

1 **Title:** Automated, phylogeny based genotype delimitation of the Hepatitis Viruses HBV and
2 HCV

Deleted: ic

4 **Authors:** Dora Serdari¹, Evangelia-Georgia Kostaki¹, Dimitrios Paraskevis², Alexandros
5 Stamatakis^{1,3}, Paschalia Kapli^{4*}

Formatted: Portuguese (Brazil)

7 **Affiliations**

8 The Exelixis Lab, Scientific Computing Group, Heidelberg Institute for Theoretical Studies,
9 Heidelberg, Germany

Deleted: Heidelberg Institute for Theoretical Studies,
Heidelberg, Germany

10 ¹Department of Hygiene Epidemiology and Medical Statistics, School of Medicine, National
11 and Kapodistrian University of Athens

12 ²Karlsruhe Institute of Technology, Institute for Theoretical Informatics, Karlsruhe, Germany

13 ³Centre for Life's Origins and Evolution, Department of Genetics Evolution and Environment,
14 University College London, University of London, London, United Kingdom

Deleted: Centre for Life's Origins and Evolution,
Department of Genetics Evolution and Environment,
University College London, London, United Kingdom

25 *Corresponding Author:

26 Paschalia Kapli

27 Darwin Building, Gower Street, London, WC1E 6BT, UK

28 **Email address:** p.kapli@ucl.ac.uk

29

36 **Abstract**

37

38 **Background:** The classification of hepatitis viruses still predominantly relies on ad hoc
39 criteria, i.e., phenotypic traits and arbitrary genetic distance thresholds. Given the subjectivity
40 of such practices coupled with the constant sequencing of samples and discovery of new
41 strains, this manual approach to virus classification becomes cumbersome and impossible to
42 generalize.

43

44 **Methods:**

45 Using two well-studied hepatitis virus datasets, HBV and HCV, we assess if computational
46 methods for molecular species delimitation that are typically applied to barcoding biodiversity
47 studies can also be successfully deployed for hepatitis virus classification. For comparison,
48 we also used ABGD, a tool that in contrast to other distance methods attempts to
49 automatically identify the barcoding gap using pairwise genetic distances for a set of aligned
50 input sequences.

51

52 **Results - Discussion:**

53 We find that, the mPTP species delimitation tool identified ~~even~~ without adapting its default
54 parameters ~~taxonomic clusters that either correspond to the currently acknowledged~~
55 genotypes or to known subdivision of genotypes (subtypes or ~~subgenotypes~~). In the cases
56 where the delimited cluster corresponded to subtype or subgenotype, there were previous
57 concerns that their status may be underestimated. The clusters obtained from the ABGD
58 analysis differed depending on the parameters used. However, under certain values the results
59 were very similar to the taxonomy and mPTP which indicates the usefulness of distance based
60 methods in virus taxonomy under ~~appropriate parameter settings~~. The overlap of predicted
61 ~~acknowledged~~ clusters among methods and taxonomically ~~genotypes~~ implies that virus
62 classification can be successfully automated.

63

64

Deleted: -

Deleted: -

Deleted: ,

Formatted: Font: Times, 12 pt

Deleted:

Formatted: Font: Times, 12 pt

Deleted: well informed parameter

Deleted:

Formatted: Font: Times, 12 pt

Deleted: acknowledge

72 Introduction

73

74 The continuous advances in next generation sequencing technologies lead to an increasingly
75 easier and inexpensive production of genome and metabarcoding data. The wealth of
76 available data has triggered the development of new models of molecular evolution,
77 algorithms, and software, that aim to improve molecular sequence analyses in terms of
78 biological realism, computational efficiency, or a trade-off between the two. In response to
79 such technological and technical advancements, several fields of biology have undergone a
80 substantial transformation. Sequence-based species delimitation and identification, in the
81 framework of DNA-(meta)barcoding constitutes a representative example that revived
82 taxonomy and systematics (Tautz *et al.* 2003; Moritz & Cicero, 2004; Savoleinen *et al.*, 2005;
83 Waugh, 2007; Bucklin *et al.*, 2010; Valentini *et al.*, 2009; Li *et al.*, 2014), while **it also**
84 **provided** a new means of analysis in several fields (Galimberti *et al.*, 2013; Mishra *et al.*,
85 2015; Lewray & Knowlton, 2015; Bell *et al.*, 2016; Batovska *et al.*, 2017). Among others, the
86 development of novel species delimitation tools has substantially advanced the study of
87 biodiversity of microorganism that are often hard to isolate and study (Taberlet *et al.*, 2012;
88 Gibson *et al.*, 2014; Thomsen & Willerslev, 2015). The sequencing of environmental samples
89 in conjunction with algorithms for genetic clustering has led to the identification of a plethora
90 of previously unknown organisms and a re-assessment of the microbial biodiversity in several
91 settings.

92 In a similar context, genetic information has been a rich source of information for viral
93 species. Several studies show how phylogenetic information can be deployed for identifying
94 the spatial and temporal origin of a virus, potential factors that trigger its dispersal, and other
95 key epidemiological parameters (Stadler *et al.*, 2012a; Stadler *et al.*, 2014b; Gire *et al.*,
96 2014). In an era of high human mobility, such methods are important, as the increase of
97 emerging and re-emerging epidemics is even more prominent than in the past (Balcan *et al.*,
98 2009; Meloni *et al.*, 2011; Pybus *et al.*, 2015). Nevertheless, phylogenetic information is still
99 not used in the context of virus species classification or identification. As we have witnessed
100 for other microorganisms, using or adapting already available methods for fast and automated
101 delimitation or identification of virus species can greatly contribute to better understand their
102 evolution.

103 To date, the official taxonomy of viruses (ICTV, i.e., International Committee on
104 Taxonomy of Viruses) has mainly been based on established biological classification criteria
105 as used for other life forms, such as plants or animals. An analogous hierarchical

Deleted: provide

107 classification system containing orders, families, subfamilies, genera, and species is being
108 applied (*Simmonds, 2015*). The ICTV is typically based on phenotypic criteria, such as
109 morphology, nucleic acid type (i.e., DNA or RNA), hosts, symptoms, mode of replication,
110 geographical data, or presence of antigenic epitopes, to name a few. Generally, such criteria,
111 despite being informative, can be subjective, require highly specialized knowledge, and are
112 time consuming to apply. In contrast, sequence evolution takes into account the evolutionary
113 history of life forms and, thus, may offer a more objective source of information for
114 taxonomic classification. An important difference in viruses compared to other organisms is
115 that they lack a common set of universal genes such as the 18S rRNA in eukaryotes or the
116 16S rRNA in prokaryotes. Therefore, we cannot infer a comprehensive virus tree of life
117 (*Simmonds et al., 2017*), and, more importantly for species delimitation, we cannot rely upon
118 barcoding markers that are universally suitable for all viruses. We can nonetheless gain
119 valuable insights for their systematics by utilizing phylogenetic information at lower
120 taxonomic ranks (e.g., families, genera, species), using appropriate genes for each dataset. In
121 this context, methods using genetic-distance thresholds (*Bao et al., 2014, Lauber &*
122 *Gorbalenya, 2012, Yu et al., 2013*) have been suggested as a complementary method to the
123 traditional virus classification for accelerating new species identification.

124 In this study, we explore whether a recently developed algorithm for molecular species
125 delimitation on barcoding or marker gene phylogenies can be deployed for ICTV. In contrast
126 to genetic distance-based methods the multi-rate Poisson Tree Processes (mPTP, *Kapli et al.,*
127 *2017*) infers the number of genetic clusters given a phylogenetic input tree. Such trees can
128 easily be inferred using both, Maximum Likelihood (*Stamatakis, 2014*), or Bayesian
129 approaches (*Ronquist et al., 2012*) on single-gene or multi-gene multiple sequence
130 alignments. The fundamental assumption of the model is that variance in the data, as
131 represented by the phylogeny, is greater among species than within a species (*Zhang et al.,*
132 *2013*). The additional assumption of mPTP, that the genetic variation may differ substantially
133 among species allows to accurately delimit species in large (meta-) barcoding datasets
134 comprising multiple species of diverse life histories (*Kapli et al., 2017*). Experiments using
135 empirical data for several animal phyla (*Kapli et al., 2017*) and recently also viruses (*Thézé et*
136 *al., 2018; Modha et al, 2018*) show that the method consistently provides extremely fast and
137 sensible species estimates on ‘classic’ phylogenetic marker and barcoding genes.

138 To assess whether mPTP can be deployed as a quantitative ICTV method we analyze
139 two medically important viruses, Hepatitis B (HBV) and Hepatitis C (HCV) that are leading
140 global causes of human mortality (*Stanaway et al., 2016*). Both viruses cause liver

Deleted: distance based

Deleted: ,

Deleted: a

144 inflammation, but are substantially different from each other. HBV has a partially double-
145 stranded circular DNA genome with a length of about 3.2 kb while HCV is a single-stranded,
146 positive-sense RNA virus, with a genome length of approximately 10kb (*Radziwill et al.,*
147 *1990; Tang et al., 2001; Martell et al., 1992*). Both virus types comprise at least two
148 taxonomic levels (HBV: genotypes, sub-genotypes; HCV: genotypes, subtypes). Besides the
149 significance of the two viruses for human health, we selected them as test cases since due to
150 the substantial amount of taxonomic research that has been conducted and that we can hence
151 use to assess the efficiency of genetic clustering (e.g., *Simmonds et al., 2005; Schaefer, 2007;*
152 *Smith et al., 2014; Messina et al., 2015*).

153
154

155 **Materials and methods**

156

157 *Datasets*

158 We generated multiple sequence alignments (MSAs) corresponding to two virus types: HBV
159 and HCV from two sets of full-length genomic sequences downloaded from publicly available
160 databases (NCBI: <http://www.ncbi.nlm.nih.gov>, Accession Numbers provided in Suppl.
161 Appendix).

162 The HBV dataset comprises 110 sequences corresponding to eight genotypes (i.e., A-H) and
163 31 subgenotypes. The genotypes (A through D, F, and H) have been further divided into
164 subgenotypes indexed by numbers for the corresponding genotype (e.g., A₁, A₂, B₁, B₂, B₃, etc.;
165 *Kramvis et al., 2014*). The inter-genotypic and inter-subgenotypic divergence exceeds 8% and
166 4%-8%, respectively across the genome. No sub-genotypes have been reported for genotypes
167 E, G and H which shows that they are of lower levels of genetic divergence than the rest. The
168 distribution of HBV genotypes differs greatly with respect to the geographical origin.

169 Moreover, they differ in their natural history, response to treatment and disease progression
170 (*Huang, 2013; Biswas, 2013; Moura, 2013; Shi, 2013*). For our study we included the
171 sequences of the eight genotypes (A-H) that form part of the oldest identified HBV groups.

172 The HCV dataset I) comprises 213 sequences corresponding to seven major taxonomic
173 units named after genotypes (1, 2, 3, 4, 5, 6, and 7) and numerous subtypes (*Smith et al.,*
174 *2014*). The HCV classification into genotypes and subtypes was based on genetic-distance
175 thresholds that were verified by the fact that they formed monophyletic clades in an inferred
176 phylogeny (*Smith et al., 2014*). Therefore, the HCV classification serves as an appropriate test
177 case for assessing whether a similar clustering can be identified with a more objective and

Deleted: obtained

Deleted: two previously published

Deleted: (*Kramvis, 2014; Smith et al., 2014*, respectively)

Deleted: .

Deleted: differ

183 automated method, such as mPTP, that does not require any user input apart from a
184 phylogeny.

185

186 *Genetic Cluster delimitation*

187 To delimit the putative species, additionally to mPTP, we used the distance-based “*Automatic*
188 *Barcode Gap Discovery*” tool (ABGD, *Puillandre et al., 2012*). ABGD is a popular distance-
189 based barcoding method that, compared to other distance-based methods attempts to
190 automatically identify the threshold value for the transition from intra-specific variation to
191 inter-specific divergence (*Puillandre et al., 2012*).

192 For the mPTP delimitation, a fully binary (bifurcating) rooted phylogeny is required.
193 Therefore, using the aligned sequences we inferred the phylogenetic relationships under the
194 GTR+*I* model of nucleotide substitution using RAxML-NG (*Kozlov et al., 2018*). We rooted
195 the phylogenetic trees according to the originally published phylogenies (i.e., using the branch
196 leading to genotypes F/H for HBV and genotype 7 for HCV). Using heuristic search
197 algorithms for finding the ‘best’ delimitation given the rooted phylogeny and without any
198 further prior assumptions. We performed the mPTP delimitation under Maximum Likelihood
199 (ML) and calculated the support of the delimited clusters using Markov-chain Monte Carlo
200 (MCMC) sampling (*Kapli et al., 2017*). We conducted the MCMC sampling twice for 10⁶
201 generations, to identify potential lack of convergence with a sampling frequency of 0.1.

202 For ABGD, the user has to define two important parameters, i) the prior maximum
203 divergence of intraspecific diversity (*P*), which implies that the barcode gap is expected to
204 exceed this value and should not be confused with the genetic thresholds assumed to define
205 the inter-specific relationships, ii) a proxy for the minimum gap width (*X*), which indicates
206 that the barcoding gap is expected to be *X* times larger than any intraspecific gap (*Puillandre*
207 *et al., 2012*). For both, HBV, and HCV, we used 10 prior maximum thresholds in the range of
208 $p = 0.001$ and $P = 0.05$. The proxy for the minimum gap width (*X*) was set to the default value
209 ($X = 1.5$) for HCV, while for HBV the default value did not yield any delimitation and we
210 therefore set it to a lower value ($X = 0.5$).

211

212

213 **Results & Discussion**

214

215 The biodiversity of viruses is tremendous and it is broadly accepted that our understanding of
216 their ecology and evolution is constrained to a small fraction of species (*Paez-Espino et al.,*

217 [2016](#)). In just a kilo of marine sediment there can be a million of different viral genotypes
 218 (*Breitbart & Rohwer, 2005*), while on a global scale the number of viruses is 10 million-fold
 219 higher than the number of stars in the universe (*Suttle, 2013*). The classification of such a
 220 diverse set of organisms constitutes a challenging task and is impossible to accomplish within
 221 reasonable time using phenotypic characters. Quantitative computational methods could
 222 provide a viable alternative, particularly for large scale clustering and fast identification of
 223 viral strains (*Simmonds et al., 2017; Modha et al., 2018*). Using empirical data of the HBV
 224 and HCV viruses we show that by applying phylogeny-aware and distance-based tools to
 225 classify the strains of the two virus types, the corresponding genetic clustering closely
 226 recovers their currently accepted taxonomy.

227

228 **HCV Clustering**

229 The current taxonomy of HCV comprises seven genotypes, while mPTP yielded 16 genetic
 230 clusters (Fig. 1, Suppl. Fig. 1 [and 2](#), Suppl. Appendix). From the 16 clusters, five were
 231 congruent with the current taxonomy, i.e., genotypes 1, 2, 4, 5 and 7. On the contrary,
 232 genotype 3 and genotype 6 were further split into three and eight sub-clusters correspondingly
 233 (Fig. 1), which corroborates former views that divergent variants of these genotypes may
 234 qualify as separate major genotypes (*Simmonds et al., 2005, Smith et al, 2014*). In particular,
 235 the additional clusters identified by mPTP correspond to previously identified groups of
 236 subtypes (Suppl. Fig. 1). For genotype 6, these clusters consisted of the following subtype
 237 groups: 6a; 6b and 6xd; 6c, 6d, 6e, 6f, 6g, 6o, 6p, 6q, 6r, 6s, 6t, 6u, 6w, 6xc and 6xf; 6h, 6i, 6j,
 238 6k, 6l, 6m, 6n, 6xb, 6xe; 6xa; 6v (Suppl. Fig. 1, Suppl. Appendix). Similarly, for genotype 3,
 239 the delimited clusters were i) 3g, 3b, 3i, 3a, 3e, 3d, ii) 3k, and iii) 3h and 3. All clusters were
 240 substantially supported by the MCMC sampling, except the split of 3k subtype from its sister
 241 group (Fig. 1), which may be due to [the limited amount of corresponding sequences](#),

242 The number of clusters inferred with ABGD ranged from 1 to 208 depending on the
 243 value of the maximum intraspecific divergence threshold (Fig. [2](#)). The most reasonable result
 244 (i.e., the one closest to the current standard taxonomy) comprised 19 clusters and was
 245 obtained for a minimum of intraspecific genetic diversity of 5.99% (i.e., $p=0.0599$). Under
 246 this threshold, the delimitation is largely identical to the delimitation obtained with mPTP
 247 (Fig. 1), with three differences: i) that genotype 3 was split into four clusters, instead of three,
 248 [ii\) genotype six](#) was divided into nine clusters instead of eight, and, iii) genotype 7 is divided
 249 into two clusters. When the prior intraspecific divergence was increased to a higher minimum
 250 of 10%, all sequences were grouped in a single cluster. When the threshold was set to a lower

Deleted: '

Deleted: lack of adequate

Deleted: for the subtype

Deleted: 3

255 value (3.6%) the number of clusters increased to 135 (Fig. 2). Nevertheless, the delimitation
256 with the 5.99% threshold is largely congruent to current taxonomy and the clusters obtained
257 with mPTP, thus indicating the usefulness of distance-based methods in virus taxonomy under
258 well informed parameters.

Deleted: 3

Deleted: distance based

259 The so far classification of HCV into genotypes and subtypes has been defined mostly
260 by visual identification of clades in phylogenetic inference of HCV sequences (Simmonds *et al.*,
261 2005; Smith *et al.*, 2014). Specifically, the genotypes correspond to the seven major
262 highly-supported phylogenetic HCV clusters while subtypes were defined as the secondary
263 hierarchical clusters found within each genotype (Smith *et al.*, 2014). This classification
264 scheme has been widely adopted (Combet *et al.*, 2007; Yusim *et al.*, 2015) and has been
265 shown to be robust (in terms of stability of the HCV phylogeny) and relevant for clinical
266 practice, since response rates to immunomodulatory treatment for the chronic hepatitis C
267 differs across genotypes. Nevertheless, new, unassigned lineages are often discovered from
268 understudied areas (Sulbaran *et al.*, 2010; Nakano *et al.*, 2011; Lu *et al.*, 2013; Tong *et al.*,
269 2015) and it is challenging to assign them a taxonomic status, given that the genetic distance
270 cut-off among intra and inter-specific relationships is arbitrary and variable for different parts
271 of the HCV phylogeny (Simmonds *et al.*, 2005). The greatly overlapping mPTP and ABGD
272 clusters with the HCV genotypes shows that the classification, and, consequently, the
273 identification, of the genotypes can be easily automated utilizing objective, transparent, and
274 unifying approaches. Embracing such alternatives can be crucial for viruses like HCV, taking
275 into account that the correct identification of the HCV genotypes, could be of clinical
276 importance in providing the appropriate medical treatment (Strader *et al.*, 2004; Ge *et al.*,
277 2009).

Deleted: ying

Deleted: is

Deleted: for

278

279 HBV clustering

280 In the case of HBV, the mPTP clustering is almost identical to the current classification
281 (Norder *et al.*, 2004; Kramvis *et al.*, 2007) of the virus that comprises eight genotypes, except
282 for subgenotype C4 which formed a new cluster (Fig. 3, Suppl. Fig. 3 and 4, Suppl.
283 Appendix). This is in line with the greater genetic divergence of C4 compared to the other
284 subgenotypes due to its ancient origin in native populations in Oceania (Paraskevis *et al.*,
285 2013). However, the split of C4 from its sister cluster (genotype C) is not supported by the
286 MCMC sampling, potentially reflecting the lack of adequate sampling. On the other hand, the
287 number of clusters identified by ABGD varied from 1 to 85 under different thresholds of
288 minimum intraspecific divergence, while the delimitation for a threshold of 1.29% exactly

Deleted: 2

Deleted: 2

296 matched the eight genotypes of the HBV classification (Fig. 2 and 3). Both ABGD and mPTP
297 identified seven of the genotypes (A-F) as distinct genetic clusters. The only difference was
298 that mPTP split genotype C into two distinct clusters (Fig. 3), i.e., subtype C4 was recovered
299 as a distinct cluster from the remaining seven subtypes of genotype.

301 **Conclusions**

302 The application of mPTP to the HCV and HBV data sets shows that automated viral species
303 delimitation using phylogeny-aware methods yields clusters that largely agree with the current
304 standard taxonomy. The additional clusters identified for HCV by mPTP is not surprising as
305 they have been previously considered divergent sub-clusters within the genotypes 3 and 6.
306 Analogously, for HBV, mPTP yielded almost identical results to the current nomenclature
307 system with the exception of a single sub-genotype, C4, that was previously mentioned to be
308 more genetically divergent within genotype C (Paraskevis et al, 2013). In both cases, these
309 new clusters indicate the potential need for taxonomic revision. However, given the wide use
310 of the current nomenclature in the medical field, and the lack of other sources of information
311 such as recombination, particularly for HBV, and, response to treatment, we wouldn't suggest
312 taxonomic changes at present. Regarding distance methods, the example of HCV and HBV,
313 shows that meaningful parameter values for distance- based methods may differ substantially
314 among datasets, and, therefore, establishing global thresholds is impossible. On the contrary,
315 mPTP can be seamlessly applied to taxa of substantially different life histories (e.g., variable
316 population sizes, evolution rates), as it does not require any input parameters except a
317 phylogeny. Overall, the ease-of-use of mPTP in conjunction with its computational efficiency
318 on phylogenies with hundreds of samples render it a useful tool for viral biodiversity
319 estimates, initial classification of understudied taxa, and accelerating the viral species
320 identification process.

321

322

323 ▼

324

325

326 **Acknowledgments**

327 The authors gratefully acknowledge the support of the Klaus Tschira Foundation.

328

329 **References**

Deleted: 2

Deleted: The application of mPTP to the HCV and HBV data sets shows that automated viral strain clustering using phylogeny-aware methods yields clusters that largely agree with the current standard taxonomy

Deleted: The major advantage of mPTP over distance-based approaches is that it can be seamlessly applied to taxa of substantially different life histories (e.g., variable population sizes, evolution rates), as it does not require any input parameters except a phylogeny. On the contrary, the example of HCV and HBV, shows that meaningful parameter values for distance based methods may differ substantially among datasets, and, therefore, establishing global thresholds is impossible. ...

Deleted: Acknowledgements

Formatted: Portuguese (Brazil)

345
346 Balcan, D., Colizza, V., Gonçalves, B., Hu, H., Ramasco, J. J., & Vespignani, A. (2009).
347 Multiscale mobility networks and the spatial spreading of infectious diseases. *Proceedings of*
348 *the National Academy of Sciences*, 106(51), 21484-21489.
349
350 Bao, Y., Chetvernin, V., & Tatusova, T. (2014). Improvements to pairwise sequence
351 comparison (PASC): a genome-based web tool for virus classification. *Archives of virology*,
352 159(12), 3293-3304.
353
354 Batovska, J., Lynch, S. E., Cogan, N. O. I., Brown, K., Darbro, J. M., Kho, E. A., & Blacket,
355 M. J. (2018). Effective mosquito and arbovirus surveillance using metabarcoding. *Molecular*
356 *ecology resources*, 18(1), 32-40.
357
358 Bell, K. L., Burgess, K. S., Okamoto, K. C., Aranda, R., & Brosi, B. J. (2016). Review and
359 future prospects for DNA barcoding methods in forensic palynology. *Forensic Science*
360 *International: Genetics*, 21, 110-116.
361
362 Biswas, A., Panigrahi, R., Pal, M., Chakraborty, S., Bhattacharya, P., Chakrabarti, S., &
363 Chakravarty, R. (2013). Shift in the hepatitis B virus genotype distribution in the last decade
364 among the HBV carriers from eastern India: possible effects on the disease status and HBV
365 epidemiology. *Journal of medical virology*, 85(8), 1340-1347.
366
367 Bolotov, I. N., Vikhrev, I. V., Kondakov, A. V., Konopleva, E. S., Gofarov, M. Y., Aksenova,
368 O. V., & Tumpeesuwan, S. (2017). New taxa of freshwater mussels (Unionidae) from a
369 species-rich but overlooked evolutionary hotspot in Southeast Asia. *Scientific Reports*, 7(1),
370 11573.
371
372 Breitbart, M., & Rohwer, F. (2005). Here a virus, there a virus, everywhere the same virus?
373 *Trends in microbiology*, 13(6), 278-284.
374
375 Bucklin, A., Hopcroft, R. R., Kosobokova, K. N., Nigro, L. M., Ortman, B. D., Jennings, R.
376 M., & Sweetman, C. J. (2010). DNA barcoding of Arctic Ocean holozooplankton for species
377 identification and recognition. *Deep Sea Research Part II: Topical Studies in Oceanography*,
378 57(1-2), 40-48.
379
380 Combet, C., Garnier, N., Charavay, C., Grando, D., Crisan, D., Lopez, J., A., Dehne-Garcia,
381 C., Geourjon, E., Bettler, C., H., P., Le Mercier, R., Bartenschlager, H., Diepolder, D.,
382 Moradpour, J.-M., Pawlotsky, C., M., Rice, C., Trépo, F., Penin, G., Deléage. euHCVdb: the
383 European hepatitis C virus database. *Nucleic Acids Res* 2007; 35: D363- D366.
384
385 Galimberti, A., De Mattia, F., Losa, A., Bruni, I., Federici, S., Casiraghi, M., Martelos, S., &
386 Labra, M. (2013). DNA barcoding as a new tool for food traceability. *Food Research*
387 *International*, 50(1), 55-63.
388

Formatted: Portuguese (Brazil)

389 Ge, D., Fellay, J., Thompson, A. J., Simon, J. S., Shianna, K. V., Urban, T. J., Heinzen, E.L.,
 390 Qiu, P., Bertelsen, A.H., Muir, A.J., & Sulikowski, M. (2009). Genetic variation in IL28B
 391 predicts hepatitis C treatment-induced viral clearance. *Nature*, 461(7262), 399.
 392
 393 Gibson, J., Shokralla, S., Porter, T. M., King, I., van Konynenburg, S., Janzen, D. H.,
 394 Hallwachs, W., & Hajibabaei, M. (2014). Simultaneous assessment of the macrobiome and
 395 microbiome in a bulk sample of tropical arthropods through DNA metasytematics.
 396 *Proceedings of the National Academy of Sciences*, 111(22), 8007-8012.
 397
 398 Gire, S. K., Goba, A., Andersen, K. G., Sealfon, R. S., Park, D. J., Kanneh, L., Jalloh, S.,
 399 Momoh, M., Fullah, M., Dudas, G., Wohl, S., Moses, L. M., Yozwiak, N. L., Winnicki, S.,
 400 Matranga, C.B., Malboeuf, C. M., Qu, J., Gladden, A. D., Schaffner, S. F., Yang, X., Jiang,
 401 P., Nekoui, M., Colubri, A., Coomber, M. R., Fonnies, M., Moigboi, A., Gbakie, M., Kamara,
 402 F. K., Tucker, V., Konuwa, E., Saffa, S., Sellu, J., Jalloh, A. A., Kovoma, A., Koninga, J.,
 403 Mustapha, I., Kargbo, K., Foday, M., Yillah, M., Kanneh, F., Robert, W., Massally, J. L. B.,
 404 Chapman, S.B., Bochichio, J., Murphy, C., Nusbaum, C., Young, S., Birren, B. W., Grant, D.
 405 S., Scheiffelin, J. S., Lander, E. S., Hapfi, C., Gevao, S. M., Gnirke, A., Rambaut, A., Garry,
 406 R. F., Khan, S. H., & Sabeti, P. C. (2014). Genomic surveillance elucidates Ebola virus origin
 407 and transmission during the 2014 outbreak. *Science*, 345(6202), 1369-1372.
 408
 409 Huang, C. C., Kuo, T. M., Yeh, C. T., Hu, C. P., Chen, Y. L., Tsai, Y. L., Chen, M. L.,
 410 Chang, C., & Chang, C. (2013). One single nucleotide difference alters the differential
 411 expression of spliced RNAs between HBV genotypes A and D. *Virus research*, 174(1-2), 18-
 412 26.
 413
 414 Kapli, P., Lutteropp, S., Zhang, J., Kobert, K., Pavlidis, P., Stamatakis, A., & Flouri, T.
 415 (2017). Multi-rate Poisson tree processes for single-locus species delimitation under
 416 maximum likelihood and Markov chain Monte Carlo. *Bioinformatics*, 33(11), 1630-1638.
 417
 418 Kozlov, A., Darriba, D., Flouri, T., Morel, B., & Stamatakis, A. (2018). RAXML-NG: A fast,
 419 scalable, and user-friendly tool for maximum likelihood phylogenetic inference. *bioRxiv*,
 420 447110.
 421
 422 Kramvis, A., & Kew, M. C. (2007). Epidemiology of hepatitis B virus in Africa, its genotypes
 423 and clinical associations of genotypes. *Hepatology Research*, 37, S9-S19.
 424
 425 Kramvis, A. (2014). Genotypes and genetic variability of hepatitis B virus. *Intervirology*,
 426 57(3-4), 141-150.
 427
 428 Lauber, C., & Gorbalenya, A. E. (2012). Genetics-based classification of filoviruses calls for
 429 expanded sampling of genomic sequences. *Viruses*, 4(9), 1425-1437.
 430

431 Leray, M., & Knowlton, N. (2015). DNA barcoding and metabarcoding of standardized
432 samples reveal patterns of marine benthic diversity. *Proceedings of the National Academy of*
433 *Sciences*, 112(7), 2076-2081.

434

435 Li, X., Yang, Y., Henry, R. J., Rossetto, M., Wang, Y., & Chen, S. (2015). Plant DNA
436 barcoding: from gene to genome. *Biological Reviews*, 90(1), 157-166.

437

438 Lu, L., Li, C., Yuan, J., Lu, T., Okamoto, H., & Murphy, D. G. (2013). Full-length genome
439 sequences of five hepatitis C virus isolates representing subtypes 3g, 3h, 3i and 3k, and a
440 unique genotype 3 variant. *The Journal of general virology*, 94(Pt 3), 543.

441

442 Martell, M., Esteban, J. I., Quer, J., Genesca, J., Weiner, A., Esteban, R., Guardia, J., &
443 Gomez, J. (1992). Hepatitis C virus (HCV) circulates as a population of different but closely
444 related genomes: quasispecies nature of HCV genome distribution. *Journal of virology*, 66(5),
445 3225-3229.

446

447 Meloni, S., Perra, N., Arenas, A., Gómez, S., Moreno, Y., & Vespignani, A. (2011). Modeling
448 human mobility responses to the large-scale spreading of infectious diseases. *Scientific*
449 *reports*, 1, 62.

450

451 Messina, J. P., Humphreys, I., Flaxman, A., Brown, A., Cooke, G. S., Pybus, O. G., & Barnes,
452 E. (2015). Global distribution and prevalence of hepatitis C virus genotypes. *Hepatology*,
453 61(1), 77-87.

454

455 Mishra, P., Kumar, A., Nagireddy, A., Mani, D. N., Shukla, A. K., Tiwari, R., & Sundaresan,
456 V. (2016). DNA barcoding: an efficient tool to overcome authentication challenges in the
457 herbal market. *Plant biotechnology journal*, 14(1), 8-21.

458

459 Modha, S., Thanki, A. S., Cotmore, S. F., Davison, A. J., & Hughes, J. (2018). ViCTree: an
460 automated framework for taxonomic classification from protein sequences. *Bioinformatics*,
461 34(13), 2195-2200.

462

463 Moritz, C., & Cicero, C. (2004). DNA barcoding: promise and pitfalls. *PLoS biology*, 2(10),
464 e354.

465

466 Moura, I. F., Lopes, E. P., Alvarado-Mora, M. V., Pinho, J. R., & Carrilho, F. J. (2013).
467 Phylogenetic analysis and subgenotypic distribution of the hepatitis B virus in Recife, Brazil.
468 *Infection, Genetics and Evolution*, 14, 195-199.

469

470 Nakano, T., Lau, G. M., Lau, G. M., Sugiyama, M., & Mizokami, M. (2012). An updated
471 analysis of hepatitis C virus genotypes and subtypes based on the complete coding region.
472 *Liver international*, 32(2), 339-345.

473

Formatted: Portuguese (Brazil)

474 Norder, H., Couroucé, A. M., Coursaget, P., Echevarria, J. M., Shou-Dong, L., Mushahwar, I.
 475 K., & Magnius, L. O. (2004). Genetic diversity of hepatitis B virus strains derived worldwide:
 476 genotypes, subgenotypes, and HBsAg subtypes. *Intervirology*, 47(6), 289.
 477
 478 Paez-Espino, D., Eloie-Fadrosch, E. A., Pavlopoulos, G. A., Thomas, A. D., Huntemann, M.,
 479 Mikhailova, N., Rubin, E., Ivanova, N. N. & Kyrpides, N. C. (2016). Uncovering Earth's
 480 virome. *Nature*, 536(7617), 425.
 481
 482 Paraskevis, D., Magiorkinis, G., Magiorkinis, E., Ho, S. Y., Belshaw, R., Allain, J. P., &
 483 Hatzakis, A. (2013). Dating the origin and dispersal of hepatitis B virus infection in humans
 484 and primates. *Hepatology*, 57(3), 908-916.
 485
 486 Puillandre, N., Lambert, A., Brouillet, S., & Achaz, G. (2012). ABGD, Automatic Barcode
 487 Gap Discovery for primary species delimitation. *Molecular ecology*, 21(8), 1864-1877.
 488
 489 Pybus, O. G., Tatem, A. J., & Lemey, P. (2015). Virus evolution and transmission in an ever
 490 more connected world. *Proceedings of the Royal Society B: Biological Sciences*, 282(1821),
 491 20142878.
 492
 493 Radziwill, G., Tucker, W., & Schaller, H. (1990). Mutational analysis of the hepatitis B virus
 494 P gene product: domain structure and RNase H activity. *Journal of virology*, 64(2), 613-620.
 495
 496 Ronquist, F., Teslenko, M., Van Der Mark, P., Ayres, D. L., Darling, A., Höhna, S., Larget,
 497 B., Liu, L., Suchard, M. A., & Huelsenbeck, J. P. (2012). MrBayes 3.2: efficient Bayesian
 498 phylogenetic inference and model choice across a large model space. *Systematic biology*,
 499 61(3), 539-542.
 500
 501 Rulik, B., Eberle, J., von der Mark, L., Thormann, J., Jung, M., Köhler, F., Apfel, W., Weigel,
 502 A., Kopetz, A., Köhler, J., Fritylar, F., Hartmann, M., Hadulla, K., Schmidt, J., Hörren, T.,
 503 Krebs, D., Theves, F., Eulity, U., Skale, A., Rohwedder, D., Kleeberg, A., Astrin, I. I., Geiger,
 504 M. F., Wägele, W., Grobe, P., & Ahrens, D. (2017). Using taxonomic consistency with semi-
 505 automated data pre-processing for high quality DNA barcodes. *Methods in Ecology and*
 506 *Evolution*, 8(12), 1878-1887.
 507
 508 Savolainen, V., Cowan, R. S., Vogler, A. P., Roderick, G. K., & Lane, R. (2005). Towards
 509 writing the encyclopaedia of life: an introduction to DNA barcoding. *Philosophical*
 510 *Transactions of the Royal Society B: Biological Sciences*, 360(1462), 1805-1811.
 511
 512 Schaefer, S. (2007). Hepatitis B virus genotypes in Europe. *Hepatology Research*, 37, S20-
 513 S26.
 514
 515 Shi, W., Zhang, Z., Ling, C., Zheng, W., Zhu, C., Carr, M. J., & Higgins, D. G. (2013).
 516 Hepatitis B virus subgenotyping: history, effects of recombination, misclassifications, and
 517 corrections. *Infection, genetics and Evolution*, 16, 355-361.

518
519 Simmonds, P. (2015). Methods for virus classification and the challenge of incorporating
520 metagenomic sequence data. *Journal of General Virology*, 96(6), 1193-1206.
521
522 Simmonds, P., Adams, M. J., Benkő, M., Breitbart, M., Brister, J. R., Carstens, E. B.,
523 Davidson, A. J., Delwart, E., Gorbalenya A. E., Harrach, B., Hull, R., King, A. M. Q., Koonin,
524 E. V., Krupovic, M., Kuhn, J. H., Lefkowitz E. J., Nibert, M. L., Orton, R., Roossinck, M. J.,
525 Sabanadyovic, S., Sullivan, M. B., Suttle, C. A., Tesh, R. B., van der Vlugt, R. A., Varsani,
526 A., & Zerbini, F. M. (2017). Consensus statement: virus taxonomy in the age of
527 metagenomics. *Nature Reviews Microbiology*, 15(3), 161.
528
529 Simmonds, P., Bukh, J., Combet, C., Deléage, G., Enomoto, N., Feinstone, S., Halfon, P.,
530 Inchauspe, G., Kuiken, C., Maertens, G., Miyokami, M., Murphy, D. G., Okamoto, H.,
531 Pawlotsky, J. M., Penin, F., Sablon, E., Shin-I, T., Stuyver, L. J., Thiel, H. J., Viayov, S.,
532 Weiner, A. J., & Widell, A. (2005). Consensus proposals for a unified system of nomenclature
533 of hepatitis C virus genotypes. *Hepatology*, 42(4), 962-973.
534
535 Smith, D. B., Bukh, J., Kuiken, C., Muerhoff, A. S., Rice, C. M., Stapleton, J. T., &
536 Simmonds, P. (2014). Expanded classification of hepatitis C virus into 7 genotypes and 67
537 subtypes: updated criteria and genotype assignment web resource. *Hepatology*, 59(1), 318-
538 327.
539
540 Strader, D. B., Wright, T., Thomas, D. L., & Seeff, L. B. (2004). Diagnosis, management, and
541 treatment of hepatitis C. *Hepatology*, 39(4), 1147-1171.
542
543 Stadler, T., Kouyos, R., von Wyl, V., Yerly, S., Böni, J., Bürgisser, P., Klimkait, T., Joos, B.,
544 Rieder, P., Xie, D., Günthard, H. F., Drummond, A. J., & Bonhoeffer, S., (2011). Estimating
545 the basic reproductive number from viral sequence data. *Molecular biology and evolution*,
546 29(1), 347-357.
547
548 Stadler, T., Kühnert, D., Rasmussen, D. A., & du Plessis, L. (2014). Insights into the early
549 epidemic spread of Ebola in Sierra Leone provided by viral sequence data. *PLoS currents*, 6.
550
551 Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis
552 of large phylogenies. *Bioinformatics*, 30(9), 1312-1313.
553
554 Stanaway, J. D., Flaxman, A. D., Naghavi, M., Fitzmaurice, C., Vos, T., Abubakar, I., Abu-
555 Raddad, L., Assadi, R., Bhala, N., Cowie, B., Forouzanfour, M. H., Groeger, J., Hanafiah, M.,
556 Jacobsen, K. H., James, S. L., MacLachlan, J., Malekyadeh, R., Martin, N. K., Mokhad, A.
557 A., Mokdad, A. H., Murray, C. J. L., Plass, D., Rana, S., Rein, D. B., Richardus, J. H.,
558 Sanabria, J., Sazyan, M., Shahrzad, S., So, S., Vlassov, V. V., Weiderpass, E., Wiersma, S. T.,
559 Younis, M., Yu, C., El Sayed Zaki, M., & Cooke, G. S. (2016). The global burden of viral
560 hepatitis from 1990 to 2013: findings from the Global Burden of Disease Study 2013. *The*
561 *Lancet*, 388(10049), 1081-1088.

Formatted: Portuguese (Brazil)

562 Sulbaran, M. Z., Di Lello, F. A., Sulbaran, Y., Cosson, C., Loureiro, C. L., Rangel, H.
563 R., Cantaloube, J. F., Campos, R. H., Moratorio, G., Cristina, H., & Pujol, F. H. (2010). Genetic
564 history of hepatitis C virus in Venezuela: high diversity and long time of evolution of HCV
565 genotype 2. *PloS one*, 5(12), e14315.
566
567 Suttle, C. A. (2013). Viruses: unlocking the greatest biodiversity on Earth. *Genome*, 56(10),
568 542-544.
569
570 Taberlet, P., Coissac, E., Pompanon, F., Brochmann, C., & Willerslev, E. (2012). Towards
571 next-generation biodiversity assessment using DNA metabarcoding. *Molecular ecology*,
572 21(8), 2045-2050.
573
574 Tang, H., & McLachlan, A. (2001). Transcriptional regulation of hepatitis B virus by nuclear
575 hormone receptors is a critical determinant of viral tropism. *Proceedings of the National*
576 *Academy of Sciences*, 98(4), 1841-1846.
577
578 Tautz, D., Arctander, P., Minelli, A., Thomas, R. H., & Vogler, A. P. (2003). A plea for DNA
579 taxonomy. *Trends in ecology & evolution*, 18(2), 70-74.
580
581
582 ~~Thézé, J., Lopez-Vaamonde, C., Cory, J. S., & Herniou, E. (2018). Biodiversity, evolution and~~
583 ~~ecological specialization of baculoviruses: a treasure trove for future applied research.~~
584 ~~*Viruses*, 10(7), 366.~~
585
586
587 Thomsen, P. F., & Willerslev, E. (2015). Environmental DNA—An emerging tool in
588 conservation for monitoring past and present biodiversity. *Biological conservation*, 183, 4-18.
589
590 Tong, Y. Q., Liu, B., Liu, H., Zheng, H. Y., Gu, J., Song, E. J., ~~Song, C.~~, & Li, Y. (2015).
591 Accurate genotyping of hepatitis C virus through nucleotide sequencing and identification of
592 new HCV subtypes in China population. *Clinical microbiology and infection*, 21(9), 874-e9.
593
594 Valentini, A., Pompanon, F., & Taberlet, P. (2009). DNA barcoding for ecologists. *Trends in*
595 *ecology & evolution*, 24(2), 110-117.
596
597 Waugh, J. (2007). DNA barcoding in animal species: progress, potential and pitfalls.
598 *BioEssays*, 29(2), 188-197.
599
600 Yu, C., Hernandez, T., Zheng, H., Yau, S. C., Huang, H. H., He, R. L., Yang, J., & Yau, S. S.
601 T. (2013). Real time classification of viruses in 12 dimensions. *PloS one*, 8(5), e64328.
602
603 Yusim, K., Korber, B. T., Brander, C., Barouch, D., de Boer, R., Haynes, B. F., & Watkins,
604 D. (2016). *HIV Molecular Immunology 2015* (No. LA-UR-16-22283). Los Alamos National
605 Lab.(LANL), Los Alamos, NM (United States).

Deleted: Thézé, J., Lopez-Vaamonde, C., Cory, J. S., & Herniou, E. A. (2018). *Viruses*, 10(7), 366

Deleted: ...

Formatted: Portuguese (Brazil)

609
610 Zhang, J., Kapli, P., Pavlidis, P., & Stamatakis, A. (2013). A general species delimitation
611 method with applications to phylogenetic placements. *Bioinformatics*, 29(22), 2869-2876.
612