruses could not be distinguished from the non-human-infecting ones.

667     (Red: human-infecting coronaviruses; Blue: non-human-infecting coronaviruses).

668     **Figure 2    The workflow of application of DeepHoF on NC_045512**

669     In the application of DeepHoF on SARS-CoV-2 NC_045512, the whole genome of

670     NC_045512 was the only input required by the pre-trained DeepHoF model and coded

671     into BOH and COH matrix for BiPathCNN network. DeepHoF output the host

672     likelihood scores of NC_045512 on five host types respectively and the corresponding

673     significance. The hosts of NC_045512 were predicted to be non-human vertebrates and

674     humans with *p*-values less than 0.05. Simultaneously, NC_045512 was characterized

675     by its host likelihood score profile. Susceptible to viruses with similar profile, *Mustela*

676     *lutreola/ Neovison vison, Rhinolophus sinicus, Canis lupus familiaris, Hipposideros*

677     *pomona* and Feline were output as the probable hosts of NC_045512. BOH: base one-

678     hot matrix, COH:    codon one-hot matrix.

679     **Figure 3    Evaluation of host likelihood scores of SARS-CoV-2**

680     The contribution of each gene in the prediction and the visualization of host likelihood

681     score profiles of SARS-CoV-2 isolates sampled in Netherlands. **A**. Host likelihood

682     scores of 17 earliest detected SARS-CoV-2 isolates and other coronaviruses on humans

683     and non-human vertebrates. SARS-CoV-2 showed high host likelihood scores on both

684     humans and non-human vertebrates with *p*-values less than 0.05. In addition, SARS-

685     CoV-2 was predicted lower score than SARS-CoV and comparable score to MERS-

686     CoV on human. As for host likelihood scores on non-human vertebrates, SARS-CoV-

687     2, SARS-CoV and MERS-CoV were close to each other. Host likelihood scores have

688     *p*-values less than 0.05 are marked 'Y (yes)'. (Red: human-infecting coronaviruses; *:

689     the 17 earliest collected SARS-CoV-2 isolates). **B**. Hierarchical clustering of early-

24

690     stage SARS-CoV-2 and other coronaviruses using five-dimensional host likelihood

691     score profiles given by DeepHoF. The profile of SARS-CoV-2 was close to that of

692     SARS-CoV and MERS-CoV (Red: SARS-CoV-2; Blue: SARS-CoV; Yellow: MERS-

693     CoV). **C**. Contributions of the protein coding genes on determining the host likelihood

694     scores of SARS-CoV-2, SARS-CoV and MERS-CoV on human. The structural genes,

695     ORF1ab and group-specific genes contributed differently in the three coronaviruses

696     (two-sided unpaired Welch Two Sample t-test, *p*-value < 0.05, see in Supplemental

697     Figure S3). S, ORF7b and ORF1ab were the most pivotal in SARS-CoV-2. ORF7b,

698     ORF9b and S were the most considerable in SARS-CoV. ORF8b, N and ORF3

699     contributed the most in MERS-CoV (S: spike glycoprotein coding gene; M:

700     membrane/matrix glycoprotein coding gene; N: nucleocapsid phosphoprotein coding

701     gene; E: envelope coding gene). **D**. Principal component analysis (PCA) of host

702     likelihood score profiles of SARS-CoV-2 detected on humans and minks in Netherlands.

703     The host likelihood score profiles of mink-derived and human-derived SARS-CoV-2

704     isolates in Netherlands are distributed in a consistent mode, containing a major cluster

705     and divergence. The host likelihood score profiles of human-derived (left) and mink-

706     derived (right) SARS-CoV-2 isolates in Netherlands distributed in a consistent mode,

707     both containing a major cluster (red) and divergence (blue). The major cluster and the

708     divergence were divided by the pam function of R package cluster.

709     **Figure 4**    **Entirety and divergence in the host likelihood score profiles of 53,759**

710     **SARS-CoV-2 isolates in the later world-wide pandemic**

711     **A**. PCA of host likelihood score profiles of 53,759 SARS-CoV-2 isolates and the

712     distribution on each principal component. All the host s likelihood core profiles of

713     53,759 SARS-CoV-2 isolates were clustered with 17 earliest sequenced isolates located

714     in the center and the density curves displayed on each principal component were

715     approximate normal distribution. **B**. Distributions of host likelihood score profiles of

716     53,759 SARS-CoV-2 isolates collected before and after 15 April 2020. When the

717     SARS-CoV-2 isolates were divided chronologically using 15 April 2020 as the split

718     date, which divided the 53,759 isolates into two parts more evenly than other dates. The

719   host likelihood score profiles of SARS-CoV-2 before and after 15 April 2020 had

720   divergent distributions on each principal component (two-sided two-sample

721   Kolmogorov-Smirnov test, $p$-value = 0, $n_{isolates}$ = 26,167 before 15 April 2020 and

722   27,592 after 15 April 2020. Blue, 26,167 isolates collected before 15 April 2020; Red,

723   27,592 isolates collected after 15 April 2020; Grey, all the 53,759 isolates). **C**. GISAID

724   clades represented in PCA of host likelihood score profiles of 53,759 SARS-CoV-2

725   genomes. All the 53,759 samples representing 53,759 host likelihood score profiles

726   were painted with six different colours corresponding to six different GISAID clades

727   of SARS-CoV-2. SARS-CoV-2 isolates fell into several clear fusiform clusters with

728   different colours according to their clades. **D**. Time series of the host likelihood scores

729   on humans and non-human vertebrates for SARS-CoV-2 in the later world-wide

730   pandemic. The host likelihood scores on humans and non-human vertebrates descend

731   gradually with time (linear regression model analysis, $R$-squared = $6.806 \times 10^{-3}$ and

732   $1.431 \times 10^{-2}$, $t_{(53,757)} = -19.22$ and $t_{(53,757)} = -27.96$, $p$-values = $5.543 \times 10^{-84}$ and

733   $3.292 \times 10^{-272}$, slopes = $-1.853 \times 10^{-6}$ and $-3.768 \times 10^{-6}$).

734

735   **Tables**

736   **Table 1   Performance metrics of DeepHoF and BLAST**

| Methods | Precision | Accuracy | TPR | FPR | AUC | F1-score |
|---|---|---|---|---|---|---|
| BLAST | 0.699 | 0.892 | 0.888 | 0.107 | 0.833 | 0.896 |
| DeepHoF | 0.968 | 0.964 | 0.865 | 0.008 | 0.987 | 0.963 |

TPR: true-positive rate; FPR: false-positive rate; AUC: area under the curve

737

738   **Table 2   Host prediction results of SARS-CoV-2**

| Prediction | Evidence of infection with SARS-CoV-2 [5] | Reported transmission to humans |
|---|---|---|

26

| | | | |
|---|---|---|---|
|  | *Mustela lutreola / Neovison vison* | - From 19 April to 1 October, 2020, out of around 120 mink farms in Netherlands, 57 have been declared infected;<br>- From 17 June to 1 October, 2020, SARS-CoV-2 has been detected in 41 mink farms in Denmark;<br>- On 16 July, 2020, 80% of the animal samples were tested positive in a Spanish farm;<br>- On 17 August, 2020, confirmed cases were reported in minks at two farms in Utah, the United States;<br>- On 9 October, 2020, 10,000 minks were dead at the United States fur farms and believed infected by SARS-CoV-2. | - Two cases that minks transmitted SARS-CoV-2 to humans in Dutch farms were reported by Nature on 1 June 2020 [11]. |
|  | *Rhinolophus sinicus / Hipposideridae* | - SARS-CoV-2 is 96% identical at the whole-genome level to a bat coronavirus. | N.A. |
|  | *Canis lupus familiaris* | - Confirmed cases in dogs were reported in Hong Kong, New York, Georgia, Texas, South Carolina, *etc*. | N.A. |
| | *Felidae* | - Laboratory confirmed cases of cats; | N.A. |

| | - Four tigers and three lions at the same facility were all confirmed with SARS-CoV-2 in New York in April, 2020; - Confirmed cases in cats in New York, Minnesota, Illinois, California. |
|---|---|

*Note*: N.A. - not available yet.

Hong Kong, Hong Kong Special Administrative Region of the People's Republic of China.

Utah, New York, Georgia, Texas, South Carolina, Minnesota, Illinois, California are states of the United States.

739

740 **Supplementary material**

741 **Supplementary material    Supplemental Figure S1-S6, Supplemental Table S1,**

742 **S3 and S6 and Supplemental Methods**

743 **Supplemental Figure S1    ROC curves and AUC values of DeepHoF and BLAST**

744 **on five host types**

745 DeepHoF performs better than BLAST on AUC of each host type.

746 **Supplemental Figure S2    The untenable linear correlations between the lengths**

747 **and the host likelihood scores for genes of SARS-CoV-2**

748 For the genes of SARS-CoV-2, there is no statistical significance in the linear

749 correlations between the lengths and the host likelihood scores on plant (**A**), germ (**B**),

750 invertebrate (**C**), vertebrate (**D**) and human (**E**).

751 **Supplemental Figure S3    Human host likelihood scores of 5 genes of SARS-**

752 **CoV-2, SARS-CoV and MERS-CoV**

753 Although all the three coronaviruses possess ORF1ab and four structural genes (S, M,

754 N, E), these genes made different contributions on human host likelihood scores in

755 these three viruses (two-sided unpaired Welch Two Sample $t$-test, $p$-value $< 0.05$). S

28

756     gene and M gene contributed more in SARS-CoV-2 and SARS-CoV, while N gene

757     and E gene were more significant in MERS-CoV.

758     **Supplemental Figure S4    Visualization of the host likelihood score profiles of**

759     **SARS-CoV-2 isolates from different GISAID clades and the manually mutated**

760     **SARS-CoV-2 isolates on two-dimensional PCA**

761     SARS-CoV-2 isolates fall into several clear fusiform clusters with different colors

762     according to their clades. Manually mutated with specific marker variants, the 17

763     earliest sequenced isolates move to the corresponding fusiform cluster of the clade

764     that is represented by the specific marker variants.

765     **Supplemental Figure S5    Six phases of an input sequence**

766     For coding the COH matrix of a given sequence, we represented it with the direct

767     conjunction of its six phases, generated from its complementary strand and itself.

768     **Supplemental Figure S6    Structure of BiPathCNN in DeepHoF**

769     BOH matrix and COH matrix are input into two paths independently and transformed

770     by the convolution and pooling layers at the beginning. A concatenation layer and a

771     normalization layer combine the output of the two paths. Five sub-paths process the

772     combined intermediate output individually. Each sub-path contains a full connection

773     layer, a normalization layer and an output layer with sigmoid activation and binary

774     cross-entropy loss function. The five sub-paths output the host likelihood scores on

775     five host types respectively.

776     **Supplemental Table S1    Comparison of performance of DeepHoF and BLAST on**

777     **each host type classification**

778     **Supplemental Table S3    Top 20 hosts predicted by DeepHoF on SARS-CoV-2**

779     **Supplemental Table S6    Subtypes in five host types**

780     **Other supplementary material for this manuscript includes the following:**

781     **Supplemental Table S2    Metadata and host likelihood scores of genes for SARS-**

782     **CoV, MERS-CoV and SARS-COV-2 isolates**

783     **Supplemental Table S4    Contributions of 11 genes in the determination of hosts**

784     **for SARS-CoV-2**

785 **Supplemental Table S5    Metadata, host likelihood score profiles, and high**

786 **frequency SNPs on 53759 SARS-CoV-2 isolates**
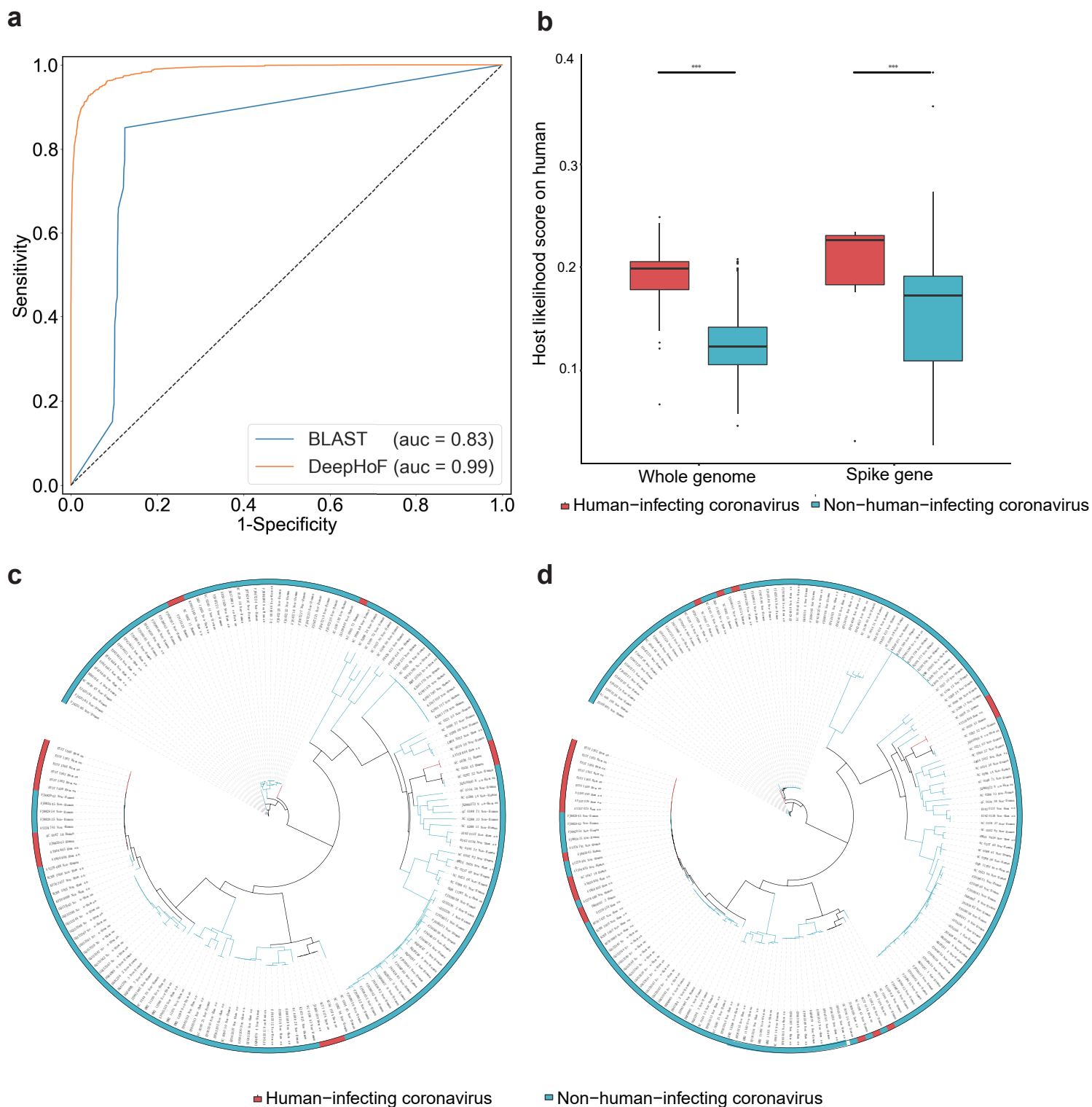
787 **Supplemental Table S7    Host information of the viral genomes in training and**
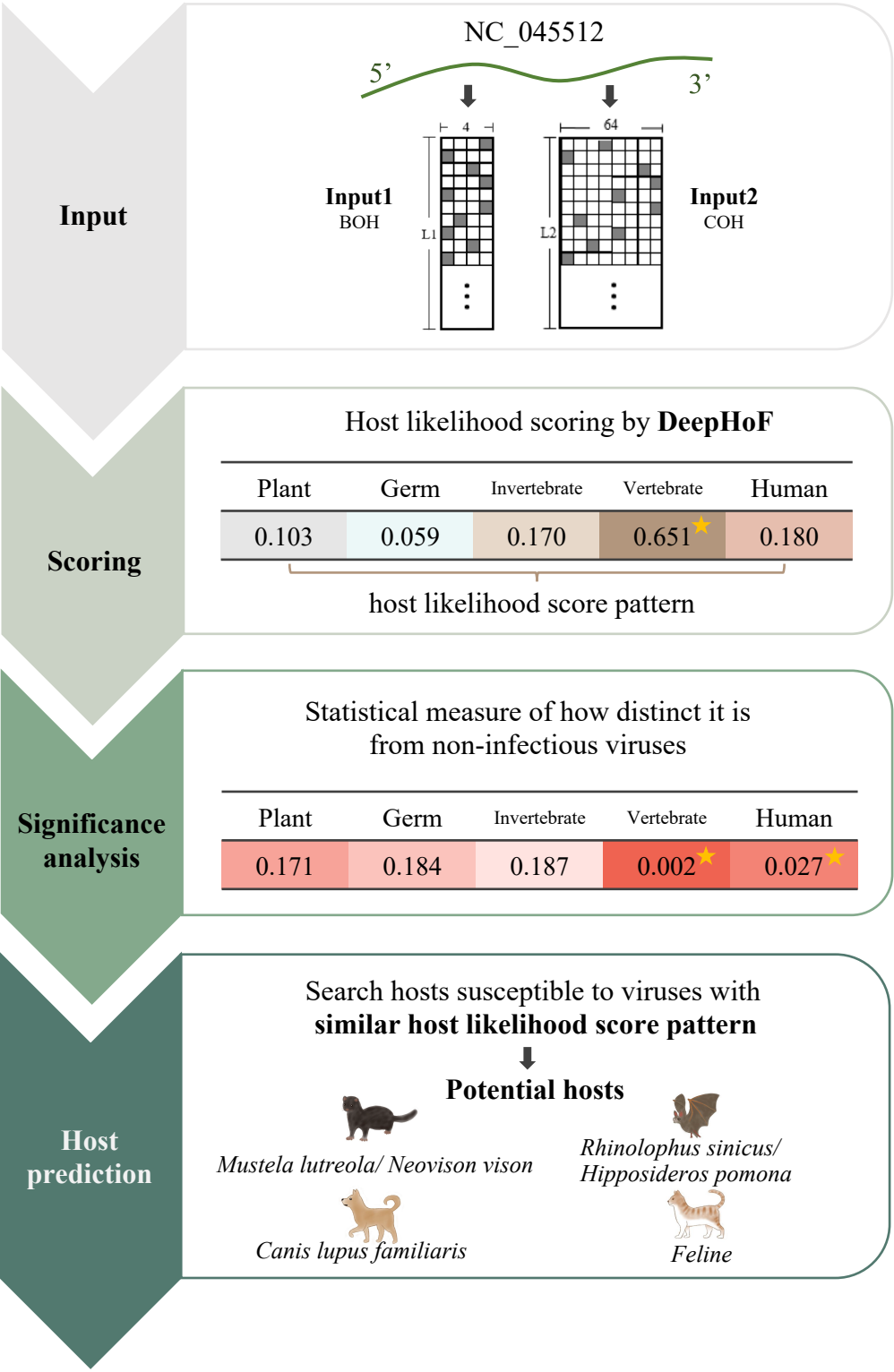
788 **test sets of DeepHoF**

789 **Supplemental Table S8    Acknowledge of sequence data of SARS-CoV-2 in**

790 **GISAID**

791

792 ## Data statement

793 Data utilized in the analysis of SARS-CoV-2, including the host likelihood score

794 profiles and the metadata of 53,759 SARS-CoV-2 isolates, are available in the main text

795 and Supplementary Information. The trimmed sequences of 53,759 isolates and the

796 training and test sets of DeepHoF have been deposited on our lab homepage

797 http://cqb.pku.edu.cn/ZhuLab/DeepHoF/.

798 The open source code utilized in this study has been deposited on GitHub

799 https://github.com/PKUbioinfo-ZhuLab/DeepHoF    and    our    lab    homepage

800 http://cqb.pku.edu.cn/ZhuLab/DeepHoF/

**a**

**b**

**c**

**d**

■ Human−infecting coronavirus    ■ Non−human−infecting coronavirus

**Input**

NC_045512

5'                    3'

**Input1**
BOH

4

L1

**Input2**
COH

64

L2

---

**Scoring**

Host likelihood scoring by **DeepHoF**

| Plant | Germ | Invertebrate | Vertebrate | Human |
|---|---|---|---|---|
| 0.103 | 0.059 | 0.170 | 0.651 ★ | 0.180 |

host likelihood score pattern

---

**Significance analysis**

Statistical measure of how distinct it is from non-infectious viruses

| Plant | Germ | Invertebrate | Vertebrate | Human |
|---|---|---|---|---|
| 0.171 | 0.184 | 0.187 | 0.002 ★ | 0.027 ★ |

---

**Host prediction**

Search hosts susceptible to viruses with **similar host likelihood score pattern**

**Potential hosts**

*Mustela lutreola/ Neovison vison*

*Rhinolophus sinicus/ Hipposideros pomona*

*Canis lupus familiaris*

*Feline*