

1 Diving Deep into Fish Bornaviruses: Uncovering Hidden Diversity and 2 Transcriptional Strategies through Comprehensive Data Mining

3 Mirette Eshak¹, Dennis Rubbenstroth¹, Martin Beer¹, Florian Pfaff^{1*}

4 ¹: Friedrich-Loeffler-Institut, Institute of Diagnostic Virology, Südufer 10, 17493 Greifswald - Insel Riems, Germany

5
6 *Corresponding author:
7 Dr. Florian Pfaff
8 Friedrich-Loeffler-Institut
9 Institute of Diagnostic Virology
10 Südufer 10 | 17493 Greifswald - Insel Riems
11 +49 38351 7 1508
12 florian.pfaff@fli.de

13

14 ABSTRACT

15 Recently, we discovered two novel orthobornaviruses in colubrid and viperid snakes using an *in*
16 *silico* data mining approach. Here, we present the results of a screening of more than 100,000
17 nucleic acid sequence datasets of fish samples from the Sequence Read Archive (SRA) for potential
18 bornaviral sequences. We discovered the potentially complete genomes of seven bornaviruses in
19 datasets from osteichthyans and chondrichthyans. Four of these are likely to represent novel
20 species within the genus *Cultervirus*, and we propose that one genome represents a novel genus
21 within the family of *Bornaviridae*. Specifically, we identified sequences of Wǔhàn sharpbelly
22 bornavirus (WhSBV) in sequence data from the widely used grass carp liver and kidney cell lines
23 L8824 and CIK, respectively. A complete genome of Murray-Darling carp bornavirus (MDCBV) was
24 identified in sequence data from a goldfish (*Carassius auratus*). The newly discovered little skate
25 bornavirus (LSBV), identified in the little skate (*Leucoraja erinacea*) dataset, contained a novel and
26 unusual genomic architecture (N-Vp1-Vp2-X-P-G-M-L), as compared to other bornaviruses. Its
27 genome is thought to encode two additional open reading frames (tentatively named Vp1 and Vp2),
28 which appear to represent ancient duplications of the gene encoding for the viral glycoprotein (G).
29 The datasets also provided insights into the possible transcriptional gradients of these
30 bornaviruses and revealed previously unknown splicing mechanisms.

31 INTRODUCTION

32 The family *Bornaviridae* belongs to the order *Mononegavirales* and includes viruses that are
33 considered zoonotic and can cause severe disease in humans, such as Borna disease virus 1
34 (BoDV-1) [1] and the variegated squirrel bornavirus 1 (VSBV-1) [2]. Other members are of veterinary
35 interest because they can cause severe disorders in birds, such as parrots [3]. Taxonomically, the
36 family *Bornaviridae* currently consists of the three genera *Orthobornavirus*, *Carbovirus* and
37 *Cultervirus* [4]. Of these, the orthobornaviruses have the widest so far known host spectrum and
38 have been identified in birds, reptiles and mammals [5]. Carbo- and culterviruses have up to now
39 only been identified in reptiles and fish, respectively [5–8]. The genus *Cultervirus* currently
40 comprises a single virus that has been discovered in fish (Wǔhàn sharpbelly bornavirus [WhSBV],
41 species *Cultervirus hemicultri*) [5, 8]. Partial genome sequences of another cultervirus, Murray-
42 Darling carp bornavirus (MDCBV), have recently been published, but its classification is still pending
43 [7].

44 The genome of bornaviruses consists of an approximately 9 kb non-segmented and single-stranded
45 RNA molecule of negative polarity (–ssRNA) [9]. Typically, six viral proteins are encoded by the viral
46 genome: nucleoprotein (N), accessory protein X, phosphoprotein (P), matrix protein (M),
47 glycoprotein (G) and the large protein (L) containing an RNA-directed RNA polymerase domain [5,
48 9]. The open reading frames (ORF) encoding these viral proteins are arranged in two known
49 genomic architectures: i) 3'-N-X-P-M-G-L-5' (genus *Orthobornavirus*) and ii) 3'-N-X-P-G-M-L-5'
50 (genera *Carbovirus* and *Cultervirus*). Bornaviral replication and transcription occur in the nucleus
51 of infected cells [10] and multiple viral transcripts are produced using conserved transcription
52 initiation and termination sites [11]. Atypically for mononegaviruses, bornaviruses use alternative
53 splicing in order to control and diversify their transcriptional capacity [12, 13].

54 Recently, we used an *in silico* data mining approach based on 'Serratus' [14] in order to screen for
55 traces of potential bornaviruses hidden in archived sequence data from public nucleic acid
56 sequences databases, such as the Sequence Read Archive (SRA) [15]. The SRA stores raw nucleic
57 acid sequence reads from next-generation sequencing runs from multidisciplinary research
58 experiments, along with extensive metadata. In these archived sequencing reads, we identified and
59 characterised two potential novel orthobornaviruses of colubrid and viperid snakes: Caribbean
60 watersnake bornavirus (CWBV) and Mexican black-tailed rattlesnake bornavirus (MRBV), in
61 datasets from a Caribbean watersnake (*Tretanorhinus variabilis*) and a Mexican black-tailed
62 rattlesnake (*Crotalus molossus nigrescens*), respectively [15].

63 In the present study, we extended the search for previously undetected bornaviruses by screening
64 116,082 transcriptomic datasets from fish samples from the orders Osteichthyes and
65 Chondrichthyes and identified seven bornavirus genomes.

66

67 MATERIAL AND METHODS

68 Selection of datasets

69 We generated a list of datasets using the European Nucleotide Archive (ENA) Browser advanced
70 search portal (<https://www.ebi.ac.uk/ena/browser/advanced-search>) and selected the data type
71 'raw reads' using the search query:

```
72 (tax_tree(1476529) OR tax_tree(7777) OR tax_tree(7898) OR tax_tree(7878)) AND  
73 (library_source="METATRANSCRIPTOMIC" OR library_source="TRANSCRIPTOMIC SINGLE CELL"  
74 OR library_source="VIRAL RNA" OR library_source="TRANSCRIPTOMIC")
```

75 Specifically, this search included the taxonomic units of jawless vertebrates (Cyclostomata;
76 NCBI:txid1476529), cartilaginous fishes (Chondrichthyes; NCBI:txid7777), ray-finned fishes
77 (Actinopterygii; NCBI:txid7898) and lungfish (Dipnomorpha; NCBI:txid7878). We further restricted
78 the search to RNA-derived datasets from (meta-) transcriptomic or viral RNA sequencing
79 experiments.

80

81 Data mining of raw reads

82 In order to identify even single reads within the selected datasets that may be related to
83 bornaviruses, we developed the bioinformatics pipeline 'SRMiner'. The 'SRMiner' pipeline is
84 based on *snakemake* [16], is multi-threading and can be run in any Linux-like environment. The
85 code for 'SRMiner' and detailed instructions on how to use it can be found at:
86 <https://gitlab.com/FPfaff/srminer>.

87 A simplified workflow of the pipeline includes the steps (i) download, (ii) blast and (iii) report: (i) A
88 subset of reads from each dataset is downloaded using *fastq-dump* (v3.0.3; SRA Toolkit). Typically,
89 a subset of 100,000 to 1,000,000 reads is sufficient to identify datasets containing sequence
90 reads of interest. (ii) Using *diamond blastx* (v2.0.15; [17]), the subset of reads is then searched
91 against a user-provided protein database. In this case, we selected and obtained the protein
92 sequences from all available members of the family *Bornaviridae* from NCBI. (iii) If at least a single
93 read matches the search criteria, additional metadata for this dataset is obtained using *ffq* (0.0.4;
94 [18]) and the results are summarised into individual reports using *R* [19] and *R markdown* [20].

95

96 Further raw read processing

97 After an initial screening of subsets of each 100,000 reads using SRMiner, the most promising
98 datasets were selected based on the number of reads matching the blastx search criteria and the
99 inferred theoretical number of reads in the full dataset. We then downloaded the full datasets of
100 these most promising SRA entries using *parallel-fastq-dump* (v0.6.7; [20]) and trimmed them for
101 low quality regions and adapter contamination using *TrimGalore!* (v0.6.10; [21]) running in
102 automatic mode. The trimmed reads were then used for *de novo* assembly with *SPAdes* genome

103 assembler (v3.15.5; [22]) running in `--rna` mode. The resulting transcripts/contigs were then
104 searched against the representative bornavirus protein database using *diamond blastx* (v2.0.15;
105 [17]). Transcripts/contigs matching the search criteria were selected and imported into Geneious
106 Prime (v2021.0.1) for further characterisation. Final genomes were additionally screened for any
107 vector or adapter contamination using the NCBI VecScreen suite
108 (<https://www.ncbi.nlm.nih.gov/tools/vecscreen>). To verify the nature of the sampled organism, we
109 selected all contigs from the assembly that matched the mitochondrial cytochrome B gene (*MT-*
110 *CYB*) and submitted them to the NCBI blastn suite (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>).

111

112 **Genomic characterization**

113 Potential ORFs were predicted using the Geneious Prime (v2021.0.1) 'Find ORFs' function and for
114 identification the deduced amino acid sequences were searched against the non-redundant blast
115 database (nr) using the NCBI blastp suite. The trimmed raw reads were mapped back to the
116 respective potential viral genome using the Geneious Prime (v2021.0.1) generic mapper (options:
117 medium sensitivity; find structural variants, short insertions, and deletions of any size) in order to
118 visualise transcriptional profiles and potential splice junctions. Potential transcription initiation and
119 termination sites were predicted based on sequence similarity to known bornaviral signal
120 sequences [11]. They were further verified by manual inspection of the read coverage at these
121 positions (e.g. transition to poly(A) at the termination sites). In addition to visual inspection of the
122 potential transcription start and termination sites, we used MEME (v5.5.2) to discover conserved
123 motifs.

124

125 **Genomic classification**

126 For the phylogenetic characterisation of potential bornavirus genomes, we used amino acid
127 alignments based on the predicted and translated N, G and L genes. The amino acid sequences of
128 these genes were individually aligned with 19 reference sequences using Muscle (v3.8.425). The
129 reference viruses were selected in order to represent all ICTV-accepted species of the family
130 *Bornaviridae* (n=12) as well as viruses below the species level (n=7). The individual alignments
131 were then concatenated into a single alignment and IQ-TREE (v2.2.2.6) was used to infer the
132 phylogenetic relationships. Specifically, a partitioning model (-Q) was used that allowed for
133 individual substitution models and evolutionary rates in each partition. The substitution model was
134 selected automatically (-m MFP+MERGE) and branch support was assessed using the ultrafast
135 bootstrap (-bb) and SH-aLRT tests (-alrt) with each 1,000,000 replicates each.

136 In addition, the Pairwise Sequence Comparison (PASC) [23, 24] was used to classify the potential
137 bornaviral genomes within the family *Bornaviridae*. PASC is based on pairwise global nucleotide
138 sequence alignments along the entire viral genome using a Blast-based approach.

139 For the prediction of the potential transmembrane domains (TM) along with the signal peptide (SP)
140 and cleavage sites (CS) within the G protein of the new genomes, we used DeepTMHMM (pybiolib,
141 version 1.1.944 [25]) and ProP-1.0 [26], respectively.

142

143 RESULTS

144 Data mining

145 During data mining, we analysed subsets of 116,078 raw transcriptomic SRA datasets from fish
146 (jawless vertebrates, cartilaginous fish, ray-finned fish and lungfish; see Supplementary Table S1).
147 In 72 of the 116,078 SRA datasets, we found at least one single read that matched one of the
148 bornavirus protein references. For all of these 72 datasets, a *de novo* assembly of all available data
149 was performed and the resulting contigs were scored (see Supplementary Table S2). In 8 of the *de*
150 *novo* assembled datasets, we identified endogenous bornavirus-like elements (EBLs), which were
151 not further analysed. In 4 datasets, we identified complete genomes from members of the viral
152 family *Chuviridae*. In a further 15 datasets, none of the resulting contigs showed any sequence
153 similarity to the bornavirus reference database. In 44 *de novo* assembled SRA datasets, full or
154 nearly full-length bornaviral genomes were identified. As some of these SRA datasets represented
155 either different organ samples from the same animal or multiple replicates belonging to a single
156 study or were based on the very same cell line, we selected only representative genomes for further
157 characterisation.

158 As a result, 7 complete and unique bornaviral genomes were assembled from SRA datasets
159 SRR10323915, SRR6207428, SRR1299086, SRR13236436, SRR9592747, SRR17661348,
160 and SRR17441645 (**Table 1**). The *MT-CYB* sequences assembled from each of these datasets
161 matched those of the specified sampled organisms (Supplementary Table S3).

162

163 Taxonomic relationship and classification

164 Phylogenetic analysis of the predicted viral proteins N, G, and L revealed that the potential
165 bornaviruses clustered with viruses of the genus *Cultervirus*, represented by WhSBV (NC_055169)
166 and MDCBV (MW645025-7), rather than carbo- or orthobornaviruses (**Figure 1**).

167 Specifically, the full genome derived from a dataset derived from the grass carp kidney cell line CIK
168 (SRR10323915) had 87.9% nucleotide identity to WhSBV (NC_055169) and was therefore
169 considered to be a variant of WhSBV. We identified the nearly identical WhSBV genome sequence
170 in 36 SRA datasets, all derived from RNA sequencing of either grass carp kidney (CIK; n=10) or liver
171 (L8824; n=26) cell lines (Supplementary Table S4).

172 In contrast, the full bornavirus genome from a goldfish tissue pool dataset (SRR6207428) showed
173 99.5% nucleotide identity to partial sequences of MDCBV (MW645025-7). This sequence can
174 therefore be considered the be the first complete genome of MDCBV.

175 Additional bornaviral sequences from Bombay duck fish (SRR17441645), electric eel
176 (SRR1299086), Pará molly (SRR17661348), finepatterned puffer (SRR13236436), and little skate
177 (SRR9592747) formed distinct taxonomic units. Hence, we tentatively named these potential
178 viruses based on the origin of the underlying sampling material: Bombay duck fish bornavirus
179 (BDBV; BK063658), electric eel bornavirus (EEBV; BK063519), Pará molly bornavirus (PMBV;
180 BK063657), finepatterned puffer bornavirus (FPBV; BK063517), and little skate bornavirus (LSBV;
181 BK063518). BDBV, EEBV, PMBV, and FPBV maintained between 42% and 66% PASC identity to
182 the known culterviruses WhSBV and MDCBV and to each other (Supplementary Figure S1). At
183 65.8%, BDBV and PMBV were more closely related to each other than to any other virus. LSBV
184 showed the greatest genetic divergence, with PASC identities ranging from 38.2% to 39.9% relative
185 to all other viruses.

186

187 **Genome architecture**

188 The genome architecture of MDCBV, BDBV, EEBV, PMBV, and FPBV was analogous to that of known
189 culter- and carboviruses, characterised by the arrangement of genes as 3'-N-X/P-G-M-L-5'
190 (**Figure 2**). The identified grass carp WhSBV variant, as well as the goldfish MDCBV variant, closely
191 resembled the WhSBV reference NC_055169 in structure and length (8,989 – 8,990 nt). In
192 contrast, BDBV, EEBV, PMBV, and FPBV had genome lengths of 9,110, 9,148, 9,324, and 9,397
193 nt, respectively. Notably, the genome structure of LSBV differed from the other bornaviral genomes,
194 as it was significantly longer, spanning 11,090 nt, and contained two additional ORFs designated
195 viral proteins 1 and 2 (Vp1 and Vp2): 3'-N-Vp1-Vp2-X/P-G-M-L-5'.

196

197 **Transcriptional profiles, motifs and alternative splicing**

198 The transcriptional profiles and splice sites of the discovered bornavirus genomes were
199 investigated by aligning/mapping the corresponding raw sequence data to the *de novo* assembled
200 genomes (**Figure 2**). The observed sequence coverage was not uniform across the genomes and
201 abrupt increases or decreases were observed within some of the potential intergenic regions. These
202 changes in genome coverage colocalised with predicted transcription start and termination motifs.
203 Specifically, the predicted start sites were characterised by a large increase in read coverage,
204 whereas the termination sites correlated with decrease in read coverage and the presence of reads
205 transitioning to poly(A) at the respective termination site. The respective positions of these
206 predicted regulatory sites were highly conserved between the different viruses. In detail, start sites
207 were present immediately upstream of the N, X/P, and M ORFs. The potential termination sites were

208 located downstream of the N, G and L ORFs. An additional termination site T3 was present within
209 the L ORF (**Figure 3**).

210 Genomic regions that showed homogeneous coverage and were flanked by adjacent start and
211 termination sites were interpreted as belonging to the same viral RNA transcripts or mRNA
212 (**Figure 3**). The overall pattern of viral transcription was highly conserved among all fish
213 bornaviruses analysed. In detail, the N protein appeared to be expressed from a monocistronic
214 mRNA, whereas X/P and G were expressed from a polycistronic mRNA. The M and L transcripts
215 appeared to share a single transcription start site (S3), but their expression levels were very
216 different, with L being expressed at low levels and M at relatively high levels.

217 Interestingly, LSBV showed an additional start and termination site, that were located adjacent to
218 the hypothetical ORFs of Vp1 and Vp2, suggesting that both proteins may be expressed from a
219 bicistronic mRNA. An additional intron was identified between the Vp1 and Vp2 ORFs at
220 nucleotide positions 1,868-2,476, which would result in an in-frame hybrid of the Vp1 and Vp2
221 ORFs, tentatively named Vp3 (see results below).

222 In addition, we identified an alternative splice site at the beginning of the L ORF, that was present
223 in all viruses analysed. The identified splice site was supported by multiple reads missing the
224 intronic sequence. The intron had a size of 110-176 nt and was located 23-53 nt downstream of
225 the M ORF stop codon. The coverage depth of the unspliced RNA was comparable to that of the L
226 ORF, while the spliced RNA had a coverage comparable to that of the M ORF (**Figure 3**). It could be
227 speculated that M is expressed from an RNA that undergoes alternative splicing and uses the T3
228 transcription termination site located within the L ORF (**Figure 4A**). The viral RNA for L on the other
229 hand is expressed from the same S3 transcription start as M but does not undergo splicing and
230 uses a the T4 termination site. The intronic sequences of all viruses analysed showed the canonical
231 dinucleotides GU and AG for donor and acceptor sites, respectively (**Figure 4A**).

232 Motif prediction revealed conserved sequence patterns for transcription termination and start sites
233 (**Figure 4B**). The termination sites T1-4 shared the conserved nucleotide sequence pattern
234 'AYUUWAKAAAAACAU', whereas the start sites S1, S2 and S3 shared the conserved nucleotide
235 sequence pattern 'GAM'. S2 and S3 were immediately adjacent to T1 and T2, respectively.

236

237 **LSBV Vp1 and Vp2 are homologues of the glycoprotein G**

238 When analysed by pairwise alignment, the hypothetical viral proteins Vp1 and Vp2 of LSBV shared
239 amino acid similarity with the glycoprotein G of LSBV (**Figure 5A**). In detail, the pairwise amino acid
240 identity between the G and Vp1 was 28%, between the G and Vp2 it was 21% and between Vp1
241 and Vp2 it was 41% (**Figure 5B**). While Vp1 shares the N-terminus of the G protein, it lacks the C
242 terminus. Vp2 shares only the central region of the G protein and lacks both, the respective N- and
243 C-terminal regions of G. Both, Vp1 and Vp2, have no detectable transmembrane domain, as they

244 lack the respective C-terminal part of the G protein (470-491 aa; **Figure 5A** and Supplementary
245 Table S5).

246 A phylogenetic tree was constructed based on an amino acid alignment of glycoproteins from
247 selected members of the *Bornaviridae* family, supplemented by LSBV Vp1 and Vp2 (see
248 Supplementary Figure S2). The tree provided evidence for the occurrence of a duplication event of
249 the LSBV G gene, with Vp1 sharing the last common ancestor with G and Vp2 sharing the last
250 common ancestor with Vp1. One possible scenario could be that initially a large part of the
251 glycoprotein gene G was duplicated to form Vp1 and later only the part encoding the C terminus of
252 Vp1 (representing the central part of the G) was duplicated to form Vp2 (**Figure 5C**). We also
253 predicted potential furin endoprotease cleavage sites within Vp1, Vp2 and G, following the amino
254 acid consensus motif 'RS (K/R) R' (**Figure 5C** and Supplementary Table S5).

255 As noted above, the predicted mRNA encoding both Vp1 and Vp2, may also undergo splicing,
256 resulting in a hybrid ORF, tentatively named Vp3* (**Figure 5D**). The potential Vp3* protein would
257 consist of the N-terminal portion of Vp1 and the C-terminal portion of Vp2, including the protease
258 cleavage site. Similar to Vp1 and Vp2, the potential Vp3* would lack a transmembrane domain
259 (Supplementary Table S5).

260

261 **DISCUSSION**

262 Knowledge on fish bornaviruses has been limited to a single full-length genome of WhSBV [8] and
263 a partial genome of MDCBV [7]. To identify additional and more diverse fish bornaviruses, we used
264 an *in silico* data mining approach that screened publicly available SRA raw sequence datasets from
265 fish (Osteichthyes and Chondrichthyes) samples. Using a similar approach, we had previously
266 successfully identified and characterised two novel snake orthobornaviruses, CWBV and MRBV, as
267 well as novel EBLs in reptile datasets [15]. Here, the screening combined with *de novo* assembly
268 led to the identification of five putative complete bornavirus genomes from different samples.

269 We found additional sequences of WhSBV (87.9% nt identity to the previously published sequence)
270 and MDCBV (99.5% nt identity) in fish other than the originally reported host species. The first full-
271 length genome sequence of MDCBV presented here matched that of WhSBV in overall structure,
272 sequence identity, and length, indicating that MDCBV and WhSBV are closely related. According to
273 the criteria defined by the ICTV *Bornaviridae* Study Group [5], they are thought to be viruses of the
274 same virus species (*Cultervirus hemicultri*). WhSBV was previously identified by RNA sequencing of
275 the gut, liver, and gill tissues from a sharpbelly or wild carp (*Hemiculter leucisculus*; family
276 Cyprinidae) from China [8], whereas MDCBV was discovered in a liver and gill tissue pool of a
277 common carp (*Cyprinus carpio*; family Cyprinidae) during a meta-transcriptomic survey of
278 freshwater species in the Murray-Darling Basin in Australia [7]. Here, we identified WhSBV in
279 multiple datasets from cell lines derived from the kidney and liver of a grass carp

280 (*Ctenopharyngodon idella*; family Cyprinidae) and MDCBV in a dataset from goldfish (*Carassius*
281 *auratus*; family Cyprinidae) brain samples. Both, WhSBV and MDCBV thus appear to be members
282 of a group of bornaviruses that are particularly common in fishes of the family Cyprinidae.
283 Cyprinidae includes a wide range of carp and is an ancient evolutionary lineage [27]. With a global
284 production of ~30 million tonnes [28], carps are of great economic interest and are often cultivated
285 in large-scale aquaculture farms. Therefore, the impact of these bornaviruses on animal health
286 needs to be carefully assessed and the genome sequences identified in this study may provide
287 valuable information to further investigate the distribution and variability of these viruses.

288 Using the data mining approach, identical WhSBV genomes were identified in datasets from the
289 grass carp cell lines CIK (kidney) and L8824 (liver). Both cell lines originate from the Freshwater
290 Fisheries Research Center of Chinese Academy of Fishery Sciences (formerly the Yangtze River
291 Fisheries Research Institute) [29]. The CIK and L8824 cell lines have been repeatedly used to study
292 viral transcriptional changes during infection, e.g. with grass carp reovirus (GCRV), and immune
293 regulation. The presence of WhSBV in samples labelled 'mock infection' or 'cell control' (see
294 Supplementary Table S4) indicates that both cell lines may be persistently infected and
295 allel experimental results from experiments should be interpreted with caution. It remains unclear
296 whether the WhSBV found in these cell lines originated from the individual(s) from which the two
297 cell lines were derived, or whether both cell lines may have been subsequently contaminated.

298 We also identified four additional bornaviral genomes in non-cyprinid ray-finned fishes, and one in
299 a cartilaginous fish. These viruses were related to WhSBV and MDCBV, but formed clearly separate
300 taxonomic entities based on a phylogenetic analysis of N, G and L protein sequences. Despite clear
301 differences at the nucleotide and amino acid level, four of these viruses shared the same overall
302 genomic structure with the culterviruses WhSBV and MDCBV, and with the viruses of the genus
303 *Carbovirus* [6]. The genome arrangement of reptilian carboviruses and these novel fish
304 bornaviruses is peculiar in that it does not follow the standard N-X-P-M-G-L pattern of
305 mononegaviruses in general and of orthobornaviruses in particular. This could indicate that the N-
306 X-P-G-M-L genome arrangement evolved independently in reptile and fish bornaviruses, or that they
307 share an ancient common ancestor that already had this genome architecture. However, assuming
308 a virus/host co-evolution, the question arises as to why orthobornaviruses (hosts: birds, mammals,
309 reptiles) have conserved the typical N-X-P-M-G-L pattern of mononegaviruses, although they should
310 be comparatively younger than culterviruses (hosts: fish). It is therefore reasonable to assume that
311 bornavirus evolution did not follow a strict virus/host co-evolution and that ancient ancestors of
312 orthobornaviruses infected a wider range of vertebrates than extant orthobornaviruses, as it has
313 been suggested by analysis of endogenous bornavirus-like elements (EBLs) [30].
314 Orthobornaviruses may therefore represent a more ancient lineage of bornaviruses, as they show
315 the typical N-X-P-M-G-L pattern of mononegaviruses, while the N-X-P-G-M-L pattern of carbo- and
316 culterviruses may be a more recent development.

317 The rearrangement of G and M may have resulted in a favourable regulation of gene expression for
318 these viruses. By analysing the transcriptional profiles of the novel fish bornaviruses, we found that
319 X/P and G are most likely co-expressed from the same polycistronic mRNA and M is transcribed
320 from a spliced mRNA. In contrast, orthobornaviruses express X and P from a bicistronic mRNA
321 starting from transcription start site S2, whereas M and G are expressed from different splice
322 variants of mRNAs starting from S3 [31].

323 Genomic rearrangements do not seem to be an isolated event in bornaviruses, as illustrated by the
324 unique genome architecture of LSBV, which encoded two more possible ORFs (Vp1 and Vp2). Both
325 appeared to be the result of at least two independent duplication events: First, a large part of the
326 G gene appears to have been copied into the intergenic region between N and X/P, forming Vp1.
327 Subsequently, a part of Vp1 was duplicated into the intergenic region between Vp1 and X/P,
328 forming Vp2 (**Figure 5**). Comparable duplication events in RNA viruses are considered very rare
329 [32], but have been reported for other mononegaviruses, such as rhabdoviruses [32–
330 35]. Exceptionally long branches in the phylogenetic analysis indicated an accelerated evolution for
331 Vp1 and Vp2 after the duplication events, possibly as a result of changing evolutionary context and
332 selection pressure [36].

333 In addition, the Vp1 and Vp2 genes may produce a hybrid gene product Vp3* by alternative splicing,
334 further extending the coding potential of LSBV even further. The function of Vp1, Vp2 and the splice
335 hybrid Vp3* is currently unknown, but conserved furin cleavage sites suggest that these proteins
336 undergo some form of post-translational modification, similar to the glycoprotein of other
337 bornaviruses [37]. As Vp1, Vp2 and Vp3* lack a detectable transmembrane domain, it can be
338 speculated that they may could function as soluble glycoproteins, similar to that of vesicular
339 stomatitis virus [38]. The predicted cleavage site within the Vp1, Vp2, and Vp3* sequences may
340 have functional significance for the virus, and future experimental investigations are needed to gain
341 deeper insights into the unique genome architecture of this bornavirus. It would be very interesting
342 to investigate whether other bornaviruses from cartilaginous fish share this unique genome
343 structure, or whether LSBV is the result of an isolated evolutionary event.

344 Based on the phylogenetic analysis and PASC, we propose that LSBV does not belong to any of the
345 existing genera within the family *Bornaviridae*. We have therefore submitted a taxonomic proposal
346 to the ICTV to establish a new genus within this family. In this proposal, EEBV and FPBV were also
347 tentatively classified as four new species within the genus *Cultervirus*.

348 Although the combination of gene arrangement, expression profile and potential hosts was
349 plausible for these potential viruses, it cannot be excluded that these genomes were based on
350 contaminated samples or inaccurate datasets and therefore did not originate from the reported
351 host species. However, the identification of known viruses such as WhSBV and MDCBV in fish
352 datasets related to the originally reported host, may support the credibility of our findings.
353 Confirmation by standard methods, such as PCR and virus isolation, using independent samples

354 from the same species would nevertheless be required to fully confirm the existence of these
355 interesting new viruses in the reported hosts.

356

357 **CONCLUSION**

358 Until now, WhSBV and the closely related MDCBV were the only viruses in the genus *Cultervirus*.
359 Screening of 116,078 fish datasets from the SRA led to the identification of six tentative cultervirus
360 genomes, including a variant of WhSBV and the first complete genome of MDCBV. These viruses
361 may primarily infect carp, as all variants of WhSBV and MDCBV have so far been found in cyprinid
362 samples. In addition, BDBV, EEBV, PMBV, and FPBV were discovered in a dataset from a Bombay
363 duck, an electric eel, a Pará molly and a finned pufferfish, respectively. They had comparatively
364 longer genomes, they all shared the same genome organisation 5'-N-X/P-G-M-L-3', similar
365 transcriptional profiles and regulatory sites, thus, suggesting a common ancestor of these fish
366 bornaviruses. In addition, LSBV was identified in a little skate dataset and showed a distinct
367 genome organisation with two additional genes that may be the result of ancient duplication events
368 of the glycoprotein gene. The LSBV had the largest genome length (11,090 nt) of any bornavirus
369 known to date and, to our knowledge, the presence of duplicated genes within a virus of the
370 *Mononegaviridae* family is quite unique and has so far only been reported for a few rhabdoviruses.
371 The study demonstrates the power of *in silico* SRA data screening and its ability to advance the
372 knowledge of viral diversity and evolution. The screening can easily be applied to the discovery of
373 novel viruses from other viral families, or to the identification of known viruses in datasets from
374 previously unknown potential host species.

375

376 **ACKNOWLEDGEMENT**

377 We are grateful to the global scientific community for their invaluable contribution in sharing raw
378 sequencing data. These datasets, originally intended for specific research purposes, contain
379 valuable information that goes beyond their original purposes. We would also like to thank Robin
380 Garcia, Jörg Linde and Michael Weber for their help with the 'SRAMiner' pipeline.

381

382 **DATA AVAILABILITY**

383 Nucleotide sequence data reported are available in the Third Party Annotation Section of the
384 DDBJ/ENA/GenBank databases under the TPA accession numbers: BK063517-BK063521,
385 BK063657 and BK063658.

386

387

388 **FUNDING**

389 The investigations were supported by a Friedrich-Loeffler-Institut internal PhD program, grant
390 number FLI-IVD-XX-2021-83, granted to F.P.

391

392 **CONFLICT OF INTEREST**

393 The authors declare no conflict of interest. The funders had no role in the design of the study, in
394 the collection, analyses, or interpretation of data, in the writing of the manuscript, or in the decision
395 to publish the results.

396

397 **REFERENCES**

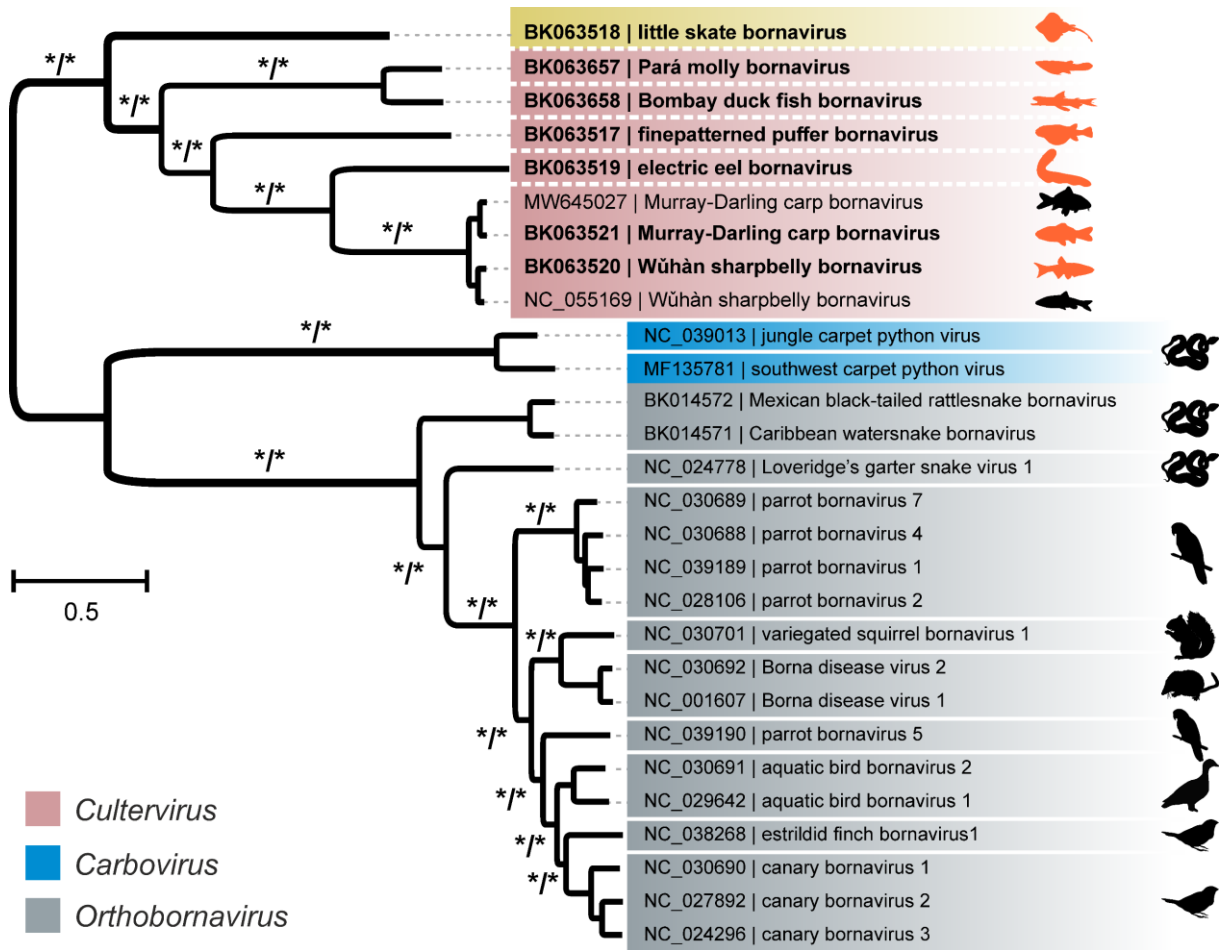
- 398 1. Niller HH, Angstwurm K, Rubbenstroth D et al. (2020) Zoonotic spillover infections with
399 Borna disease virus 1 leading to fatal human encephalitis, 1999-2019: an epidemiological
400 investigation. *Lancet Infect Dis* 20:467–477. [https://doi.org/10.1016/S1473-](https://doi.org/10.1016/S1473-3099(19)30546-8)
401 [3099\(19\)30546-8](https://doi.org/10.1016/S1473-3099(19)30546-8)
- 402 2. Hoffmann B, Tappe D, Höper D et al. (2015) A Variegated Squirrel Bornavirus Associated
403 with Fatal Human Encephalitis. *N Engl J Med* 373:154–162.
404 <https://doi.org/10.1056/NEJMoa1415627>
- 405 3. Rubbenstroth D (2022) Avian Bornavirus Research-A Comprehensive Review. *Viruses* 14.
406 <https://doi.org/10.3390/v14071513>
- 407 4. Kuhn JH, Dürrwald R, Bào Y et al. (2015) Taxonomic reorganization of the family
408 *Bornaviridae*. *Arch Virol* 160:621–632. <https://doi.org/10.1007/s00705-014-2276-z>
- 409 5. Rubbenstroth D, Briese T, Dürrwald R et al. (2021) ICTV Virus Taxonomy Profile:
410 *Bornaviridae*. *J Gen Virol* 102. <https://doi.org/10.1099/jgv.0.001613>
- 411 6. Hyndman TH, Shilton CM, Stenglein MD et al. (2018) Divergent bornaviruses from Australian
412 carpet pythons with neurological disease date the origin of extant *Bornaviridae* prior to the
413 end-Cretaceous extinction. *PLoS Pathog* 14:e1006881.
414 <https://doi.org/10.1371/journal.ppat.1006881>
- 415 7. Costa VA, Mifsud JCO, Gilligan D et al. (2021) Metagenomic sequencing reveals a lack of
416 virus exchange between native and invasive freshwater fish across the Murray-Darling Basin,
417 Australia. *Virus Evol* 7:veab034. <https://doi.org/10.1093/ve/veab034>
- 418 8. Shi M, Lin X-D, Chen X et al. (2018) The evolutionary history of vertebrate RNA viruses.
419 *Nature* 556:197–202. <https://doi.org/10.1038/s41586-018-0012-7>

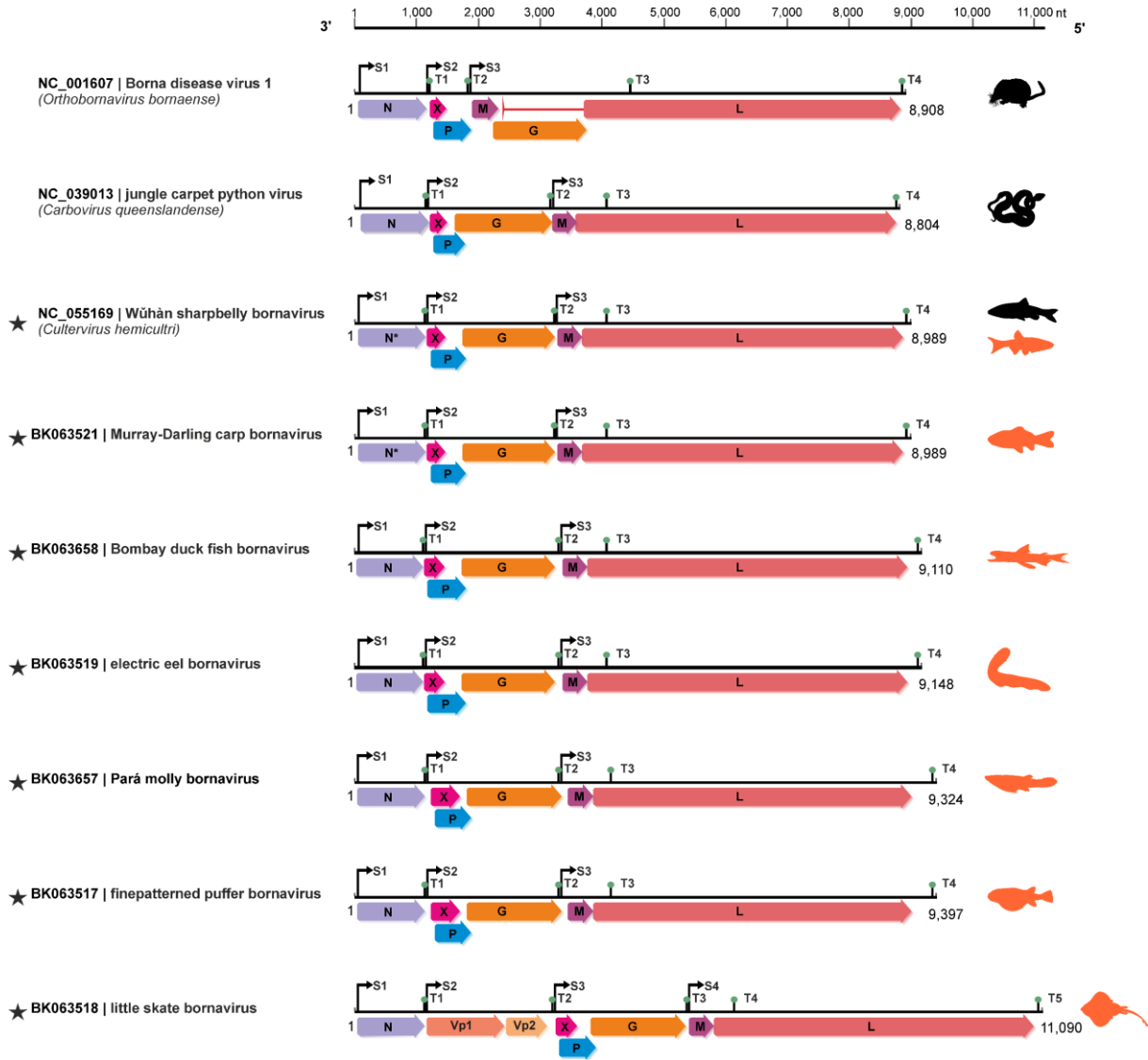
- 420 9. Briese T, Schneemann A, Lewis AJ et al. (1994) Genomic organization of Borna disease virus.
421 Proc Natl Acad Sci U S A 91:4362–4366. <https://doi.org/10.1073/pnas.91.10.4362>
- 422 10. Briese T, La Torre JC de, Lewis A et al. (1992) Borna disease virus, a negative-strand RNA
423 virus, transcribes in the nucleus of infected cells. Proc Natl Acad Sci U S A 89:11486–
424 11489. <https://doi.org/10.1073/pnas.89.23.11486>
- 425 11. Schneemann A, Schneider PA, Kim S et al. (1994) Identification of signal sequences that
426 control transcription of borna disease virus, a nonsegmented, negative-strand RNA virus. J
427 Virol 68:6514–6522. <https://doi.org/10.1128/JVI.68.10.6514-6522.1994>
- 428 12. Schneider PA, Schneemann A, Lipkin WI (1994) RNA splicing in Borna disease virus, a
429 nonsegmented, negative-strand RNA virus. J Virol 68:5007–5012.
430 <https://doi.org/10.1128/JVI.68.8.5007-5012.1994>
- 431 13. Tomonaga K, Kobayashi T, Lee BJ et al. (2000) Identification of alternative splicing and
432 negative splicing activity of a nonsegmented negative-strand RNA virus, Borna disease virus.
433 Proc Natl Acad Sci U S A 97:12788–12793. <https://doi.org/10.1073/pnas.97.23.12788>
- 434 14. Edgar RC, Taylor J, Lin V et al. (2022) Petabase-scale sequence alignment catalyses viral
435 discovery. Nature 602:142–147. <https://doi.org/10.1038/s41586-021-04332-2>
- 436 15. Pfaff F, Rubbenstroth D (2021) Two novel bornaviruses identified in colubrid and viperid
437 snakes. Arch Virol 166:2611–2614. <https://doi.org/10.1007/s00705-021-05138-3>
- 438 16. Mölder F, Jablonski KP, Letcher B et al. (2021) Sustainable data analysis with Snakemake.
439 F1000Res 10:33. <https://doi.org/10.12688/f1000research.29032.1>
- 440 17. Buchfink B, Reuter K, Drost H-G (2021) Sensitive protein alignments at tree-of-life scale
441 using DIAMOND. Nat Methods 18:366–368. <https://doi.org/10.1038/s41592-021-01101-x>
- 442 18. Gálvez-Merchán Á, Min KHJ, Pachter L et al. (2023) Metadata retrieval from sequence
443 databases with ffq. Bioinformatics 39. <https://doi.org/10.1093/bioinformatics/btac667>
- 444 19. R Core Team (2022) R: A Language and Environment for Statistical Computing. R Foundation
445 for Statistical Computing, Vienna, Austria
- 446 20. (2023) rmarkdown: Dynamic Documents for R
- 447 21. Trim Galore
- 448 22. Bushmanova E, Antipov D, Lapidus A et al. (2019) rnaSPAdes: a de novo transcriptome
449 assembler and its application to RNA-Seq data. Gigascience 8.
450 <https://doi.org/10.1093/gigascience/giz100>
- 451 23. Bao Y, Chetvernin V, Tatusova T (2014) Improvements to pairwise sequence comparison
452 (PASC): a genome-based web tool for virus classification. Arch Virol 159:3293–3304.
453 <https://doi.org/10.1007/s00705-014-2197-x>

- 454 24. Bao Y, Kapustin Y, Tatusova T (2008) Virus Classification by Pairwise Sequence Comparison
455 (PASC). In: Encyclopedia of Virology. Elsevier, pp 342–348
- 456 25. Hallgren J, Tsigos KD, Pedersen MD et al. (2022) DeepTMHMM predicts alpha and beta
457 transmembrane proteins using deep neural networks
- 458 26. Duckert P, Brunak S, Blom N (2004) Prediction of proprotein convertase cleavage sites.
459 Protein Eng Des Sel 17:107–112. <https://doi.org/10.1093/protein/gzh013>
- 460 27. Cavender TM The fossil record of the Cyprinidae. In: Winfield, Nelson (Hg.) 1991 – Cyprinid
461 fishes, pp 34–54
- 462 28. FAO (ed) (2021) FAO yearbook: Fishery and Aquaculture Statistics 2019/FAO annuaire.
463 Statistiques des pêches et de l'aquaculture 2019/FAO anuario. Estadísticas de pesca y
464 acuicultura 2019, Rome/Roma
- 465 29. Wengong Z, Huaxin Q, Yingfang X et al. (1986) A cell line derived from the kidney of grass
466 carp (*Ctenopharyngodon idellus*). Journal of Fisheries of China 10:10–17
- 467 30. Kawasaki J, Kojima S, Mukai Y et al. (2021) 100-My history of bornavirus infections hidden
468 in vertebrate genomes. Proc Natl Acad Sci U S A 118.
469 <https://doi.org/10.1073/pnas.2026235118>
- 470 31. (2012) Bornaviridae. In: Virus Taxonomy. Elsevier, pp 658–664
- 471 32. Simon-Loriere E, Holmes EC (2013) Gene duplication is infrequent in the recent evolutionary
472 history of RNA viruses. Mol Biol Evol 30:1263–1269.
473 <https://doi.org/10.1093/molbev/mst044>
- 474 33. Wang Y, Walker PJ (1993) Adelaide river rhabdovirus expresses consecutive glycoprotein
475 genes as polycistronic mRNAs: new evidence of gene duplication as an evolutionary process.
476 Virology 195:719–731. <https://doi.org/10.1006/viro.1993.1423>
- 477 34. Gubala AJ, Proll DF, Barnard RT et al. (2008) Genomic characterisation of Wongabel virus
478 reveals novel genes within the Rhabdoviridae. Virology 376:13–23.
479 <https://doi.org/10.1016/j.virol.2008.03.004>
- 480 35. Gubala A, Davis S, Weir R et al. (2010) Ngaingan virus, a macropod-associated rhabdovirus,
481 contains a second glycoprotein gene and seven novel open reading frames. Virology
482 399:98–108. <https://doi.org/10.1016/j.virol.2009.12.013>
- 483 36. Lynch M, Conery JS (2000) The evolutionary fate and consequences of duplicate genes.
484 Science 290:1151–1155. <https://doi.org/10.1126/science.290.5494.1151>
- 485 37. Richt JA, Fürbringer T, Koch A et al. (1998) Processing of the Borna disease virus
486 glycoprotein gp94 by the subtilisin-like endoprotease furin. J Virol 72:4528–4533.
487 <https://doi.org/10.1128/JVI.72.5.4528-4533.1998>

- 488 38. Graeve L, Garreis-Wabnitz C, Zauke M et al. (1986) The soluble glycoprotein of vesicular
489 stomatitis virus is formed during or shortly after the translation process. *J Virol* 57:968–975.
490 <https://doi.org/10.1128/JVI.57.3.968-975.1986>
- 491 39. Chen G, He L, Luo L et al. (2018) Transcriptomics Sequencing Provides Insights into
492 Understanding the Mechanism of Grass Carp Reovirus Infection. *Int J Mol Sci* 19.
493 <https://doi.org/10.3390/ijms19020488>
- 494 40. Shan B, Liu Y, Yang C et al. (2021) Comparative Transcriptome Analysis of Female and Male
495 Fine-Patterned Puffer: Identification of Candidate Genes Associated with Growth and Sex
496 Differentiation. *Fishes* 6:79. <https://doi.org/10.3390/fishes6040079>
- 497 41. Da Fonte DF, Martyniuk CJ, Xing L et al. (2017) Secretoneurin A regulates neurogenic and
498 inflammatory transcriptional networks in goldfish (*Carassius auratus*) radial glia. *Sci Rep*
499 7:14930. <https://doi.org/10.1038/s41598-017-14930-8>
- 500 42. Gallant JR, Traeger LL, Volkening JD et al. (2014) Nonhuman genetics. Genomic basis for the
501 convergent evolution of electric organs. *Science* 344:1522–1525.
502 <https://doi.org/10.1126/science.1254432>
- 503

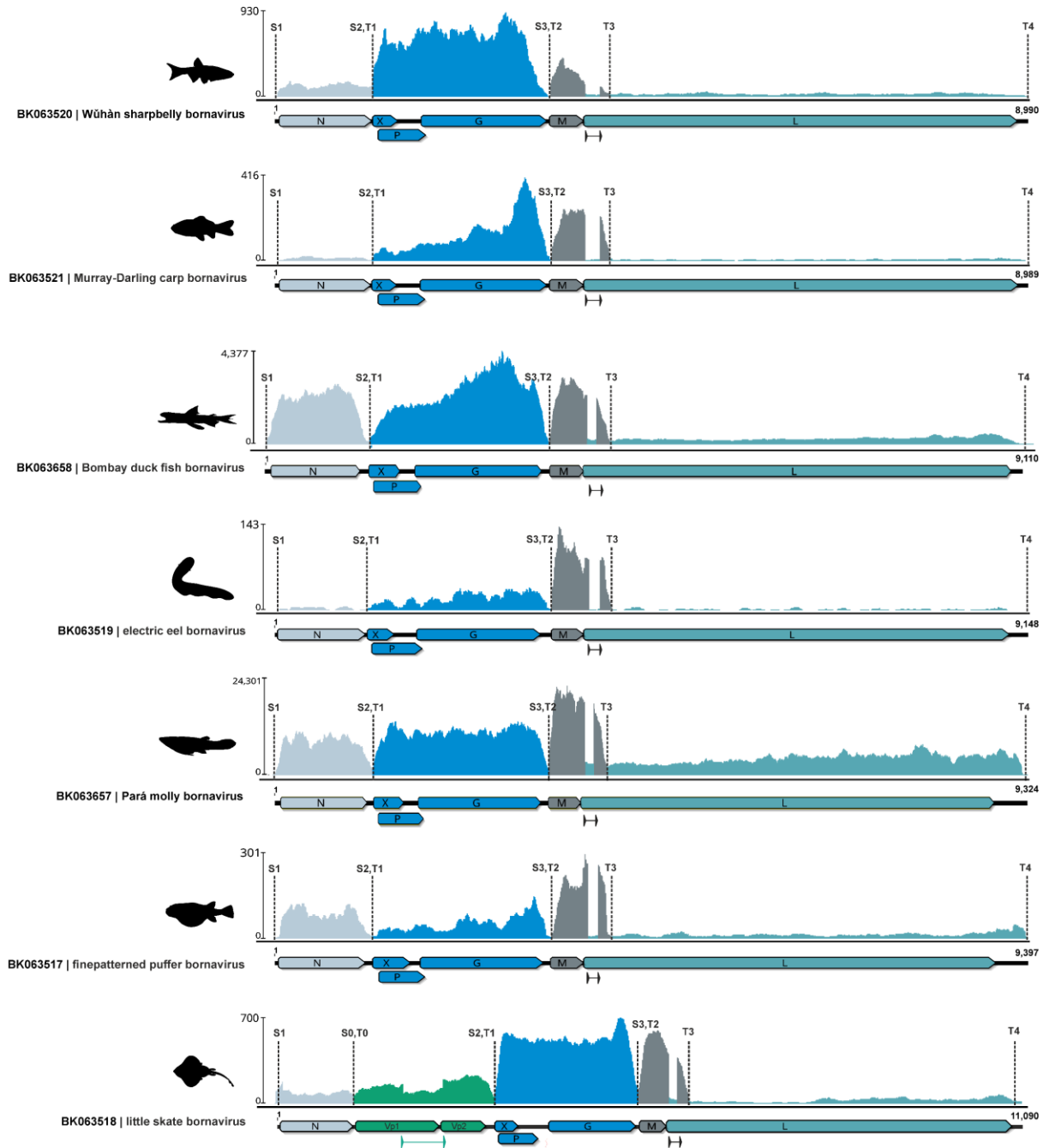
504 FIGURES





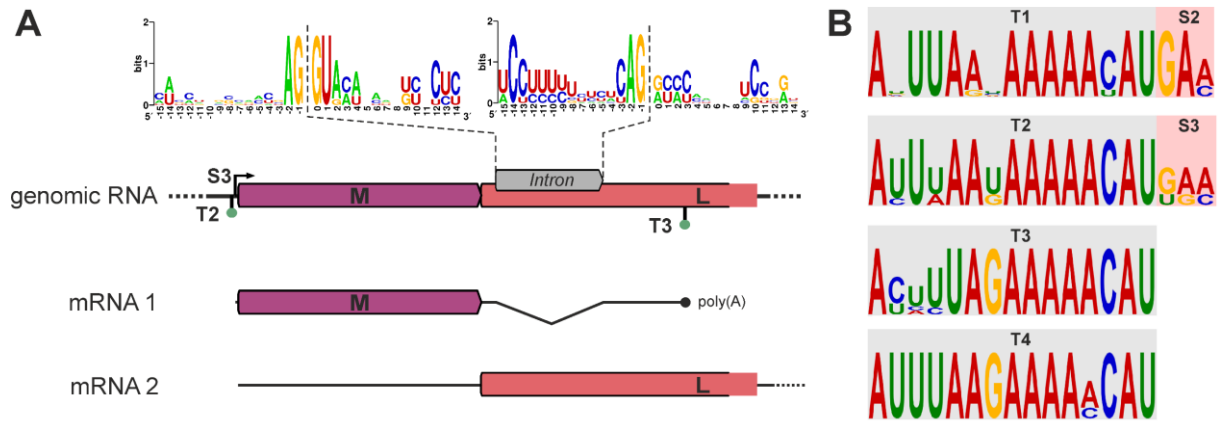
515

516 **Figure 2: Genome architectures of current and potential novel bornaviruses.** Representative overall genome
 517 organisations are shown for representative bornaviruses along with the potential novel viruses (star). The
 518 (predicted) open reading frames (ORF) are shown as arrows together with the predicted transcription start
 519 (S) and transcription termination (T) sites. For each of the genomes, the potential hosts/sources for each
 520 virus are shown. Note the different genomic arrangements: 3'-N-X-P-M-G-L-5' (genus *Orthobornavirus*) and
 521 3'-N-X-P-G-M-L-5' (genera *Carbovirus* and *Cultervirus*). The little skate bornavirus shares the genomic
 522 structure of carbo- and culterviruses, but encodes two additional predicted ORFs: 3'-N-Vp1-Vp2-X-P-M-G-L-5'.



523

524 **Figure 3: Transcriptional profiles of novel bornaviruses.** Raw reads were mapped to the novel bornavirus
 525 genomes and the coverage was plotted. Open reading frames (ORFs) are shown as arrows and predicted
 526 transcription start (S) and transcription termination (T) motifs are indicated as dashed lines. S and T sites
 527 collocate with large increases and decreases in coverage, respectively. Regions that with similar coverage
 528 and are bordered by S and T sites were considered to represent individual RNA transcripts. These viral
 529 transcripts and their corresponding ORFs are highlighted in different colours. Alternative splicing was
 530 detected within in all viruses for potential M transcript (intron shown as line arrow). In addition, a potential
 531 intron was identified in the bicistronic transcript encoding Vp1 and Vp2 of little skate bornavirus.



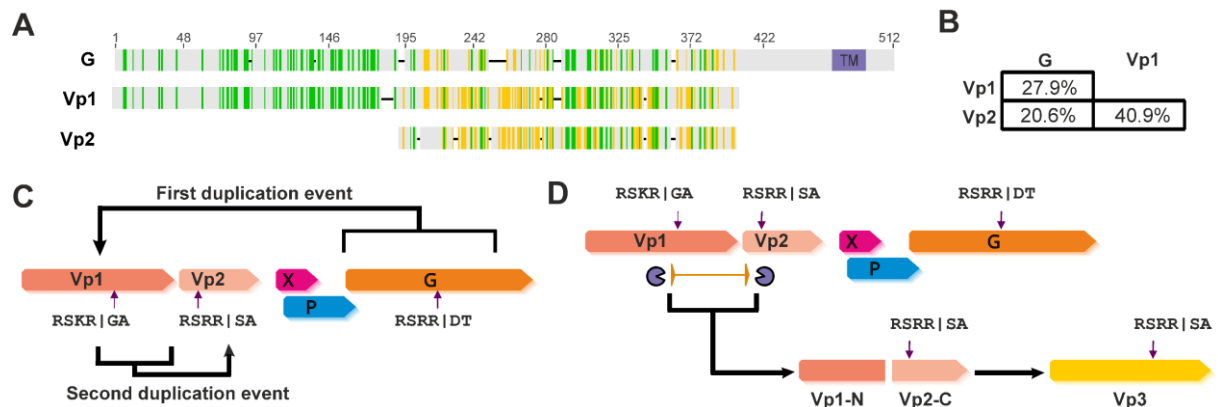
532

533 **Figure 4: Conserved splicing mechanisms and transcriptional motifs among fish bornaviruses.** (A) Alternative
 534 splicing was detected in the M/L ORF region for all viruses analysed. The genomic arrangement in this region
 535 is shown with ORFs indicated by arrows. The potential intron is shown as an arrow within the L ORF region.
 536 The sequence motifs of the splice acceptor and donor sites from all analysed viruses are shown and the
 537 dashed lines indicate the position of the splicing. The canonical GU/AG splice site is present in all viruses
 538 analysed. Two possible mRNAs are shown: The spliced mRNA contains only the M ORF and is terminated at
 539 the T3 site, while the unspliced mRNA will contains the full L ORF. (B) The conserved motifs of the
 540 transcription termination and start sites of the analysed viruses are shown. Note that T1/S2 and T2/S3 are
 541 directly adjacent to each other.

542

543

544



545

546 **Figure 5: Little skate bornavirus encodes two proteins that may be the result of an ancient duplication event**
 547 **of the glycoprotein.** (A) The amino acid alignment of LSBV viral proteins 1 and 2 (Vp1 and Vp2) together with
 548 the glycoprotein (G) shows, that they share homology. Vp1 and Vp2 lack the corresponding transmembrane
 549 domain (TM) of G, but each contain a predicted furin protease cleavage site (highlighted by scissors symbol).
 550 (B) Pairwise amino acid identities indicate, that Vp1 and Vp2 are more closely related to each other than to
 551 G. Therefore, in (C), supported by phylogenetic analysis (see also the Supplementary Figure S3), we
 552 hypothesised that Vp1 was first duplicated from G, followed by a second duplication of Vp1, which gave rise
 553 to Vp2. Predicted cleavage sites are indicated by arrows. (D) Transcriptional profiling suggested the possibility
 554 of alternative splicing of Vp1 and Vp2, resulting in a hybrid of the Vp1 C-terminus, including its cleavage site,
 555 and the Vp2 N-terminus, tentatively named Vp3*.

556 TABLES

557 **Table 1:** Summary of SRA datasets that were selected for *de novo* assembly of complete bornaviral genomes

SRA Accession	Sampled organism	Sampled material	<i>de novo</i> assembled virus	Reads matching viral genome
SRR10323915	grass carp <i>Ctenopharyngodon Idella</i> (Valenciennes, 1844)	permanent kidney cell line (CIK) [29, 39]	Wùhàn sharpbelly bornavirus WhSBV BK063520	311,433 (0.515%)
SRR6207428	goldfish <i>Carassius auratus</i> (Linnaeus, 1758)	tissue pool of adult male [40]	Murray-Darling carp bornavirus MDCBV BK063521	90,150 (0.123%)
SRR1299086	electric eel <i>Electrophorus electricus</i> (Linnaeus, 1766)	ampullae of Lorenzini tissue of an adult female	electric eel bornavirus EEBV BK063519	132,721 (0.046%)
SRR13236436	finepatterned puffer <i>Takifugu poecilonotus</i> (Temminck & Schlegel, 1850)	radial glial cells from the brain of an adult female [41]	finepatterned puffer bornavirus FPBV BK063517	9,469 (0.022%)
SRR9592747	little skate <i>Leucoraja erinacea</i> (Mitchill, 1825)	kidney tissue [42]	little skate bornavirus LSBV BK063518	37,573 (0.217%)
SRR17661348	Pará molly <i>Poecilia parae</i> (Eigenmann, 1894)	head of an adult female	Pará molly bornavirus PMBV BK063657	615,854 (0.38%)
SRR17441645	Bombay duck fish <i>Harpadon nehereus</i> (Hamilton, 1822)	gill tissue	Bombay duck fish bornavirus BDBV BK063658	65,661 (0.154%)

558