

1 **Title: The evolution, distribution and diversity and of endogenous circoviral**
2 **elements.**

3

4 **Running title: Endogenous circoviral elements**

5

6 **Authors:** Tristan P.W. Dennis^{1*}, William Marciel de Souza^{1,2*}, Soledad Marsile-
7 Medun^{1,3}, Joshua B. Singer¹, Sam J. Wilson¹, and Robert J. Gifford

8

9 **Affiliations:**

10 1. MRC-University of Glasgow Centre for Virus Research, 464 Bearsden Road,
11 Glasgow, UK

12 2. Virology Research Center, School of Medicine of Ribeirão Preto of University
13 of São Paulo, Ribeirão Preto, Brazil

14 3. Agrocampus Ouest, 65 Rue de Saint-Brieuc, 35000 Rennes, France

15

16 **Corresponding author:**

17

18 Robert J. Gifford,

19 MRC-University of Glasgow Centre for Virus Research,

20 Bearsden,

21 Glasgow, UK.

22 e-mail: robert.gifford@glasgow.ac.uk

23

1 **Abstract**

2 Circoviruses (family *Circoviridae*) are small, non-enveloped viruses that have
3 short, single-stranded DNA genomes. Circovirus sequences are frequently recovered
4 in metagenomic investigations, indicating that these viruses are widespread, yet they
5 remain relatively poorly understood. Endogenous circoviral elements (CVe) are DNA
6 sequences derived from circoviruses that occur in vertebrate genomes. CVe can
7 provide unique, retrospective information about the biology and evolution of
8 circoviruses. In this study, we screened 362 vertebrate genome assemblies *in silico*
9 to generate a catalog of CVe loci. We identified a total of 179 CVe sequences, most
10 of which have not been reported previously. We show that these CVe loci reflect at
11 least 19 distinct germline integration events. We determine the structure of CVe loci,
12 identifying some that show evidence of potential functionalization. We also identify
13 orthologous copies of CVe in snakes, fish, birds, and mammals, allowing us to add
14 new calibrations to the timeline of circovirus evolution. Finally, we observed that
15 some ancient CVe group robustly with contemporary circoviruses in phylogenies,
16 with all sequences within these groups being derived from the same host class or
17 order, implying a hitherto underappreciated stability in circovirus-host relationships.
18 The openly available dataset constructed in this investigation provides new insights
19 into circovirus evolution, and can be used to facilitate further studies of circoviruses
20 and CVe.

21

22 **Keywords:** Circovirus; Evolution; Endogenous; Paleovirology; Taxonomy;
23 Phylogeny;

24

25 **Abbreviations:**

26 CVe = endogenous circoviral element

27 ORF = open reading frame

28 Cap = Capsid

29 Rep = Replicase

30 ssDNA = single-stranded DNA

1 **1. Introduction**

2 Circoviruses (family *Circoviridae*, genus *Circovirus*) are small, non-enveloped
3 viruses with single-stranded DNA (ssDNA) genomes. Circovirus genomes are
4 typically ~2 kilobases (kb) in length and contain only two open reading frames
5 (ORFs): one encoding a non-structural replicase (Rep) protein, and a second
6 encoding the viral capsid (Cap). The family contains two genera: *Circovirus* and
7 *Cyclovirus* [1]. Over recent years, sequencing-based virus discovery efforts have
8 identified many new members of these two genera [2]. However, little or nothing is
9 known about most of the novel viruses that have been identified using these
10 approaches. Only a handful of circoviruses have been investigated at a level beyond
11 sequencing: porcine circoviruses 1 and 2 (PCV-1 and PCV-2), which infect swine,
12 and beak and feather disease virus (BFDV), which infects various avian species [3].

13 Circovirus sequences are frequently recovered from tissue and environmental
14 samples in metagenomic investigations, indicating that these viruses are widespread,
15 yet they remain relatively poorly understood [4]. Endogenous circoviruses (CVe)
16 provide an unconventional but useful source of information about circovirus
17 distribution, diversity and evolution. These sequences are derived from the genomes
18 of circoviruses that circulated millions of years ago, and became integrated into the
19 host germline [5, 6]. Relatively robust minimum age estimates can be obtained for
20 CVe via the identification of orthologous copies in distinct host lineages. On this
21 basis, we now know that the association between circoviruses and vertebrates
22 extends back millions of years before the present day [7, 8].

23 In this study, we screened vertebrate genomes *in silico* to generate a
24 comprehensive catalog of CVe. We used these data to: (i) extract information about
25 the long-term evolution of circoviruses; (ii) generate an openly accessible data
26 resource that can facilitate the further investigation of CVe and circoviruses.

27

28 **2. Material & Methods**

29 2.1. Identification and analysis of CVe sequences

30 We used similarity searches to systematically screen genome assemblies of
31 362 chordate species (**Table S1**) for sequences homologous to circovirus proteins.
32 Vertebrate genome assemblies and circovirus reference genomes were obtained
33 from the NCBI genomes resource. Screening *in silico* was performed using the
34 database-integrated genome-screening tool. The DIGS procedure used to identify
35 CVe comprises two steps. In the first, a circovirus probe sequence (e.g. a Cap or
36 Rep protein sequence) is used to search a particular genome assembly file using the
37 basic local alignment search tool (BLAST) program [9]. In the second, sequences

1 that produce statistically significant matches to the probe are extracted and classified
2 by BLAST-based comparison to a set of virus reference genomes (see **Table S2**).
3 Results are captured in a MySQL database.

4 We inferred the ancestral ORFs of C_{Ve} (and the number of stop codons and
5 frameshifts interrupting these ORFs) via a combination of automated alignment and
6 manual adjustment. Multiple sequence alignments were constructed using MUSCLE
7 [10] and PAL2NAL [11]. Manual inspection and adjustment of alignments was
8 performed in Se-AL [12]. Phylogenies were constructed using maximum likelihood as
9 implemented in RaxML [13], and the VT protein substitution model [14] as selected
10 using ProTest [15].

11

12 2.2. Construction of C_{Ve} sequence data resource.

13 We used GLUE - an open, data-centric software environment specialized in
14 capturing and processing virus genome sequence datasets – to collate the
15 sequences, alignments and associated data used in this investigation. The aim was
16 to create a standardized data C_{Ve} resource that would be openly accessible, and
17 would facilitate the further use and development of the dataset assembled here. The
18 project includes all the C_{Ve} identified by our *in silico* screen, as well as a set of
19 representative reference sequences for the *Circovirus* genus (**Table S2**). All of these
20 sequences are linked to the appropriate auxiliary data; for the virus sequences, this
21 includes information about the sample from which the sequence was obtained; for
22 C_{Ve}, it includes the name of genome assembly and contig in which the C_{Ve}
23 sequence was identified, and its coordinates and orientation within that contig.

24 The project also includes the key alignments constructed in this study, linked
25 together using the GLUE ‘alignment tree’ data structure. These include: (i) ‘tip’
26 alignments in which all taxa are C_{Ve} that are known or putative orthologs of one
27 another; (ii) a ‘root’ alignment constructed to represent proposed homologies
28 between the genomes of representative viruses in the genus *Circovirus* and the C_{Ve}
29 recovered by our screen. Because each of these alignments is constrained to a
30 standard reference sequence, alignments are linked to one another.

31 We applied a systematic approach to naming C_{Ve}. Each element was
32 assigned a unique identifier (ID) constructed from a defined set of components. The
33 first component is the classifier ‘C_{Ve}’. The second is a composite of two distinct
34 subcomponents separated by a period: the name of C_{Ve} group (usually derived from
35 the host group in which the element occurs in (e.g. Carnivora), and the second is a
36 numeric ID that uniquely identifies the insertion. Orthologous copies in different
37 species are given the same number, but are differentiated using the third component

1 of the ID that uniquely identifies the species from which the sequence was obtained.
2 An additional unique numeric ID may be added to this component in cases where a
3 Cve element has expanded via duplication.

4

5 **3. Results**

6 3.1. Identification and phylogenetic analysis of vertebrate Cve

7 We systematically screened 362 vertebrate genome assemblies for Cve, and
8 identified a total of 179 Cve sequences (**Table S3**), in 52 distinct species (**Table 2**).
9 For each Cve sequence, we determined the regions of the circovirus genome
10 represented, and attempted to identify genomic flanks. Where genomic flanks were
11 present, we compared these with one another to identify potentially orthologous Cve
12 loci. In several cases, it was not possible to determine whether multiple Cve loci
13 within the same species (or group of closely-related species) represented the
14 outcome of distinct incorporation events, or the fragmented remains of a single,
15 ancestrally acquired element. The main causes of uncertainty were; (i) lack of
16 flanking sequences due to short contig length, or undetermined DNA sequences
17 flanking Cve, and; (ii) the presence of multiple Cve that spanned non-overlapping
18 regions of the circovirus genome. Since Cve are comparatively rare in vertebrate
19 genomes [7, 16], we conservatively assumed a single incorporation event had taken
20 place except in cases where it could be demonstrated otherwise. On this basis, we
21 estimate that the 179 Cve identified here represent at least 19-26 distinct germline
22 incorporation events (**Table 1, Table 2, Figure 1**), depending on whether Cve in ray-
23 finned fish are taken to represent a single incorporation event, or seven distinct
24 incorporation events, each in a different order (see **section 3.3**). The large
25 discrepancy between the number of elements versus the number of events reflects
26 the fact that 101 of the 179 Cve identified in our study (57%) belong to a group of
27 highly duplicated Cve loci in carnivore genomes, all of which derive from a single
28 germline incorporation event.

29 We only identified four cases where Cve encoding both *rep* and *cap* were
30 present in the same species or species group. In most, only *rep*-derived sequences
31 appear to have been incorporated/retained, and in one case only *cap* (**Table 1**). We
32 constructed a multiple sequence alignment (MSA) that spanned the entire circovirus
33 genome and contained both reference sequences for Cve (these could be based on
34 individual loci, or a consensus), and representative circovirus reference taxa (**Table**
35 **S2**). We used this 'root' MSA (see **section 2.2**) to infer which regions of the
36 circovirus genome had been incorporated as Cve. Where Cve spanned coding
37 sequence, we inferred the putative ancestral reading frame by comparing Cve and

1 circovirus sequences, and attempting to identify likely frameshifting mutations. Most
2 Cve represent only fragments of the genome (**Figure 1**), and many are relatively
3 degraded, containing multiple frameshifting indels and stop codons.

4 Where we identified several Cve from the same species, we compared
5 genomic regions to search for evidence of homology and thereby identify orthologs.
6 Where we were able to identify orthologous Cve insertions, we used these data to
7 create a timeline of circovirus evolution (**Figure 2**). In addition, we identified sets of
8 ‘potentially orthologous’ Cve, where sequence similarity and phylogenetic
9 relationships were consistent with orthology, but this could not be confirmed or ruled
10 out based on flanking sequences.

11 A range of distinct partitions were derived from the virtually translated root
12 MSA (with frameshifts removed), and used to construct bootstrapped ML phylogenies
13 (**Figure 3**). In general, support for the deeper branching relationships between Cve
14 and circoviruses was weak, irrespective of which genomic region was used to
15 construct trees. This reflects the fact that most Cve are short and/or highly degraded,
16 and these sequences tend to group distantly from other taxa. However, in
17 phylogenies based on Rep (**Figure 3**), several robustly supported subgroupings were
18 observed, three of which – referred to here as mammal 1, cyprinid 1, and avian 1 -
19 included a mixture of Cve and contemporary circoviruses. Notably, in all three of
20 these clades, Cve and circovirus sequences were obtained from the same hosts of
21 the same taxonomic class. The sections that follow describe the distribution and
22 diversity of Cve within individual classes in the subphylum Vertebrata (chordates with
23 backbones).

24

25 3.2. Cve in jawless vertebrates

26 Extant vertebrates are divided into the jawed vertebrates (Gnathostomata)
27 and jawless vertebrates (Agnatha). The Agnatha represent the most basal group of
28 vertebrates and includes the hagfishes (myxinooids) and lampreys (petromyzontids).
29 We identified seven sequences exhibiting homology to *rep* in the genome assembly
30 of the inshore hagfish (*Eptatretus burgeri*). These sequences are relatively distinct
31 from other circoviruses, and also showed relatively high genetic diversity relative to
32 one another, forming three distinct groups in phylogenetic trees (**Figure S1**). Notably,
33 the putative Rep polypeptides encoded by these sequences contained several in-
34 frame indels relative to one another. Because such a pattern of variation is unlikely to
35 arise through neutral accumulation of mutations in the germline, this suggests the
36 occurrence of at least three distinct genome incorporation events, each involving
37 distinct, but relatively closely related viruses. However, since we were unable to

1 identify unambiguous genomic flanking sequences for any of these loci, their
2 classification as C_{Ve} should for now be considered tentative.

3

4 3.3. C_{Ve} in ray-finned fish (class Actinopterygii)

5 Circoviruses are thought to infect barbel fish (*Barbus barbus*) and European
6 catfish (*Silurus glanis*), based on (i) the observation of viral particles in tissues, and
7 the recovery of circovirus sequences from these tissues via nested PCR [17, 18]. In
8 addition, C_{Ve} have been reported in one fish species - the Indian rohu (*Labeo rohita*)
9 [19]. We identified numerous additional C_{Ve} sequence in the genome assemblies of
10 ray-finned fishes (Class Actinopterygii) (**Table 2, Table S3**). We established that at
11 least two of these C_{Ve} - occurring in the common carp (*Cyprinus carpio*) and golden-
12 line barbell (*Sinocyclocheilus grahami*) genomes - were orthologs of one another,
13 indicating they were incorporated into the germline of cyprinid fish more than 39
14 million years ago [28, 29]. These C_{Ve} were comprised of multiple complete circovirus
15 genomes arranged in tandem, and intriguingly, were observed C_{Ve} group as sister
16 taxa to barbel circovirus (BarbCV) in phylogenetic trees, sharing ~70% nucleotide
17 identity (across 1654 nucleotides) with the BarbCV genome.

18 We also identified matches to *rep* in eight other species of ray-finned fish
19 (**Table 2**). We could not determine with certainty how many integration events these
20 C_{Ve} represented. Interestingly, however, all of these sequences group together in
21 phylogenies (**Figure 3**), and the phylogeny constructed for these elements - when
22 rooted on the C_{Ve} from the most basal host - the European eel (*Anguilla anguilla*),
23 approximately follows that of the host species, consistent with a single ancestral
24 integration event >200 million years ago (**Figure 2**). Alternatively, the C_{Ve} observed
25 in distinct orders might represent distinct incorporation events. This is supported by
26 the placement of C_{Ve}.anura in phylogenies, in which it splits the fish C_{Ve} from one
27 another, albeit with weak support (**Figure S1**). In addition, the observation that C_{Ve}
28 elements in order cypriniforme fish (golden-line barbell and carp) occur as full-length
29 tandem genomes, whereas those in Perciformes are derived from more divergent
30 fragments of *rep*, is suggestive of at least two separate incorporation events. Notably
31 one C_{Ve} in the mangrove rivulus (*Kryptolebias marmoratus*) encoded a complete
32 intact *rep* gene (**Figure 1**) that is predicted to be expressed, suggesting it may have
33 been functionalized in some manner.

34

35 3.4. C_{Ve} in amphibians

36 Sequences homologous to circovirus *rep* genes have previously been
37 identified in the Western clawed frog (*Xenopus tropicalis*) [20]. We identified C_{Ve} in

1 the genome of the American bullfrog (*Rana catesbeiana*) that partially overlaps that
2 identified in *Xenopus*. Potentially, these sequences could be orthologs of one
3 another, which would imply a minimum age of ~204 MYA [21, 22] (**Figure 2**).
4 However, we were unable to confirm this based on analysis of flanking genomic
5 sequences.

6

7 3.5. C_{Ve} in reptiles

8 A pair of orthologous C_{Ve}, each covering about 75% of the circovirus genome,
9 have previously been recovered from rattlesnake genomes (*Crotalus spp*) [23]. We
10 identified C_{Ve} in four additional snake species (**Table 2**). Examination of aligned
11 snake C_{Ve} sequences indicated that all are likely to be orthologs of those previously
12 reported in rattlesnakes (see **Figure S2**), implying that this C_{Ve} integrated into the
13 serpentine germline ~72-90 million years ago (Mya) (**Figure 2**).

14

15 3.6. C_{Ve} in birds (class Aves)

16 C_{Ve} have previously been reported in the genomes of several avian species:
17 the little egret (*Egretta garzetta*), white-throated tinamou (*Tinamus guttatus*), medium
18 ground-finch (*Geospiza fortis*), and kea (*Nestor notabilis*) [16, 20]. We identified C_{Ve}
19 in eight additional species. Some of these appeared likely to be orthologs of C_{Ve}
20 reported previously. For example, we identified C_{Ve} in two species of psittacine bird
21 that appeared represented orthologs of one another, and possibly of those previously
22 identified in the kea (*Nestor notabilis*) [16] (**Table 2**), which would imply integration
23 into the psittacine germline prior to the divergence of the major extant lineages within
24 the order Psittaciformes (estimated to have occurred 30-60 Mya [13, 24]) (**Figure 2**).

25 We also identified orthologs of the *rep*-derived insertion previously described
26 in the medium ground finch in several additional species in the avian order
27 Passeriformes (songbirds) (**Table 2**). Identification of these orthologs demonstrates
28 that this particular C_{Ve} predates the radiation of avian sub-order Passeroidea ~38
29 Mya [13, 25] (**Figure 2**).

30 In addition to identifying the previously reported C_{Ve} in the genomes of the
31 white-throated tinamou (*Tinamus guttatus*) and little egret (*Egretta garzetta*) [16], we
32 identified previously unreported C_{Ve} in the Japanese rail (*Gallirallus okinawae*: order
33 Gruiformes) and downy woodpecker (*Picoides pubescens*: order Piciformes) (**Table**
34 **2**). Both these sequences were relatively short and divergent, and consequently we
35 could not determine their relationships to other C_{Ve} and circoviruses with confidence.

36

37 3.6. C_{Ve} in mammals (class Mammalia)

1 The majority of C_{Ve} identified in our screen were recovered from carnivore
2 genome assemblies. As far as we are able to discern from phylogenetic and
3 comparative analysis, all of these C_{Ve} derive from a 1-4 germline incorporation
4 events involving an ancient carnivore *rep* gene. However, the copy number of these
5 elements has expanded subsequent to their incorporation into the germline, in some
6 cases quite dramatically. The grouping of carnivore C_{Ve} in phylogenies (**Figure 4**)
7 indicates that at least four C_{Ve} insertions were present in the carnivore germline prior
8 to the divergence of extant families within this order. The copy number of one
9 particular element (referred to here as C_{Ve}-Carnivora-4) has expanded in some
10 carnivore lineages. As shown in **Figure 4**, the phylogenetic relationships between
11 duplicates in the group C_{Ve}-Carnivora-4 indicate that these expansions have
12 occurred independently in ursids (bears), pinnipeds (seals and walruses), and
13 mustelids. C_{Ve} in this lineage are flanked by sequences that disclose homology to
14 non-LTR retrotransposons. Thus, one plausible explanation for the elevated copy
15 number in certain carnivore lineages is that C_{Ve} have become embedded into
16 retroelements and copied along with these sequences when they undergo
17 transposition.

18 A novel, relatively well-preserved *rep*-derived C_{Ve} was identified in the
19 genome of the Ryukyu mouse (*Mus caroli*) that grouped closely with circoviruses
20 genome recovered from dogs [26, 27]. This element presumably arose after this
21 species diverged from the house mouse (*Mus musculus*) ~6-7 Mya, since it is absent
22 from this species.

23 In the cape golden mole (*Chrysochloris asiaticus*) matches to both *cap* and
24 *rep* were identified. However, these occurred on distinct contigs and did not overlap.
25 Furthermore, both C_{Ve} were relatively short and degraded, and were highly
26 divergent relative to other C_{Ve}. C_{Ve} derived from *cap* were also identified in the
27 genome of Hoffmann's two-toed sloth (*Choloepus hoffmanni*) (**Figure 1**).

28 C_{Ve} have previously been identified in the genome of the short-tailed
29 opossum (*Monodelphis domestica*), an American marsupial [7]. In phylogenies based
30 on *rep*, this sequence groups together with the porcine circoviruses, canine
31 circovirus, and the C_{Ve} we identified in *Mus caroli*. We identified the first examples of
32 C_{Ve} from the genomes of Australian marsupial species: the Tasmanian devil
33 (*Sarcophilus harrisii*) and the koala (*Phascolarctos cinereus*). Both these sequences
34 derived from circovirus *rep* genes, and grouped together in phylogenetic trees
35 (**Figure S1**). However, their placement relative to other taxa was not supported with
36 confidence, reflecting their short and degraded nature. Several other short and
37 degraded matches to Rep probes were identified in other mammalian species (**Table**

1 **1, Table 2, Figure 1**). These sequences were relatively distantly related to one
2 another and to contemporary circoviruses.

3

4 **4. Discussion**

5 4.1. CVe provide retrospective information about circovirus evolution.

6 In this study, we recovered CVe from published vertebrate genomes,
7 determined their genomic structures, and examined their phylogenetic relationships
8 to contemporary circoviruses. Our analysis is the first to examine such a large set of
9 CVe sequences, and to screen so widely within vertebrates. We show that CVe are
10 relatively widespread in vertebrate genomes, though it appears they are absent from
11 some lineages (e.g. primates, in which genome coverage is relatively high).

12 Several of the CVe loci identified here have been reported previously [7, 16,
13 19, 20], and the majority of novel CVe sequences recovered by our screen were
14 orthologs or duplicates of these loci. Nevertheless, we identified 17 CVe loci that
15 have not been reported before (**Table 1, Table S3**). These sequences provide the
16 first evidence of (ancestral) circovirus infection for several species (**Table 2**). In
17 addition, the identification and characterisation of novel orthologs allowed us to
18 establish the first minimum age estimates for some CVe loci, and to markedly
19 extended those of others. Thus, we were able to derive a more accurately calibrated
20 timeline of evolution for the *Circovirus* genus, spanning multiple geological eras
21 (**Figure 2**). Furthermore, we observed that CVe in fish, birds and mammals cluster
22 phylogenetically with exogenous circoviruses identified from the same host class.
23 This implies that there is a stability to the relationship between circovirus and host
24 relationships, at least at higher taxonomic levels.

25

26 4.1. Impact of CVe on host genome evolution

27 The majority of CVe are derived from *rep* genes. To the extent that CVe have
28 been exapted or co-opted, the predominance of CVe derived from *rep* might reflect
29 that these sequences are more readily functionalised than those derived from *cap*.
30 Furthermore, we identified one sequence. Notably one CVe in the mangrove rivulus
31 (*Kryptolebias marmoratus*) encoded an intact *rep* gene that is predicted to express
32 mRNA, suggesting it may have been functionalized in some manner (**Figure 1**).

33 Notably, several examples have now been described of endogenous viral
34 elements (EVEs) that are derived from replicase genes, are expressed, and encode
35 intact ORFs [7, 30, 31]. These elements are derived from a range of different viruses,
36 and have clearly arisen in distinct events, suggesting there might be a common
37 mechanism causing EVEs derived from the replicases of distinct viruses to be

1 selected and maintained in different species. Alternatively, it is possible that the
2 discrepancy in numbers simply reflects that *cap*-derived sequences are less
3 conserved and therefore harder to detect.

4 Curiously, it is rare for more than one C_{Ve} to occur in the germline of any
5 jawed vertebrate lineage. Carnivores are an obvious exception, since C_{Ve} have been
6 amplified to relatively high copy number (10-20 copies) in several carnivore lineages
7 (**Figure 4**), apparently via retrotransposon-mediated duplication. Further investigation
8 of how these C_{Ve} have been amplified may reveal if their presence within an actively
9 replicating retrotransposon lineage has impacted on the fixation of transposable
10 elements derived from that lineage.

11

12 **Conclusions**

13 We identified the complete repertoire of C_{Ve} sequences in published
14 vertebrate genome assemblies. Through comparative analysis of these sequences,
15 we provide the most complete picture yet of how viruses in the genus *Circovirus* have
16 evolved and interacted with their hosts over the course of their evolution. The
17 sequence-based resource implemented here can facilitate further characterisation of
18 circovirus distribution, diversity and evolution as new C_{Ve} and circovirus sequence
19 data become available.

1 Acknowledgements

2 RJG was funded by the Medical Research Council of the United Kingdom
3 (MC_UU_12014/12). WMS is supported by the Fundação de Amparo à Pesquisa do
4 Estado de São Paulo, Brazil (Scholarships NO. 17/13981-0).

6 References

- 8 1. Li, L., et al., *Multiple diverse circoviruses infect farm animals and are*
9 *commonly found in human and chimpanzee feces.* J Virol, 2010. **84**(4): p.
10 1674-82.
- 11 2. Delwart, E. and L. Li, *Rapidly expanding genetic diversity and host range of*
12 *the Circoviridae viral family and other Rep encoding small circular ssDNA*
13 *genomes.* Virus Res, 2012. **164**(1-2): p. 114-21.
- 14 3. Amery-Gale, J., et al., *A high prevalence of beak and feather disease virus in*
15 *non-psittacine Australian birds.* J Med Microbiol, 2017. **66**(7): p. 1005-1013.
- 16 4. Shulman, L.M. and I. Davidson, *Viruses with Circular Single-Stranded DNA*
17 *Genomes Are Everywhere!* Annu Rev Virol, 2017. **4**(1): p. 159-180.
- 18 5. Holmes, E.C., *The evolution of endogenous viral elements.* Cell Host
19 Microbe, 2011. **10**(4): p. 368-77.
- 20 6. Feschotte, C. and C. Gilbert, *Endogenous viruses: insights into viral evolution*
21 *and impact on host biology.* Nat Rev Genet, 2012. **13**(4): p. 283-96.
- 22 7. Katzourakis, A. and R.J. Gifford, *Endogenous viral elements in animal*
23 *genomes.* PLoS Genet, 2010. **6**(11): p. e1001191.
- 24 8. Belyi, V.A., A.J. Levine, and A.M. Skalka, *Sequences from ancestral single-*
25 *stranded DNA viruses in vertebrate genomes: the parvoviridae and*
26 *circoviridae are more than 40 to 50 million years old.* J Virol, 2010. **84**(23): p.
27 12458-62.
- 28 9. Altschul, S.F., et al., *Gapped BLAST and PSI-BLAST: a new generation of*
29 *protein database search programs.* Nuc. Acids Res., 1997. **25**: p. 3389-3402.
- 30 10. Edgar, R.C., *MUSCLE: multiple sequence alignment with high accuracy and*
31 *high throughput.* Nucleic Acids Res, 2004. **32**(5): p. 1792-7.
- 32 11. Suyama, M., D. Torrents, and P. Bork, *PAL2NAL: robust conversion of*
33 *protein sequence alignments into the corresponding codon alignments.*
34 Nucleic Acids Res, 2006. **34**(Web Server issue): p. W609-12.
- 35 12. Rambaut, A., *SE-AL Sequence Alignment Editor.* 2002, University of Oxford:
36 Oxford, UK.
- 37 13. Jarvis, E.D., et al., *Whole-genome analyses resolve early branches in the tree*
38 *of life of modern birds.* Science, 2014. **346**(6215): p. 1320-31.
- 39 14. Muller, T. and M. Vingron, *Modeling amino acid replacement.* J Comput Biol,
40 2000. **7**(6): p. 761-76.
- 41 15. Darriba, D., et al., *ProtTest 3: fast selection of best-fit models of protein*
42 *evolution.* Bioinformatics, 2011. **27**(8): p. 1164-5.
- 43 16. Cui, J., et al., *Low frequency of paleoviral infiltration across the avian*
44 *phylogeny.* Genome Biol, 2014. **15**(12): p. 539.
- 45 17. Lorincz, M., et al., *First detection and analysis of a fish circovirus.* J Gen Virol,
46 2011. **92**(Pt 8): p. 1817-21.
- 47 18. Lorincz, M., et al., *Novel circovirus in European catfish (Silurus glanis).* Arch
48 Virol, 2012. **157**(6): p. 1173-6.
- 49 19. Feher, E., et al., *Integrated circoviral rep-like sequences in the genome of*
50 *cyprinid fish.* Virus Genes, 2013. **47**(2): p. 374-7.
- 51 20. Liu, H., et al., *Widespread horizontal gene transfer from circular single-*
52 *stranded DNA viruses to eukaryotic genomes.* BMC Evol Biol, 2011. **11**: p.
53 276.

- 1 21. Cannatella, D., *Xenopus in Space and Time: Fossils, Node Calibrations, Tip-*
2 *Dating, and Paleobiogeography*. Cytogenet Genome Res, 2015. **145**(3-4): p.
3 283-301.
- 4 22. Roelants, K., A. Haas, and F. Bossuyt, *Anuran radiations and the evolution of*
5 *tadpole morphospace*. Proc Natl Acad Sci U S A, 2011. **108**(21): p. 8731-6.
- 6 23. Gilbert, C., et al., *Endogenous hepadnaviruses, bornaviruses and*
7 *circoviruses in snakes*. Proc Biol Sci, 2014. **281**(1791): p. 20141122.
- 8 24. Jetz, W., et al., *The global diversity of birds in space and time*. Nature, 2012.
9 **491**(7424): p. 444-8.
- 10 25. Prum, R.O., et al., *A comprehensive phylogeny of birds (Aves) using targeted*
11 *next-generation DNA sequencing*. Nature, 2015. **526**(7574): p. 569-73.
- 12 26. Li, L., et al., *Circovirus in tissues of dogs with vasculitis and hemorrhage*.
13 *Emerg Infect Dis*, 2013. **19**(4): p. 534-41.
- 14 27. Decaro, N., et al., *Genomic characterization of a circovirus associated with*
15 *fatal hemorrhagic enteritis in dog, Italy*. PLoS One, 2014. **9**(8): p. e105909.
- 16 28. Wang, X., et al., *Cyprinid phylogeny based on Bayesian and maximum*
17 *likelihood analyses of partitioned data: implications for Cyprinidae*
18 *systematics*. Sci China Life Sci, 2012. **55**(9): p. 761-73.
- 19 29. Ren, Q. and R.L. Mayden, *Molecular phylogeny and biogeography of African*
20 *diploid barbs, 'Barbus', and allies in Africa and Asia (Teleostei:*
21 *Cypriniformes)*. Zoologica Scripta, 2016. **45**: p. 642–649.
- 22 30. Arriagada, G. and R.J. Gifford, *Parvovirus-derived endogenous viral elements*
23 *in two South American rodent genomes*. J Virol, 2014. **88**(20): p. 12158-62.
- 24 31. Horie, M., et al., *An RNA-dependent RNA polymerase gene in bat genomes*
25 *derived from an ancient negative-strand RNA virus*. Sci Rep, 2016. **6**: p.
26 25873.
27
28

1 **Table 1. CVe detected in published vertebrate genome assemblies**

EVE name	Reference	Genes	# Seqs	# Species
Agnatha				
CVe- <i>Eptatretus</i> *	This study	Rep	7	1
Bony Fish				
CVe- <i>Anguilla</i> *	This study	Rep		1
CVe- <i>Characiformes</i> *	This study	Rep		1
CVe- <i>Clupeiformes</i> *	This study	Rep		1
CVe- <i>Cypriniformes</i>	[19]	Rep-Cap		2
CVe- <i>Cyprinodontiformes</i> *	This study	Rep		1
CVe- <i>Perciformes</i> *	This study	Rep		3
CVe- <i>Salmoniformes</i> *	This study	Rep		1
Amphibians				
CVe- <i>Anura</i>	[16]	Rep	2	2
Reptiles				
CVe- <i>Viperidae</i>	[23]	Rep-Cap	16	6
Birds				
CVe- <i>Tinamou</i>	[16]	Rep-Cap	2	1
CVe- <i>Psittaciformes</i>	[16]	Rep-Cap	4	3
CVe- <i>Passeriformes</i>	[16]	Rep	7	5
CVe- <i>Egretta</i>	[16]	Rep	1	1
CVe- <i>Gallirallus</i> *	This study	Rep	1	1
CVe- <i>Picoides</i> *	This study	Rep	1	1
Mammals				
CVe- <i>Chrysochloris</i> *	This study	Rep, Cap	3	1
CVe- <i>Carnivora</i>	[7]	Rep	101	13
CVe- <i>Mus.caroli</i> *	This study	Rep	1	1
CVe- <i>Heterocephalus</i> *	This study	Rep	1	1
CVe- <i>Phascolarctos</i> *	This study	Rep	1	1
CVe- <i>Sarcophilus</i> *	This study	Rep	1	1
CVe- <i>Monodelphis</i>		Rep	1	1
CVe- <i>Galeopterus</i> *	This study	Rep	2	1
CVe- <i>Manis</i> *	This study	Rep	1	1
CVe- <i>Choloepus</i> *	This study	Cap	2	1
Totals			179	53

2
3
4

Footnote: Asterisks indicate newly identified CVe loci.

1 **Table 2. Vertebrate species with CVE**
2

Latin binomial	Common name	EVE name	1st. ^a	Copies ^b
Agnatha				
<i>Eptatretus burgeri</i>	Inshore hagfish	CVE- <i>Eptatretus</i> *	x	7
Bony Fish				
<i>Anguilla anguilla</i>	European eel	CVE- <i>Anguilla</i>	x	1
<i>Pygocentrus nattereri</i>	Red-bellied piranha	CVE- <i>Characiformes</i>		1
<i>Clupea harengus</i>	Atlantic herring	CVE- <i>Clupeiformes</i>		1
<i>Cyprinus carpio</i>	Common carp	CVE- <i>Cypriniformes</i>		4
<i>Sinocyclocheilus grahami</i>	Golden-line barbel	CVE- <i>Cypriniformes</i>		2
<i>Kryptolebias brasiliensis</i>	Killifish	CVE- <i>Cyprinodontiformes</i>	X	4
<i>Micropterus floridanus</i>	American black bass	CVE- <i>Perciformes</i>	X	1
<i>Neolamprologus brichardi</i>	Princess of Burundi	CVE- <i>Perciformes</i>	X	2
<i>Acanthochromis polyacanthus</i>	Spiny chromis damselfish	CVE- <i>Perciformes</i>	X	5
<i>Salmo salar</i>	Atlantic salmon	CVE- <i>Salmoniformes</i>	X	3
Amphibians				
<i>Xenopus tropicalis</i>	Western clawed frog	CVE- <i>Xenopus</i>		1
<i>Rana catesbeiana</i>	American bullfrog	CVE- <i>Rana</i>	x	1
Reptiles				
<i>Pantherophis guttatus</i> *	Corn snake	CVE- <i>Viperidae</i>		1
<i>Python molurus</i> *	Indian python	CVE- <i>Viperidae</i>		1
<i>Crotalus horridus</i>	Timber rattlesnake	CVE- <i>Viperidae</i>		1
<i>Crotalus mitchellii pyrrhus</i>	Mitchell's rattlesnake	CVE- <i>Viperidae</i>		1
<i>Protobothrops mucrosquamatus</i> *	Brown spotted pit viper	CVE- <i>Viperidae</i>		1
<i>Ophiophagus hannah</i> **	King cobra	CVE- <i>Viperidae</i>		1
Birds				
<i>Serinus canaria</i> *	Atlantic canary	CVE- <i>Passeriformes</i>		1
<i>Setophaga coronata</i> *	Yellow-rumped warbler	CVE- <i>Passeriformes</i>		1
<i>Sporophila hypoxantha</i> *	Lined seedeater	CVE- <i>Passeriformes</i>		1
<i>Zonotrichia albicollis</i> *	White-throated sparrow	CVE- <i>Passeriformes</i>		1
<i>Geospiza fortis</i>	Medium ground finch	CVE- <i>Passeriformes</i>		1
<i>Agapornis roseicollis</i> *	Rosy-faced lovebird	CVE- <i>Psittaciformes</i>		1
<i>Amazona aestiva</i> *	Turquoise-fronted amazon	CVE- <i>Psittaciformes</i>		2
<i>Nestor notabilis</i>	Kea	CVE- <i>Psittaciformes</i>		1
<i>Tinamus guttatus</i>	White-throated tinamou	CVE- <i>Tinamou</i>		1
<i>Egretta garzetta</i>	Little egret	CVE- <i>Egretta</i>		1
<i>Gallirallus okinawae</i> *	Okinawa rail	CVE- <i>Gallirallus</i>	x	1
<i>Picoides pubescens</i> *	Downy woodpecker	CVE- <i>Picoides</i>	x	1
Mammals				
<i>Ailurus fulgens</i> *	Red panda	CVE- <i>Carnivora</i>		4
<i>Canis familiaris</i>	Domestic dog	CVE- <i>Carnivora</i>		4
<i>Lycaon pictus</i> *	Cape hunting dog	CVE- <i>Carnivora</i>		3
<i>Acinonyx jubatus</i> *	Cheetah	CVE- <i>Carnivora</i>		1
<i>Felis catus</i>	Domestic cat	CVE- <i>Carnivora</i>		3
<i>Panthera tigris altaica</i> *	Siberian tiger	CVE- <i>Carnivora</i>		1
<i>Enhydra lutris</i> *	Sea otter	CVE- <i>Carnivora</i>		15
<i>Mustela putorius furo</i>	Ferret	CVE- <i>Carnivora</i>		32
<i>Odobenus rosmarus</i> *	Walrus	CVE- <i>Carnivora</i>		15
<i>Leptonychotes weddellii</i> *	Weddell seal	CVE- <i>Carnivora</i>		4
<i>Neomonachus schauinslandi</i> *	Monk seal	CVE- <i>Carnivora</i>		11
<i>Ailuropoda melanoleuca</i>	Panda	CVE- <i>Carnivora</i>		7
<i>Ursus maritimus</i> *	Polar bear	CVE- <i>Carnivora</i>		8
<i>Heterocephalus glaber</i> *	Naked mole rat	CVE- <i>Heterocephalus</i>	x	1
<i>Mus caroli</i> *	Ryuku mouse	CVE- <i>Mus</i>	x	1
<i>Manis pentadactyla</i> *	Chinese pangolin	CVE- <i>Manis</i>	x	1
<i>Monodelphis domestica</i>	Opossum	CVE- <i>Monodelphis</i>		1
<i>Sarcophilus harrisii</i> *	Tasmanian devil	CVE- <i>Sarcophilus</i>	x	1
<i>Phascogaleos cinereus</i> *	Koala	CVE- <i>Phascogaleos</i>	x	1
<i>Choloepus hoffmanni</i>	Hoffmann's two-toed sloth	CVE- <i>Choloepus</i>		2
<i>Galeopterus variegatus</i> *	Sunda flying lemur	CVE- <i>Galeopterus</i>	x	2
<i>Chrysochloris asiatica</i> *	Cape golden mole	CVE- <i>Chrysochloris</i>	x	3

3 **Footnote:** Asterisks indicate newly identified circovirus EVEs.^a CVE that provide the first evidence of
4 circovirus infection in the host order in which they occur are marked with an 'x'.^b Tandem repeated
5 elements were considered to

1 **Figure Legends**

2

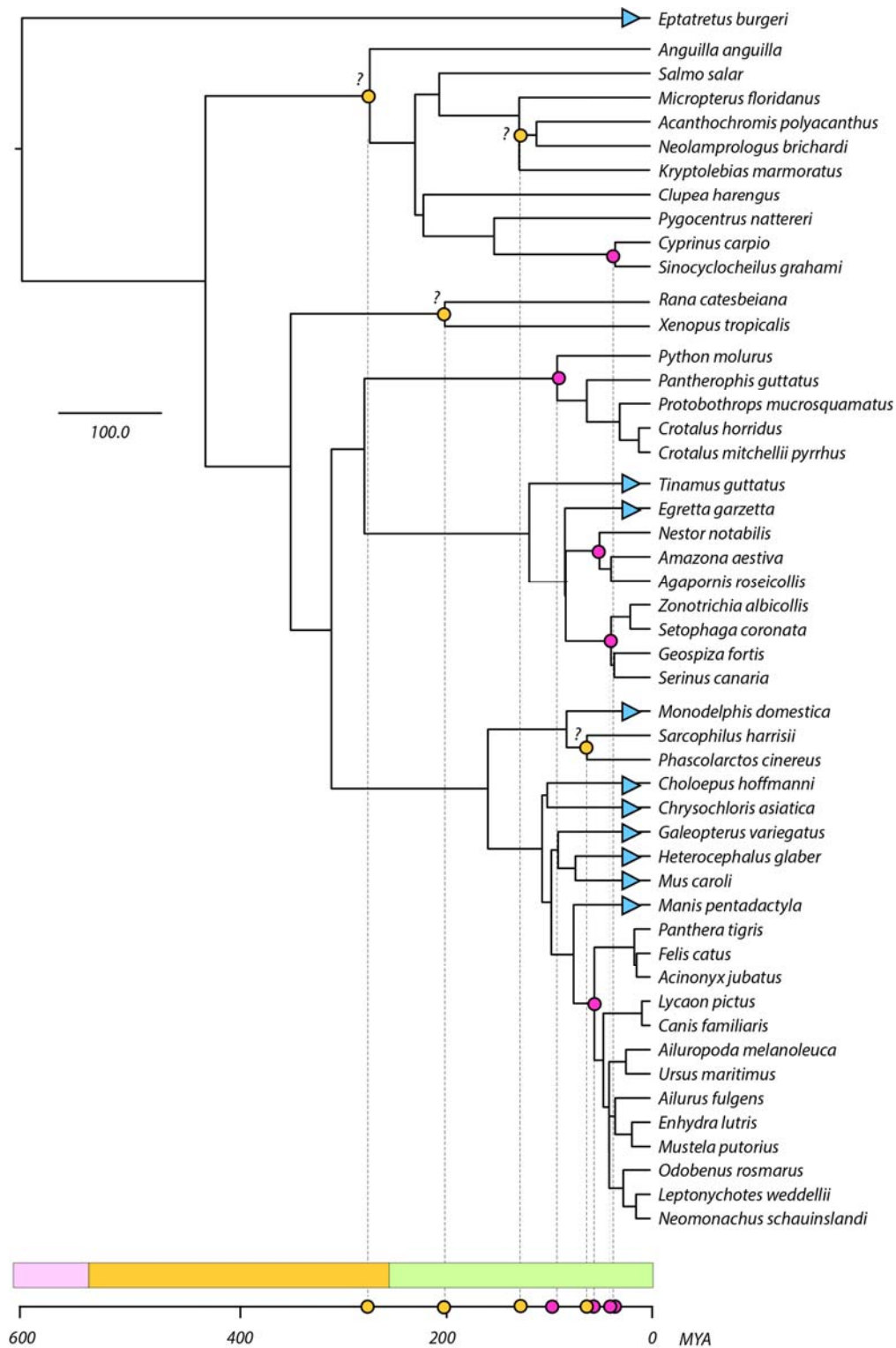
3



4

5

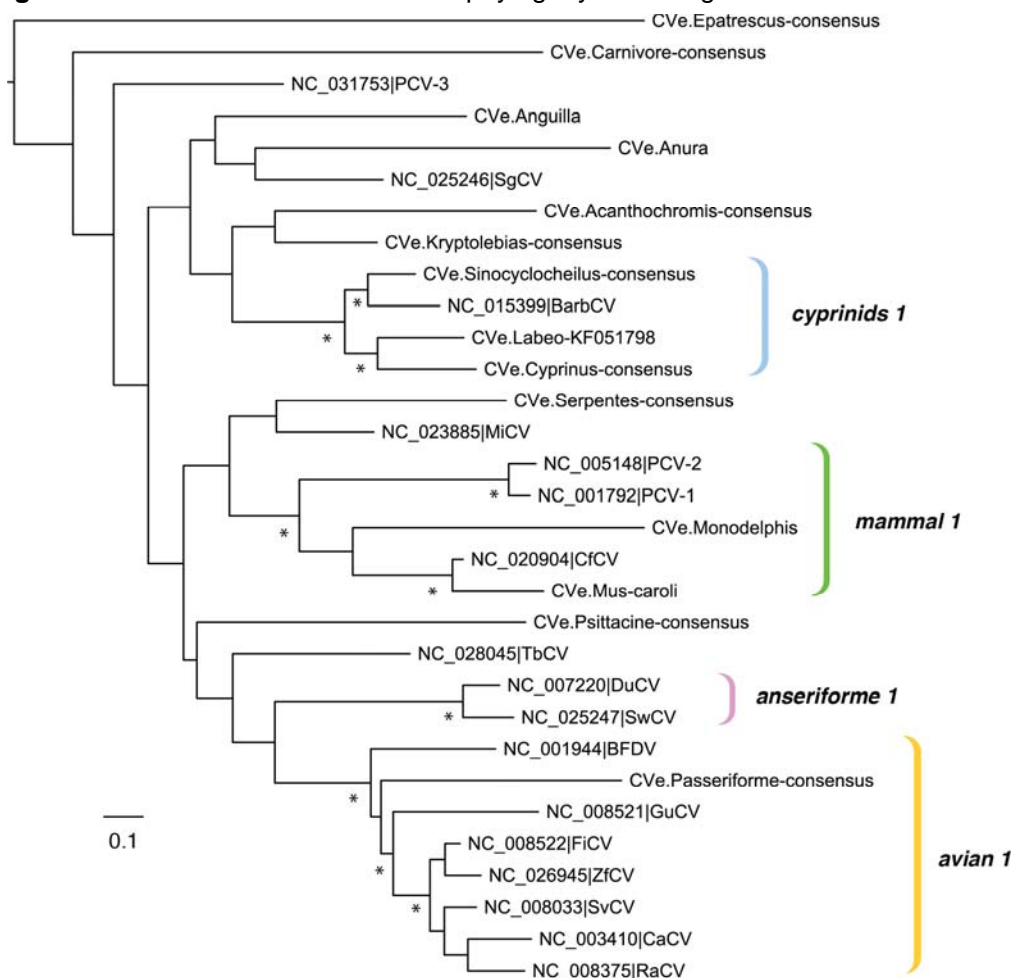
6 **Figure 1.** Genome structures of 21 endogenous circovirus (CVe) elements identified
7 in vertebrate genomes. CVe coding sequences are represented schematically as
8 green and yellow bars relative to a porcine circovirus 1 (PCV1) reference genome
9 (accession # NC_001792.2).



1
2 **Figure 2.** Evolutionary relationships of vertebrate species in which Cve have been
3 identified, and timeline of Cve evolution. Pink circles indicate confirmed orthologs.
4 Yellow circles indicate the presence of potential orthologs that have not been
5 confirmed. Blue triangles indicate where Cve loci are present, but no information
6 about their ages could be obtained.

1

2 **Figure 3.** A maximum likelihood phylogeny showing estimated evolutionary



3 relationships between endogenous circovirus (CVe) elements and exogenous
4 circoviruses. The phylogeny constructed from an alignment spanning ~200 amino
5 acids in Rep. The scale bar shows evolutionary distance in substitutions per site.



1
2
3
4
5
6
7
8
9

Figure 4. Phylogeny of CVe sequences recovered from carnivore genome assemblies. At least four distinct CVe loci are present in the carnivore germline (clades I-IV) as indicated by the coloured brackets. Within group IV, distinct copy number expansions appear to have occurred in ursids (bears), pinnipeds (seals and walrus), and mustelids. The scale bar shows evolutionary distance in substitutions per site. The tree is midpoint rooted for display purposes.