

1 **ABSTRACT**

2 Chapparoviruses are a highly divergent group of parvoviruses (family
3 *Parvoviridae*) first identified in 2013. Interest in these poorly characterized
4 viruses has been raised by recent studies indicating that they are the cause of
5 chronic kidney disease that arises spontaneously in laboratory mice. In this
6 study, we investigate the biological and evolutionary characteristics of
7 chapparoviruses via comparative analysis of genome sequence data. Our
8 analysis, which incorporates sequences derived from endogenous viral
9 elements (EVEs) as well as exogenous viruses, reveals that
10 chapparoviruses are an ancient lineage within the family *Parvoviridae*,
11 clustering separately from members of both currently established parvoviral
12 subfamilies. Consistent with this, they exhibit a number of characteristic
13 genomic and structural features, i.e. a large number of putative auxiliary
14 protein-encoding genes, capsid protein genes non-homologous to any hitherto
15 parvoviral *cap*, as well as a putative capsid structure lacking the canonical fifth
16 strand of the ABIDG sheet comprising the luminal side of the jelly roll. Our
17 findings demonstrate that the chapparovirus lineage infects an exceptionally
18 broad range of host species, including both vertebrates and invertebrates.
19 Furthermore, we observe that chapparoviruses found in fish are more closely
20 related to those from invertebrates than they are to those that infect amniote
21 vertebrates. This suggests that transmission between distantly related host
22 species may have occurred in the past. Our study provides the first integrated
23 overview of the chapparovirus group, and revises current views of parvovirus
24 evolution

25

26 **AUTHOR SUMMARY**

27 Chapparoviruses are a recently identified group of viruses about
28 which relatively little is known. However, recent studies have shown that these
29 viruses cause disease in laboratory mice and are prevalent in the fecal virome
30 of pigs and poultry, raising interest in their potential impact as pathogens, and
31 utility as experimental tools. We examined the genomes of chapparoviruses
32 and endogenous viral elements ('fossilized' virus sequences derived from
33 ancestral viruses) using a variety of bioinformatics-based approaches. We

1 show that the chapparvoviruses have an ancient origin and are evolutionarily
2 distinct from all other related viruses. Accordingly, their genomes and virions
3 exhibit a range of distinct characteristic features. We examine the distribution
4 of these features in the light of chapparvovirus evolutionary history (which we
5 can also infer from genomic data), revealing new insights into chapparvovirus
6 biology.

7

8 INTRODUCTION

9 Parvoviruses (family *Parvoviridae*) are small, non-enveloped viruses
10 with T=1 icosahedral symmetry and linear, single-stranded DNA (ssDNA)
11 genomes ~4-6 kilobases (kb) in length. The family has historically been
12 divided into two subfamilies; *Parvovirinae* and *Densovirinae*, containing
13 viruses that infect vertebrate and invertebrate hosts, respectively (1). Despite
14 exhibiting great variation in expression and transcription strategies, they have
15 a relatively conserved overall genome structure: a non-structural (NS)
16 expression cassette is located at the left side of the genome, while the
17 structural viral proteins (VPs) are encoded by the right, and complex, hairpin-
18 like DNA secondary structures are present at both genomic termini (2). Small
19 satellite proteins and an assembly-activating protein have been discovered as
20 products of open reading frames overlapping the right hand expression
21 cassette (3, 4)

22 Numerous novel parvoviruses have been identified in recent years,
23 primarily via approaches based on high throughput sequencing (HTS) (5-11).
24 In addition, progress in whole genome sequencing of eukaryotes has revealed
25 that sequences derived from parvoviruses occur relatively frequently in animal
26 genomes (12-17). These endogenous parvoviral elements (EPVs) are derived
27 from the genomes of ancient parvoviruses that were incorporated into the
28 gene pool of ancestral host species. This can presumably occur when
29 infection of a germline cell leads to parvovirus-derived DNA becoming
30 integrated into host chromosomes, and the cell containing the integrated
31 sequences then goes on to develop into a viable organism (18). Many EPVs
32 are millions of years old, and are genetically “fixed” in the genomes of host
33 species (i.e. all members of the species have the integrated EPV in their

1 genomes). Such ancient EPV sequences are in some ways analogous to
2 “parvovirus fossils”, since they preserve information about the ancient
3 parvoviruses that infected ancestral animals.

4 Among the novel parvovirus groups identified via sequencing, one -
5 provisionally labeled ‘chapparvovirus’ - stands out as being particularly
6 unusual. These viruses, which have been primarily reported via metagenomic
7 sequencing of animal feces, derive their name from an acronym (CHAP),
8 referring to the host groups in which they were first identified (CHiropteran-
9 Avian-Porcine) (15, 16, 19, 20). Subsequently several additional
10 chapparvovirus sequences have been reported, including some that were
11 identified in whole genome sequence (WGS) data derived from vertebrates,
12 including reptiles, mammals, and birds (9). These sequences were picked up
13 by *in silico* screens designed to detect EPV. However, since all the
14 chapparvovirus sequences identified in WGS data lack clear evidence of
15 genomic integration, it is thought likely that they actually derive from infectious
16 chapparvovirus genomes that contaminated WGS samples, rather than from
17 EPVs (9).

18 Until relatively recently, evidence that the chapparvoviruses detected
19 via sequencing actually infect vertebrate hosts has been lacking. However, a
20 recent study has claimed to demonstrate that a chapparvovirus called ‘mouse
21 kidney parvovirus’ (MKPV) circulates among laboratory mice populations, in
22 which it causes a kidney disease known as ‘inclusion body nephropathy’ (21).
23 These findings imply that chapparvoviruses represent a potential disease
24 threat to humans and domestic species. In addition, they have raised interest
25 in the use of these viruses as experimental tools.

26 In this study, we perform comparative analysis of ChPV genomes and
27 ChPV-derived EPVs, revealing new insights into the biology and evolution of
28 this poorly understood group.

1 RESULTS

2 Identification and characterization of novel chapparvovirus sequences

3 We systematically screened published WGS data and identified a total
4 of fifteen previously unreported, chapparvovirus-derived DNA sequences.
5 Two were identified in vertebrates and thirteen in invertebrates (**Table 1**). The
6 majority of the novel chapparvovirus sequences identified in our screen were
7 derived from the replicase (*rep*) gene, but we identified complete sequences
8 derived from both the *rep* and capsid (*cap*) genes in two species: the Gulf
9 pipefish (*Syngnathus scovelli*) and the black widow spider (*Latrodectus*
10 *hesperus*). Partial *cap* genes were identified in the scarab beetle (*Oryctes*
11 *borbonicus*), taurus scarab (*Onthophagus taurus*) and Chinese golden
12 scorpion (*Mesobuthus martensii*) elements (**Figure 1**).

13 Two of the novel sequences were identified in WGS assemblies of fish
14 in the family Syngnathidae, including the tiger tail seahorse (*Hippocampus*
15 *comes*) and the Gulf pipefish. The seahorse sequence clearly represented an
16 EPV (**Figure 2a**), while the pipefish sequence lacked flanking genomic
17 sequences and appeared to derive from an exogenous virus. In addition, all
18 hits of invertebrate origin clearly represented EPVs (**Figure 2a**). However,
19 none shared homologous flanking sequences, suggesting that they each
20 derive from distinct germline incorporation events.

21 We used phylogenetic approaches to reconstruct the evolutionary
22 relationships of newly identified chapparvovirus sequences to previously
23 reported parvoviruses. We reconstructed evolutionary relationships using
24 maximum likelihood-based approaches and an alignment spanning the
25 tripartite helicase domain of the NS protein (**Figure 3a**). Strikingly,
26 reconstructions indicated that the family *Parvoviridae* four major clades, rather
27 than the two that have historically been recognised (1). Of the four major
28 lineages, one corresponds to the subfamily *Parvovirinae* as recognised by
29 current taxonomic schemes. However, the subfamily *Densovirinae* appeared
30 to be split into two clades; one encompassing all ambisense densoviruses
31 along with viruses of the genus *Iteradensovirus* (which have monosense
32 genomes) and the second comprised by genera *Hepan-*, *Brevi-* and

1 *Penstyldensovirus*. Moreover, a fourth parvovirus lineage was evident,
2 comprised of the ChPVs and ChPV-derived EPVs.

3 The branching relationships among ChPVs were not fully resolved by
4 phylogenetic analysis of the helicase domain. The putative large replicase
5 proteins (NS1) of ChPVs displayed a high degree of amino acid variability,
6 particularly toward their N- and C-termini. However, a region ~500-aa-long
7 could be aligned reliably throughout all complete and partial entries previously
8 proven to cluster within the chapparvovirus lineage in case of the NS
9 helicase-based inference. According to this, the ChPV clade is comprised of
10 three significantly-supported monophyletic lineages. One of these includes
11 ChPVs sampled from amniotes (reptiles, birds, and mammals), in which three
12 further consistently well-supported sublineages could be observed (**Figure**
13 **3b**). The amniote ChPVs form a sister clade to EPVs found in the arthropod
14 subphyla Chelicerata (arachnids, camel spiders, scorpions, whip scorpions,
15 harvestmen, horseshoe crabs and kin) and Myriapoda (millipedes, centipedes
16 and kin) as well as with the syngnathic fish associated sequences. Within this
17 clade, however, only the grouping of spider and syngnathic fish sequences
18 was supported. A third lineage was also observed, containing sequences from
19 the arthropod subphylum Hexapoda (insects, springtail, and forcepstail).
20 Monophyly of the beetle EPVs is robustly supported (**Figure 3b**).

21

22 Comparative analysis of previously reported chapparvovirus genomes

23 We performed comparative analysis of nine previously sequenced
24 ChPV genomes so that we might: (i) identify genome features that
25 characterise these viruses, and (ii) make inferences about aspects of ChPV
26 biology and evolution (**Figure 2b**). ChPV genomes tend toward the shorter
27 end of the parvovirus genome size range (~4kb). They encode a relatively
28 long *rep* gene, and a relatively short *cap* gene. The *rep* gene product (NS) is
29 ~650 amino acids (aa) in length, with the longest example being the 668 aa
30 long protein encoded by *Desmodus rotundus* chapparvovirus (DrChPV).
31 Chapparvovirus NS proteins contain ATPase and helicase domains, but these
32 are the only regions exhibiting clear homology to those found in other
33 parvovirus groups (**Figure 2b**). Overlapping the *rep*, a predicted minor ORF,

1 ~220 aa in length, is located in a position equivalent to that of the
2 nucleoprotein (NP) ORF found in certain *Parvovirinae* genera (i.e. *Ave-* and
3 *Bocaparvovirus*). However, it should be noted that the polypeptide encoded
4 by this gene – which we tentatively refer to as NP - exhibits no significant
5 similarity to any other parvovirus NP proteins. Secondary structure predictions
6 indicate that the vast majority of the NP protein has a mostly helical structure,
7 with numerous possible phosphorylation sites as well as a potentially protein-
8 binding disordered N-term (**Figure S2**). Together, these observations suggest
9 a non-structural function. The NP ORF, although of similar length in all
10 genomes, has no canonical start codon in case of porcine parvovirus 7
11 (PPV7) and simian parvo-like virus 3 (SPV3). This would imply that in these
12 viruses, splicing of the *rep* RNA is required for expression of the NP protein.

13 All chapparvoviruses appear to be characterized by relatively short VP
14 sequences that are between 450- 500 aa in length, as opposed to ~650-820
15 aa as found in most other parvoviruses (the exception being the brevi- and
16 penstyldensoviruses which encode an even shorter VP). In all
17 chapparvoviruses, the first methionine of the VP ORF is preceded by potential
18 coding sequence, and in all published ChPV sequences, a canonical splice
19 acceptor site is located immediately upstream. Possibly, the VP ORF encodes
20 only the major capsid protein, and there may be other versions of this VP
21 protein that are elongated at the N-terminus, and are incorporated into the
22 capsid at a lower copy number, as found in the majority of parvoviruses (1).
23 However, the only splice donor sites we identified are located relatively far
24 upstream. In MKDV, however, there are two large introns present, putting
25 these upstream exons in frame with the VP encoding exon.

26 The VP proteins of chapparvoviruses share no significant sequence
27 similarity with those of other parvoviruses. Interestingly, however, structural
28 similarity with erythro-, proto, and bocaparvoviruses can be detected for VP
29 using fold recognition (**Figure 4**).

30 In addition to their fundamental NS-NP-VP genome organization,
31 chapparvoviruses encode various additional small ORFs. ORF1 is predicted
32 to encode a small protein of approximately 15 kDa that contains a putative
33 nucleus localization signal (NLS) in its C-terminal region. ORF1, which

1 partially overlaps the N-terminal region of NS, is present in all genomes
2 except PPV7. However, since the PPV7 genome also lacks the corresponding
3 region of NS, this likely reflects a 5' truncated genome sequence. The same is
4 the case for turkey parvovirus (TPV2), although the C-terminal encoding
5 region of the putative ORF1 protein could be revealed.

6 A second additional, putative ORF is present in only two of the
7 chapparviruses examined here: PPV7 and simian parvo-like virus 3. This
8 ORF, referred to as ORF2, is located downstream from ORF1 in a position
9 completely overlapping the NS ORF. The TPV2 genome also contains a
10 unique, presumably genome-specific additional ORF (ORF4) that overlaps the
11 C-term encoding region of VP, and may encode a predicted 17 kDa protein
12 (**Figure 2b**). Interestingly, this ORF was absent from the other, closely-related
13 avian entries.

14 Analysis *in silico* revealed at least three potential promoters in
15 chapparvirus genomes. One of these is conserved throughout the clade,
16 and is located upstream of all coding features, indicating that it likely drives
17 early expression of virus genes. Moreover, its presence has been confirmed
18 in MKDV by sequencing of cDNA derived from infected mouse tissue. None
19 of the other potential promoters proved to be functional in case of MKDV,
20 nevertheless. The MKDV transcriptome includes three confirmed transcripts
21 to undergo splicing. Out of these, however, only the one with the shortest
22 intron could be predicted to exist in case of all examined GenBank entries
23 (**Figure 2b**). Interestingly, DrChPV, along with further rodent-derived ChPVs
24 was predicted to possess the large intron of the putative VP transcript, hence
25 displaying a strikingly MKDV-like transcription mechanism, despite of missing
26 an acceptor site upstream the NP start codon.

27 In all ChPVs examined, with the exception of the 5' truncated entries,
28 we identified two potential polyadenylation signals in positions equivalent to
29 those found in MKDV (Roedinger et al 2018), This implies that the
30 polyadenylation strategy is a conserved feature of chapparvirus
31 transcription.

32

33 Characterization of syngnathid chapparviruses and EPVs

1 We identified two chapparvoviral sequences in syngnathid fish. One of
2 these, identified in the genome of the Gulf pipefish (*Syngnathus scovelli*)
3 occurred within a relatively short scaffold (4002 nt). The entire scaffold was
4 comprised of viral sequence, displaying truncated but nonetheless detectable
5 J-shaped terminal hairpin-like structures (**Figure S1**). This suggests it likely
6 represents a virus contaminant, as suspected for other ChPV sequences
7 recovered from vertebrate WGS data (9). The virus from which this sequence
8 was presumably derived was designated *Syngnathus scovelli* chapparvovirus
9 (ScChPV). The ScChPV genome encodes a long NS ORF (807 aa), a
10 strikingly short VP (367 aa) and a ChPV-like NP. Furthermore, a homologue
11 of the ORF2 protein found in the amniote parvoviruses PPV7 and SPV3 was
12 present. A predicted ORF was present in a genomic position equivalent to that
13 of ORF1, found in amniote ChPVs. However, the predicted protein sequence
14 did not disclose any detectable similarity to its amniote counterpart. ORF6,
15 identified in partial overlap with the VP C-term encoding region, has the
16 capability to encode a small protein of 27.2 kDa (239 aa), which demonstrated
17 no sequential similarity to any GenBank entries to date. Fold recognition,
18 however, revealed possible structural similarity to viral structural proteins,
19 including the major envelope glycoprotein of Epstein-Barr virus (PDB ID:
20 2H6O chain A, p=0.012), the minor viral protein of the Sputnik virophage
21 (PDB ID: 3J26, chain N, p=0.017) and the surface region of *Galleria*
22 *mellonella* ambidensovirus (PDB ID: 1DNL, p=0.021) (**Figure 4**). These
23 findings imply ORF6 may encode an additional structural protein in addition to
24 VP.

25 A partial ChPV-like sequence identified in the genome of the tiger
26 seahorse (*Hippocampus comes*) was flanked by extensive stretches of host
27 genomic sequence, establishing that, unlike the ScChPV sequence identified
28 in the Gulf pipefish genome assembly, it likely represents an EPV rather than
29 a virus. Interestingly, however, phylogenies showed that both sequences
30 obtained from syngnathid fish are relatively closely related, and cluster
31 together with high bootstrap support (**Figure 3a and b**).

32

33 Characterization of chapparvovirus-derived EPVs in invertebrate genomes

1 Via screening of WGS data derived from invertebrate species, we
2 identified a total of 13 EPV sequences that disclosed a relatively close
3 phylogenetic relationship to ChPVs. These elements showed varying degrees
4 of degradation (**Figure 1**). In many cases only genome fragments were
5 detected, and these usually included numerous nonsense mutations (**Table**
6 **1**). ChPV-derived elements were detected in three major arthropod clades
7 that primarily occupy terrestrial habitats, namely arachnids of Chelicerata,
8 chilopods of Myriapoda, as well as hexapod insects and entognaths.

9 Among the EPVs we identified, the most complete were identified in
10 the Western black widow spider (*Latrodectus hesperus*). Two elements
11 spanned the *rep*, *cap*, and NP genes, and a homologue of the *orf1* gene
12 found in ScChPV (**Figure 2a**). At aa level of the putative NS1, however, these
13 elements displayed only 62% identity. Although the *rep* of the element
14 ChPV.2-*Latrodectus hesperus* appeared to encode the complete, 690-aa-long
15 protein product, this ORF contained nonsense mutations in case of the other
16 arachnid element, ChPV.3-*Latrodectus hesperus* (**Table 1**). In contrast,
17 ChPV.3-*Latrodectus hesperus* displayed an undisrupted ORF1 of 113 aa, a
18 slightly smaller homologue of its ScChPV counterpart as well as an intact *cap*,
19 capable of encoding a 386-aa-long VP. In ChPV.2-*Latrodectus hesperus* only
20 an N-term truncated ORF1 was present at the length of 37aa, along a *cap*
21 disrupted by several nonsense mutations. Furthermore, the *cap* of this
22 element was found to include an insertion of 74 aa, suspected to originate
23 from a yet unknown repetitive element (revealed by sequence comparisons to
24 be interspersed throughout the *L. hesperus* genome). ChPV.3-*Latrodectus*
25 *hesperus* from the other hand, appeared to include an intact upstream region
26 of the genome, revealing an additional small ORF of 81 aa length directly
27 upstream the ScChPV ORF1 homologue, designated ORF1-Lh. This ORF
28 disclosed no detectable homology to any entries to date. Upstream this ORF
29 a potential promoter sequence could be revealed with high confidence (0.98
30 of 1). Both elements included complete, NP-encoding ORFs of 233 aa,
31 although a canonical ATG start codon could be revealed only in case of one
32 element (ChPV.3-*Latrodectus hesperus*).

1 We revealed the existence of two more elements in the Western black
2 widow spider genome, although these only spanned the *rep* gene, disrupted
3 by several nonsense mutations. However, one element (ChPV.4-*Latrodectus*
4 *hesperus*) encoded nearly complete NS and NP genes, as well as a complete
5 homologue of the ScChPV ORF1 gene. The true extent of preservation could
6 not be assessed for this EPV as it occurred on a short scaffold that terminated
7 within the EPV *rep* sequence. The putative NS1 was 0.8% identical to its
8 counterpart in ChPV.2-*Latrodectus hesperus* at aa level. Interestingly, directly
9 upstream this EPV another copy of the *rep* was present, encoding only the
10 first 221 aa of the putative NS1 protein. This element also contained a
11 complete ScChPV ORF1 homologue as well as ORF1-Lh, which did not
12 display an ATG start codon in this case. The upstream promoter could only be
13 identified with a much lower score of 0.6. ChPV.5-*Latrodectus hesperus*
14 displayed a highly divergent, partial *rep* of 216 aa with only 0.42% identity to
15 the ChPV.2-*Latrodectus hesperus* NS1 at aa level (**Figure 1**). This element
16 clustered outside the monophyletic branch comprised by the three other EPVs
17 of the same species (**Figure 3**).

18 A single chapparvovirus-derived EPV was identified in a second
19 arachnid species - the Chinese golden scorpion (*Mesobuthus martensii*). This
20 element was identified in a relatively short unplaced scaffold, and comparison
21 to the *Latrodectus* elements indicated that the contig was truncated within the
22 EPV, and consequently the true extent of its preservation could not be
23 assessed. Nevertheless, ORFs disclosing homology to the NS, NP and VP
24 proteins could be identified. While the first 100 or so codons of the NS ORF
25 were absent, a complete NP ORF was detected, along with the first 46
26 codons of VP. All three ORFs were disrupted by frameshifts and stop codons
27 (**Table 1**). No homologues of any alternative ORFs identified in other
28 chapparvovirus genomes could be revealed.

29 We identified a chapparvovirus-derived EPV in the genome of a
30 myriapod - the European centipede (*Strigamia maritima*). This element
31 displayed partial homologues of the NS and NP encoding ORFs, both of
32 which contained large deletions (**Figure 2b**) as well as numerous nonsense
33 mutations (**Table 1**). Moreover, the NS ORF was disrupted by an extensive

1 stretch of an insertion of unknown origin. No homologues of any of the above-
2 mentioned alternative chapparvoviral ORFs could be identified in this
3 endogenous sequence.

4 Seven chapparvovirus-derived EPVs were identified in hexapod
5 arthropods (subphylum Hexapoda). One occurs in the genome of a bristletail
6 species - the Northern forcepstail (*Catajapyx aquillonaris*) - belonging to the
7 entognath order Diplura. The other six were identified in three species
8 belonging to the vast insect order Coleoptera: the emerald ash borer (*Agrilus*
9 *planipennis*), the taurus scarab (*Onthophagus taurus*), and the scarab beetle
10 (*Oryctes borbonicus*). The bristletail element contains a C-terminal truncated
11 *rep* of at least 250 aa and a near full-length NP ORF. The partial *rep* was
12 intact, but the NP ORF is disrupted and highly divergent, showing significant
13 sequence similarity only in the conserved core region of the putative protein.
14 The ash borer element occurs in a scaffold that is ~1 kb in length. One end of
15 this scaffold contains 592 nt region exhibiting homology to the NS ORF, which
16 harboured an N-terminal deletion of at least 200 aa.

17 In the case of the taurus scarab element ChPV.10-*Onthophagus*
18 *taurus*, the almost complete NS ORF could be identified, disrupted by
19 numerous nonsense mutations (**Table 1**). Interestingly, another endogenous
20 element of parvoviral origin was detected at the same locus. This element
21 encodes an intact, potentially fully-expressible NS gene, homologous to the
22 NS1 of ambidensoviruses (genus *Ambidensovirus*) and discloses similarity to
23 a recently reported ambidensovirus species, detected only at cDNA level in
24 the transcriptome of two bumble bee species (*Bombus cryptarum* and *B.*
25 *terrestris*) (22). An additional ORF was present in this ambidensoviral
26 element, overlapping the reputative NS1 gene, which harboured no significant
27 similarity to any sequences deposited to GenBank to date. In its derived aa
28 sequence, however, a homeobox domain could be revealed, also in intact,
29 potentially expressible state. The other two elements of the taurus scarab
30 genome located in the same assembly scaffold, only 2540 nts apart from each
31 other. Both EPV consisted of only a partial ORF, which disclosed similarity to
32 chapparvoviral *reps*. None of these elements encompassed the tripartite
33 helicase domain, hence they were not included in phylogenetic inference.

1 Two EPVs were derived from the scarab beetle genome. One of these,
2 designated ChPV.13-*Oryctes borbonicus*, harboured a near complete *rep* at
3 402 aa as well as a short, partial *cap*, capable of encoding only the first 33 aa
4 of the putative VP. The region of *rep* homology occurred within an ORF that
5 was not disrupted by any frame shifts and extended without disruption
6 upstream and downstream, suggesting a putative longer gene product -
7 potentially encoding a longer, divergent NS protein - to be present. However,
8 these regions did not disclose sequence similarity to any proteins hitherto
9 deposited to GenBank. The *Oryctes* element (EPV-ChPV.14-*Oryctes*), in
10 contrast, included only a heavily truncated NS of 254 aa (**Figure 1**).

11

12 Structural characteristics of chapparvovirus capsids

13 We built 3D homology models to facilitate comparison of
14 chapparvovirus capsid structures to those found in other parvoviruses. The
15 derived polypeptide sequence of the complete VP ORF encoded by DrChPV
16 was subjected to fold recognition, to identify suitable templates for homology
17 modelling. This comparative analysis showed that the most similar VP protein
18 occurs in parvovirus H1 (genus *Protoparvovirus*) (PDB ID: 4G0R, $p=9e-05$),
19 and this sequence was therefore used as the template for homology
20 modelling. To overcome the stochastic aspect of model construction, due to
21 the lack of sequence identity and the non-homologous nature of the ChPV VP
22 genes to other parvoviral VPs, we used the final model obtained in this
23 analysis as a template for analysis of four further ChPVs - rat parvovirus 2,
24 PPV7, TPV2, and pit viper chapparvovirus. The pitfalls of using models as
25 templates has to be noted here, however, this method ensured that only those
26 regions showed structural variability, which would probably do so in case of
27 the actual capsid structures.

28 We examined the VP sequences of two representatives of the second
29 major ChPV clade (**Figure 3b**) – one derived from a presumably exogenous
30 virus (ScChPV) and one from an EPV (ChPV.3-*Latrodectus*). Fold recognition
31 identified, for the VP protein encoded by ScChPV dependoparvovirus VP
32 proteins as potential templates (adeno-associated virus 8, PDB ID: 2QA0,
33 $p=8e-04$; Adeno-associated virus rh32.33, PDB ID: 4IOV, $p=9e-04$), while for

1 the VP encoded by the black widow spider EPV the most reliable hit was the
2 VP4 protein of an iteradensovirus (*Bombyx mori* iteradensovirus, PDB ID:
3 3P0S, $p=8e-04$). When superimposing the obtained models with the VPs of
4 AAV8 and BmDV, however, structural similarity only covered the jelly roll core
5 and the α A helix, whereas only the EC loop out of the traditionally more
6 variable surface loops.

7 Modelling indicated that the ChPV VP monomer harbours an eight-
8 stranded β -barrel 'jelly roll' core and the α A helix at the twofold symmetry axis
9 as found in all members of the family *Parvoviridae* to date (23) (**Figure 4a**).
10 Equivalent of all short strands were present (β -C, H, E, F), as well to four out
11 of the five longer strands (β -B, D, I, G). However, no structural analogue to the
12 outmost β -A could be identified (**Figure 4**). Examining the secondary structure
13 prediction results further disproved the existence of a β -A analogue, indicating
14 β -B to be the closest to the N-term. The first strand of the *Syngnathus scovelli*
15 ChPV VP, appeared to fold outside of the jelly roll, leaving the longer side of
16 the barrel without a β -B, comprised of only three strands, namely D, I and G,
17 despite of a complete upper, CHEF sheet (**Figure 4**). When modelling the
18 complete T=1 capsid polymer, this manifested as a hole, which is normally
19 covered by β -B even in case of the smallest parvoviral capsids (**Figure 4b**).
20 All ChPVs as well as the only ChPVe VPs displayed two canonical loops
21 surrounding their fivefold axes, linking sheets D with E at the channel and
22 sheets H with I on the floor surrounding the channel. In case of the amniote
23 ChPVs, the pore displayed a tight opening. The sequence of the DE loop
24 varied to some extent among these seven entries, which also manifested in
25 the models. The HI loop was, however, highly conserved throughout,
26 containing only one variable position between the amniote ChPVs.

27 We mapped the chapparvoviral VRs identified by VP alignments
28 (**Figure S3**) to both VP monomers and complete capsids, to examine how
29 they manifest on the virion surface and make comparisons to parvoviruses of

1 known structure, represented by the minute virus of mice (MVM) capsid
2 structure, as the prototype virus of subfamily *Parvovirinae* (PDB ID: 1Z14)
3 (**Figure 4**). Out of ten chapparvoviral VRs identified (VR 1 to 10), presented
4 by **Figure S3a** only VR1, VR2, and VR9 proved to be similarly positioned,
5 hence probably analogous to the parvoviral counterparts. Some VRs
6 appeared to be positioned at luminal surface of the chapparvovirus capsid
7 (VR4 in all cases and VR8 of the amniote ChPVs, whereas VR6 of the
8 *Latrodectus* element and ScChPV), unlike all parvoviruses studied to date
9 with the exception of bovine parvovirus, a bocaparvovirus (24) in which VRVIII
10 shows this configuration. Since the ChPV VRs appeared to be non-
11 homologous to those established for either proto- or dependoparvoviruses, we
12 re-defined them by numbering from N to C-term.

13 In addition to their distinctive VRs, ChPVs ubiquitously appeared to
14 harbor a highly variable C-terminal region, with a length varied between 12 to
15 62 residues. The ChPV VP variable C-term appears to be buried in most
16 cases with the exception of ScChPV, where it is probably exposed. In case of
17 the *Latrodectus* element, it forms the luminal surface of the threefold, whereas
18 in the case of fish and amniote ChPVs it is located at the twofold (**Figure 4a**).

19 The *ScChPV* and *ChPV.3-Latrodectus hesperus* VP lacked a VR6
20 homologous to that of the amniote ChPVs albeit displayed variation in another
21 position instead, still in the sixth-place counting from the N-term (**Figure S3b**).
22 Moreover, both of them displayed truncated VRs 3, 5, and 7, compared to
23 their amniote counterparts. VR9, furthermore, was absent from the ScChPV
24 VP, whereas VR 10 was missing from the VP of *ChPV.2-Latrodectus*
25 *hesperus* (**Figure S3b**). As for the surface, the largest variable region for
26 amniote ChPVs, namely DrChPV is VR7, forming the entire three-fold
27 protrusions, with VRs 1, and 9, forming small protrusions surrounding the
28 aforementioned peaks.

29 The complete capsid models of non-amniote chapparvoviruses were
30 observed to harbor surface features that are strikingly distinct from those of
31 the amniote ones, more closely resembling the capsids of the
32 *Ambidensovirus-Iteradensovirus* clade of *Densovirinae* (see **Figure 3a**), with
33 a surface that is less spikey (**Figure 5**).

1 We also constructed a homology model of ORF6 of ScCHPV, based on
2 the minor viral protein of the Sputnik virophage (PDB ID: 3J26). This protein
3 appears to harbor multiple beta strands close to its C-term, out of which the
4 outermost could potentially fill in the aforementioned gap, caused by the lack
5 of a β B (**Figure 4b**).

6

7 **DISCUSSION**

8 Historically, the family *Parvoviridae* has always been comprised of two
9 subfamilies, with specificity for vertebrate or invertebrate hosts being the
10 major demarcation criterion (2). This division was initially supported by
11 phylogenetic inference, however, as the number of densoviral genera
12 increased, the heterogeneity of densoviruses, specifically their segregation
13 into two clades, has not gone unnoticed (1). Our study provides further
14 evidence that the traditional division of parvoviruses into vertebrate-specific
15 and invertebrate-specific subfamilies no longer holds - rather, it supports the
16 division of the *Parvoviridae* into four major subgroups: the *Parvovirinae*, a split
17 *Densovirinae*, and the Chapparvoviruses as illustrated in **Figure 3a**.

18 The data presented here show that chapparvoviruses infect an
19 exceptionally broad range of hosts, including both vertebrates and
20 invertebrates. Furthermore, we show that chapparvoviruses found in fish are
21 more closely related to those that infected ancestral arachnoid arthropods
22 than they are to those that infect amniote vertebrates such as rodents (**Figure**
23 **3b**). These findings suggest that chapparvoviruses have been transmitted
24 between distantly related host species in the past.

25 Phylogenies indicate that all amniote ChPVs have a common origin
26 (**Figure 3a**), consistent with the overall conservation of their genome
27 organization and some aspects of predicted transcriptional strategy (**Figure**
28 **2b**). Previous studies suggested that these viruses might have broadly co-
29 diverged with host species (9). The present, expanded data set implies that
30 some transmission of ChPVs between distantly related vertebrate classes
31 may have occurred among (**Figure 3b**) - though it should be kept in mind that
32 almost all amniote ChPVs have been identified via metagenomic sequencing

1 of environmental samples, mostly fecal viromes, thus their true host affiliations
2 remain uncertain.

3 Sequences derived from parvoviruses occur relatively frequently in
4 animal genomes (13, 14, 17, 25, 26). However, these sequences
5 overwhelmingly derive from a small proportion of parvovirus lineages. For
6 example, ambidensovirus-derived EPVs dominate invertebrate genomes (14),
7 whereas vertebrate EPVs almost exclusively derive from
8 dependoparvoviruses, protoparvoviruses, or viruses closely related to these
9 two genera (12, 13, 25, 26). In this study we found no trace of ChPV-derived
10 sequences in tetrapod genomes, despite recent evidence that ChPVs infect
11 this host group (21). By contrast, ChPV-derived sequences are relatively
12 common in arthropods, with the genomes of some species harboring multiple,
13 independently acquired ChPVs (**Table 1**). The tendency of EPVs to derive
14 from a subset of parvovirus genera likely has biological underpinnings. For
15 example, in vertebrates it may reflect the tendency of dependoparvoviruses to
16 integrate into host DNA, and/or the requirement of protoparvoviruses to
17 initiate DNA damage response (DDR) during replication (27, 28). Similar
18 features of the viral life cycle could account for the biased distribution of
19 ChPV-related sequences in species genomes - i.e. arthropod and fish ChPVs
20 might have adopted a replication strategy that favors germline integration,
21 whereas that of amniote ChPVs precludes it. Notably, some arthropod
22 species have integration sites containing multiple independently acquired
23 EPVs of both ChPV and ambidensovirus origin, suggesting that hotspots of
24 parvovirus integration and/or fixation might exist in their genomes.

25 Our discovery of ChPV-derived elements in fish and arthropod
26 genomes establishes that ChPVs can infect these species in addition to
27 mammals (21). Moreover, it provides evidence that the ChPVs are probably
28 an ancient lineage of parvoviruses. Although we did not identify any
29 orthologous ChPV insertions, the EPVs described here show extensive
30 evidence of germline degradation. Through comparison to studies of EPVs in
31 mammals (in which several orthologous EPVs have been described), it
32 appears likely that chapparvoviruses have been present in animals for many
33 millions of years. Moreover, as the hexapod EPVs appear to be monophyletic

1 and mirror the evolution of their host species, the age of ChPVs could
2 possibly correlate with the Insecta-Entognatha split, suggesting a minimum
3 age of 400 million years (29).

4 Through comparative analysis of EPVs and chapparvoviruses, we
5 show that chapparvovirus genomes exhibit a number of defining
6 characteristics. Firstly, all possess a short, monosense genome, encoding a
7 relatively large NS and a relatively short VP. The short VP proteins of ChPVs
8 are clearly homologous to one another, but show no similarity to those found
9 in other parvovirus lineages. Similar to those found in the penstyl-, hepan- and
10 brevidensoviruses, the VP proteins of ChPVs lack PLA2 domains. Notably
11 these are also the genera to which ChPVs are most closely related in NS-
12 based phylogenies (**Figure 3a**).

13 Secondly, chapparvoviruses typically encode multiple additional gene
14 products besides the replicase and capsid. To begin with, almost all encode a
15 nucleoprotein (NP) gene in an overlapping frame with *rep*. In this report, we
16 show that putative NP ORFs are present in ChPV-derived EVEs, suggesting it
17 is an ancestral, conserved feature of these viruses. It's absence, however,
18 from the coleopteran lineage is intriguing, as it is still present in the EPV of the
19 hexapod stem group Diplura of Entognatha. Phylogenetic reconstructions
20 (and the extensive overlap with *rep*) imply it was acquired ancestrally and
21 independently lost in the lineage derived from members of the hexapod crown
22 group, Coleoptera (**Figure 3b**).

23 ChPV genomes appear to encode several predicted auxiliary ORFs in
24 addition to NP. A functional role for auxiliary ORF1 is supported by (i) it's
25 conservation across the entire amniote ChPV clade, and; (ii) limited
26 experimental data indicating it is expressed in MKPV via a spliced transcript.
27 Auxiliary ORF2 was only identified in a small subset of ChPV genomes, but a
28 functional role for this ORF is suggested by the presence of homologues in
29 distantly related ChPVs of amniotes and fish (see **Figure 3b**). Interestingly,
30 although all ChPVs appear to express ORF1 via splicing of a small intronic
31 sequence (**Figure 2**), those harboring an ORF2 homologue are predicted to
32 lack the peculiar large introns found in expression of MKPV NP and VP
33 transcripts (21). ScChPV lacks an ORF1 homolog, but contains a predicted

1 reading frame in the corresponding position. Homologues of this ScChPV
2 ORF1 variant are present in all three arachnid EPVs, although not in the first,
3 but in the second position. As only the three *Latrodectus* EPVs possess a
4 homologue of ORF1-Lh, it is possible that this small ORF became
5 incorporated only after the split from the syngnathid fish lineage, whereas the
6 incorporation of ScChPV ORF1 predates it. The distribution of homologous
7 auxiliary genes across phylogenetic lineages of ChPVs implies that distinct
8 lineages have acquired and/or lost these genes on multiple, independent
9 occasions.

10 MKPV has been reported to possess only one promoter and two
11 polyadenylation signals, as well as an extensive number of spliced transcripts.
12 This transcription pattern, however, appears to be unique to only one lineage
13 of the amniote ChPVs, comprising of rodent, chiropteran, primate, and
14 reptilian entries. As both the avian ChPV sister clade as well as the lineage
15 including PPV7 appear to display genome organization specific for these
16 groups suggesting that ChPVs infecting host groups outside of the MKPV
17 lineage may utilize distinct transcription strategies.

18 Despite the potential pitfalls of homology modeling, and the use of
19 distinct templates to reconstruct both the VP monomer and capsid structures,
20 we obtained remarkably similar predicted structures for VP sequences found
21 in closely related viruses/EPVs. Since the viral capsid plays an important role
22 in mediating the interactions between parvoviruses and their hosts,
23 comparisons of capsid structures can potentially reveal insights into
24 parvovirus biology. Homology modelling indicated that ChPV VPs would
25 assemble into a complete, T=1 icosahedral capsid, despite their relatively
26 small size. Furthermore, their predicted structures are remarkably similar to
27 that found in other parvoviruses, despite the lack of any detectable similarity
28 in the sequences of their VP proteins. Similarities include the presence of a
29 conserved jelly roll core and α A helix, the existence of the DE and HI loops,
30 and the presence of identifiable VRs. Interestingly, the amniote ChPV
31 capsids appear to possess the same number of VRs as most of the vertebrate
32 parvoviruses of subfamily *Parvovirinae*, even if only a few of them (namely
33 VRs 1, 2 and 9) proved to be analogous features. In these virus capsids,

1 variations were most prominent among the threefold peaks and protrusions as
2 well as the two-fold depression, as observed in members of the *Parvovirinae*
3 (**Figure 5**). The tendency of some VRs to manifest at the luminal surface of
4 the capsid in models suggests these regions could play a role in intracellular
5 host-virus interactions. However, in order to become accessible to interact
6 elements of intracellular signaling pathways, either uncoating would be
7 necessary, or potential conformational changes to expose these regions.
8 Based on previous findings, however, the parvovirus capsid appears to traffic
9 into the nucleus intact (30, 31). Considering this, these buried regions might
10 play a role in processes linked to the nucleus. Interestingly, however, bovine
11 parvovirus, the only other parvovirus in which buried VRs have previously
12 been observed (32), is an enteric pathogen, and metagenomic studies
13 suggest that amniote ChPVs (with the exception of MKPV) are also enteric.

14 In addition to the VRs, all ChPVs seem to harbor highly-variable VP C-
15 terms. A similar phenomenon has been observed in case of
16 iteradenoviruses, where the last 40 C-terminal residues are disordered and
17 their structure could not be resolved (33). Although the location of the ChPV
18 C-term appears to vary, its association with regions that are overtly involved in
19 parvovirus-host interactions (e.g. the two- and threefold peaks) is certainly
20 intriguing.

21 MKPV is associated with pathology of the urogenital system, whereas
22 a related virus - murine ChPV - has been detected at very high prevalence in
23 murine liver tissue suggesting it is a gastrointestinal agent (34). The VPs of
24 the two, however, only differ in six aa residues, located within VR 3 and near
25 VR2 on the surface and in the buried VR4 as well as in the similarly buried
26 variable C-term (**Figure S3a**). Thus, these positions could constitute potential
27 determinants of tissue tropism in murine ChPVs. Parvoviruses subfamilies
28 *Parvovirinae* and *Densovirinae* utilize distinct strategies to stabilize their
29 icosahedral capsids (35). Vertebrate parvoviruses extend the longer side of
30 the jellyroll fold with an additional, N-terminal strand by folding back β -A to
31 interact with the twofold axis of the very same monomer, hence creating an
32 extended ABDIG sheet (36, 37). By contrast, the densovirus capsid preserves

1 the symmetric arrangement of the jellyroll fold, and possesses a β -A which is
2 a direct elongated N-terminal extension of the β -B instead, interacting with the
3 β -B strand of the neighboring monomer toward the fivefold axis (38, 39).
4 Strikingly, our data show that ChPV capsids lack β -A strands (and also the β -
5 B strand in the case of ScChPV). The functional implications of this are
6 unclear - possibly ChPV capsids are stabilized in the absence of β -A via a yet
7 unknown, additional VP. If ChPVs express additional structural proteins, they
8 are presumably encoded by spliced transcripts (given the unusually small size
9 of the *cap* gene). Alternatively, the ChPV capsid might assemble without the
10 incorporation of an additional β strand, perhaps at the cost of losing the
11 stability and resilience typical of parvoviruses in general. Potentially, this could
12 account for the apparent presence of buried VRs. Interestingly, in studies of
13 MKPV, viral proteins could be detected in the kidneys of infected mice even
14 though no assembled particles could be observed in inclusion body-affected
15 tubular cells (21). This, along with our structural predictions, suggests that the
16 ChPV strategy for uncoating and cellular trafficking might be very different
17 from that found in the *Parvovirinae* and *Densovirinae*.

18 Uniquely, the genome of ScChPV appears to include a putative
19 additional structural protein (ORF6), in addition to the above-mentioned
20 alternative ORFs. All parvoviruses to date - except those of genus
21 *Penstylidensovirus* (39) - have been reported to incorporate up to three
22 additional minor capsid proteins into the virion, which share a common C-
23 terminal region. To encode a structural protein on an entirely separate ORF
24 sharing no mutual coding sequence with *cap* would be unique. Possibly, this
25 unusual feature could be connected to the predicted lack of a β -B strand in
26 ScChPV VP protein.

27 Taken together, the data presented here establish that the ChPVs
28 belong to a parvovirus lineage that comprises a distinct lineage from all other
29 parvoviruses, and infects an exceptionally broad range of host species,
30 including both vertebrates and invertebrates. Consistent with this, their

1 relatively complex genomes exhibit numerous unique features, implying their
2 life cycle might significantly differ from what has been established in case of
3 other members of the family. These findings underscore the need for further
4 basic and comparative studies of ChPVs, both to assess their potential impact
5 on animal health, both wildlife and livestock. Furthermore, this is the first study
6 to imply that vertebrate parvoviruses are not monophyletic, as well as
7 members of the family must have evolved to infect vertebrates on at least two
8 separate occasions.

9

10 **METHODS**

11 Genome screening and sequence analysis

12 WGS data were screened for chapparvovirus sequences using the
13 database-integrated genome screening (DIGS) tool (40). ChPV sequences
14 were characterized and annotated using Artemis Genome Browser (41). The
15 BLAST program was used to compare sequences and investigate predicted
16 viral ORFs. To determine potential homology and sequence similarity even
17 between previously undescribed ORFs, we constructed a local database by
18 including all ORFs exceeding 100 aa in length derived from all the exogenous
19 and endogenous sequences incorporated into this study and used the local
20 BLAST P and X algorithms to conduct similarity searches in it. Two ORFs
21 were accepted as homologous if they gave a significant hit in case of an
22 expectation value threshold of 1.

23 Promoters were predicted using the neural network-based promoter
24 prediction server of the Berkeley Drosophila Genome Project and further
25 verified by the Promoter Prediction 2.0 server (42, 43). Splice sites were also
26 detected using the neural network-based applications of the Berkeley
27 Drosophila Genome project and SplicePort (43, 44). Polyadenylation signals
28 were predicted by the SoftBerry application POLYAH (45). To verify the
29 suitability of these applications to be capable of detecting the above
30 mentioned chapparvoviral transcription elements we ran MKDV through the
31 workflow pipeline.

32

33 Phylogenetic reconstructions

1 The derived aa sequences of ORFs disclosing homology to parvoviral
2 NS1 proteins were aligned with at least five representatives of each genera of
3 Parvoviridae, or with one representative of each species of given genus in
4 case the number of species did not exceed five. To ensure the correct
5 identification of the tripartite helicase domain, structural data was also
6 incorporated into alignment construction using T-coffee Espresso (46) and
7 Muscle (47). The full-length NS1 derived aa sequences of the chapparvovirus
8 clade were aligned by Muscle and the M-coffee algorithm of T-coffee (48).
9 Model selection was carried out by ProTest and the substitution models
10 RtREV+I+G in case of the helicase-based inference and LG+I+G for the
11 complete chapparvoviral NS1 tree were predicted to be the most suitable
12 based on both Akaike and Bayes information criteria. To infer the maximum
13 likelihood phylogenetic tree the PhyML-3.1 program was used with 100
14 bootstrap iterations (49), based on a guide tree previously constructed by the
15 ProtDist and Fitch programs of the Phylip 3.697 package (50).

16

17 Homology modelling and DNA structure prediction

18 Structural homology was detected by applying the pGenTHREADER
19 and pDomTHREADER algorithms of the PSIPRED Protein Sequence
20 Analysis Workbench (51). The same workbench was used to map disordered
21 regions using DISOPRED3 and to predict the secondary structure of the
22 complete chapparvoviral VP protein sequences via the PSIPRED algorithm.
23 The selected PDB structures were applied as templates for homology
24 modeling, carried out by the I-TASSER Standalone Package v.5.1 (52). To
25 guide the modelling, the predicted secondary structures were applied as a
26 restriction. 60-mers of the acquired purative VP monomer structures were
27 constructed by the Oligomer Generator feature of the Viper web database
28 (<http://viperdb.scripps.edu/>) (53). Surface images of the capsids were
29 rendered using the PyMOL Molecular Graphics System (54). Capsid surface
30 maps and VP monomer superposition were carried out by UCSF Chimera
31 (55). To predict the presence of potential DNA secondary structural elements
32 the DNA Folding Form algorithm of the mFold web server was utilized (56).

33

1 **Table 1. Novel chapparvovirus sequences identified in this study**

2

Host common name	Host scientific name	ID ^a	Gene content ^b	Nonsense Mutations (stop codons; frame shifts)
Vertebrates				
Gulf pipefish	<i>Syngnathus scovelli</i>	ScChPV	<u>rep+cap</u>	0; 0
Tiger tail seahorse	<i>Hippocampus comes</i>	ChPV.1	<u>rep</u>	2; 2
		ChPV.2		
Invertebrates				
Black widow spider	<i>Latrodectus hesperus</i>	ChPV.3	<u>rep+cap</u>	4; 1
		ChPV.4	<u>rep+cap</u>	3; 1
		ChPV.5	<u>rep</u> *	3; 3
		ChPV.6	<u>rep</u>	4; 2
Chinese scorpion	<i>Mesobuthus martensii</i>	ChPV.7	<u>rep+cap</u> *	2; 3
European centipede	<i>Strigamia maritima</i>	ChPV.8	<u>rep</u>	2; 3
Northern forcepstail	<i>Catajapyx aquilonaris</i>	ChPV.9	<u>rep</u>	0; 0
Emerald ash borer	<i>Agilus planipennis</i>	ChPV.10	<u>rep</u> *	
Taurus scarab	<i>Onthophagus taurus</i>	ChPV.11	<u>rep</u>	2; 3
		ChPV.12	<u>rep</u>	0; 0
		ChPV.13	<u>rep</u>	2; 1
Rhinoceorous beetle	<i>Oryctes borbonicus</i>	ChPV.14	<u>rep</u>	2; 0
		ChPV.15	<u>rep</u>	0; 0

3

4 **Footnote:** ^a For sequences that are presumed to derive from viruses, the proposed name of the virus is
5 shown. ^b Underlined names indicate the presence of the complete ORF. For endogenous parvoviral
6 elements (EPV) the locus name is given, following the standard nomenclature proposed for endogenous
7 retrovirus (ERV) loci (57), using the classifier 'ChPV'. Asterisks indicate contigs that were truncated
8 within the virus-derived portion of the sequence.

9

1 **Figure Legends**

2

3 **Figure 1.** Basic gene content of novel chapparvoviruses and chapparvovirus
4 EPV, shown in relation to a representative chapparvovirus genome (mouse
5 kidney parvovirus). Asterisks indicate contigs that were truncated within the
6 virus-derived portion of the sequence. Abbreviations: non-structural protein
7 (NS); capsid protein (VP); nucleoprotein (NP).

8

9 **Fig.2 legend:** Genome organization, position of open reading frames (ORFs)
10 and predicted cis transcription elements of chapparvoviruses and endogenous
11 chapparvoviral elements shown in six reading frames. ORFs are represented
12 by arrows, colored according to homology. In-frame stop codons are shown
13 as vertical lines. Splice donor sites are marked by white-colored, acceptor
14 sites by orange-colored bars. Blue colored bars show predicted
15 polyadenylation signals. Promoters are presented as small green arrows if
16 predicted with a higher score than 0.95 of 1 and pink if the score is between
17 0.9 and 0.95. Sequences marked by grey boxes are assessed to be
18 transcribed but not translated. (A) Near complete entries detected include the
19 exogenous, potentially circulating chapparvovirus of the gulf pipefish
20 (*Syngnathus scovelli*) as well as two, non-identical, endogenous
21 chaparvoviruses of the Western black widow spider (*Latrodectus hesperus*).
22 ChPV.2-*Latrodectus hesperus* contains a previously unidentified repetitive
23 element, present as multiple copies scattered in the *Latrodectus* genome,
24 marked by the white box within the VP gene. The element ChPV.10-
25 *Onthophagus taurus* shares its integration site with another endogenous
26 element disclosing similarity to ambidensovirus. (B) Representatives of the
27 three basic genome organization types of exogenous amniote
28 chapparvoviruses. ORF4, however, does not occur in any of the other avian
29 chapparvoviruses.

30

31 **Figure 3 legend:** (A) Maximum likelihood phylogeny of the family
32 *Parvoviridae* based on the tripartite helicase domain of the large non-
33 structural protein NS1. (B) Maximum likelihood phylogenetic reconstructions
34 of the Chapparvovirus clade based on the complete aligned amino acid

1 sequences of the NS1. The presence or suspected presence (marked by a
2 question mark) of predicted auxiliary protein encoding open reading frames
3 (ORFs) are mapped as boxes of various colors next to each entry. Bold taxa
4 labels indicate endogenous sequences, whereas taxa labels in italics indicate
5 sequences known or believed to derive from viruses.

6

7 **Figure 4**

8 (A) Comparison of VP monomer ribbon diagrams of the protoparvovirus
9 minute virus of mice (PDB ID: 1Z14) from subfamily *Parvovirinae* to homology
10 models of an amniote, a fish and an endogenous arthropod chapparvovirus.
11 Variable regions (VRs) of the same number are marked by the same color
12 and mapped to the surface and luminal area of the icosahedral capsid model
13 constructed of 60 monomers. In case of minute virus of mice, the VRs are
14 marked by both the traditional numbering established for
15 dependoparvoviruses (Latin numerals) and by the special one applied for
16 protoparvoviruses only (Arabic numerals). Triangles mark the position of an
17 asymmetric unit within the capsid, the fivefold symmetry axis marked by a
18 pentagon, the threefolds with triangles and the twofold with an ellipsoid. (B)
19 Homology model of ORF6, the hypothetical structural protein of *Syngnathus*
20 *scovelli* chapparvovirus (ScChPV). The trimer of the ScChPV monomer model
21 reveals a gap at each subunit interaction (arrows), unlike in case of the trimer
22 of the hitherto smallest parvoviral capsid protein, *Penaeus stylirostris*
23 densovirus. The gap might accommodate ORF6 in the assembled ScChPV
24 capsid. Symmetry axes marked by the same symbols as for panel A.

25

26 **Figure 5**

27 (A) The three different predicted chapparvovirus VP structural types
28 represented by ribbon diagrams. The first panel shows sSuperposition of
29 protein monomer homology models of amniote chapparvovirus capsids,
30 including reptilian, avian, rodent, chiropteran and ungulate representatives.
31 Black arrows show variable regions (VRs) The next two panels show the
32 homology model of a fish chapparvovirus capsid monomer and an arthropod
33 endogenous chapparvoviral one. (B) Capsid surface morphology of amniote

1 chapparvovirus homology models compared to that obtained for a prototypic
2 parvovirus, minute virus of mice (MVM) (PDB ID: 1Z14 at 3.25 Å resolution).
3 Capsids are orientated by their twofold symmetry axes, as shown in the line
4 diagram. Below the comparison of homology models of complete viral capsid
5 surface morphology of the newly-identified fish chapparvovirus and arachnid
6 endogenous chapparvoviral element is shown with that of the actual capsid
7 structure of two densoviruses (subfamily *Densovirinae*, genus
8 *Ambidensovirus*) (PDB ID: 4MGU at 3.5 Å resolution for *Acheta domestica*
9 ambidensovirus and 1DNU at 3.7 Å for *Galleria mellonella* ambidensovirus).

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
4 Pesquisa do Estado de São Paulo, Brazil (Scholarships No. 17/13981-0).
5 MAM and JP are funded by NIH R01 GM109524.

6

7 References

8

9

- 10 1. **Cotmore SF, Agbandje-McKenna M, Chiorini JA, Mukha DV, Pintel DJ,**
11 **Qiu J, Soderlund-Venermo M, Tattersall P, Tijssen P, Gatherer D,**
12 **Davison AJ.** 2014. The family Parvoviridae. *Arch Virol* **159**:1239-1247.
- 13 2. **Tijssen P, Agbandje-McKenna M, Almendral JM, Bergoin M, Flegel TW,**
14 **Hedman K, Kleinschmidt J, Li Y, Pintel DJ, Tattersall P.** 2011. Family
15 Parvoviridae, p 405–425. *In* King AMQ, Adams MJ, Carstens EB, Lefkowitz
16 EJ (ed), *Virus taxonomy—Ninth Report of the International Committee on*
17 *Taxonomy of Viruses.* Elsevier/Academic Press, London.
- 18 3. **Zadori Z, Szelei J, Tijssen P.** 2005. SAT: a late NS protein of porcine
19 parvovirus. *J Virol* **79**:13129-13138.
- 20 4. **Sonntag F, Kother K, Schmidt K, Weghofer M, Raupp C, Nieto K, Kuck A,**
21 **Gerlach B, Bottcher B, Muller OJ, Lux K, Horer M, Kleinschmidt JA.**
22 2011. The assembly-activating protein promotes capsid assembly of different
23 adeno-associated virus serotypes. *J Virol* **85**:12686-12697.
- 24 5. **Siqueira JD, Ng TF, Miller M, Li L, Deng X, Dodd E, Batac F, Delwart E.**
25 2017. ENDEMIC INFECTION OF STRANDED SOUTHERN SEA OTTERS
26 (ENHYDRA LUTRIS NEREIS) WITH NOVEL PARVOVIRUS,
27 POLYOMAVIRUS, AND ADENOVIRUS. *J Wildl Dis* **53**:532-542.
- 28 6. **Vaisanen E, Fu Y, Hedman K, Soderlund-Venermo M.** 2017. Human
29 Protoparvoviruses. *Viruses* **9**.
- 30 7. **Geoghegan JL, Pirotta V, Harvey E, Smith A, Buchmann JP, Ostrowski**
31 **M, Eden JS, Harcourt R, Holmes EC.** 2018. Virological Sampling of
32 Inaccessible Wildlife with Drones. *Viruses* **10**.
- 33 8. **de Souza WM, Dennis T, Fumagalli MJ, Araujo J, Sabino-Santos G, Maia**
34 **FGM, Acrani GO, Carrasco AOT, Romeiro MF, Modha S, Vieira LC,**
35 **Ometto T, Queiroz LH, Durigon EL, Nunes MRT, Figueiredo LTM, Gifford**
36 **RJ.** 2018. Novel Parvoviruses from Wild and Domestic Animals in Brazil
37 Provide New Insights into Parvovirus Distribution and Diversity. *Viruses* **10**.
- 38 9. **de Souza WM, Romeiro MF, Fumagalli MJ, Modha S, de Araujo J,**
39 **Queiroz LH, Durigon EL, Figueiredo LT, Murcia PR, Gifford RJ.** 2017.
40 Chapparvoviruses occur in at least three vertebrate classes and have a broad
41 biogeographic distribution. *J Gen Virol* **98**:225-229.
- 42 10. **Phan TG, Gulland F, Simeone C, Deng X, Delwart E.** 2015. Sesavirus:
43 prototype of a new parvovirus genus in feces of a sea lion. *Virus Genes*
44 **50**:134-136.
- 45 11. **Phan TG, Dreno B, da Costa AC, Li L, Orlandi P, Deng X, Kapusinszky B,**
46 **Siqueira J, Knol AC, Halary F, Dantal J, Alexander KA, Pesavento PA,**
47 **Delwart E.** 2016. A new protoparvovirus in human fecal samples and
48 cutaneous T cell lymphomas (mycosis fungoides). *Virology* **496**:299-305.
- 49 12. **Belyi VA, Levine AJ, Skalka AM.** 2010. Sequences from ancestral single-
50 stranded DNA viruses in vertebrate genomes: the parvoviridae and
51 circoviridae are more than 40 to 50 million years old. *J Virol* **84**:12458-12462.
- 52 13. **Katzourakis A, Gifford RJ.** 2010. Endogenous viral elements in animal
53 genomes. *PLoS Genet* **6**:e1001191.

- 1 14. **Liu H, Fu Y, Xie J, Cheng J, Ghabrial SA, Li G, Peng Y, Yi X, Jiang D.** 2011. Widespread endogenization of densoviruses and parvoviruses in
2 animal and human genomes. *J Virol* **85**:9863-9876.
- 3 15. **Reuter G, Boros A, Delwart E, Pankovics P.** 2014. Novel circular single-
4 stranded DNA virus from turkey faeces. *Arch Virol* **159**:2161-2164.
- 5 16. **Yang S, Liu Z, Wang Y, Li W, Fu X, Lin Y, Shen Q, Wang X, Wang H,**
6 **Zhang W.** 2016. A novel rodent Chapparvovirus in feces of wild rats. *Viol J*
7 **13**:133.
- 8 17. **Kapoor A, Simmonds P, Lipkin WI.** 2010. Discovery and characterization of
9 mammalian endogenous parvoviruses. *J Virol* **84**:12628-12635.
- 10 18. **Holmes EC.** 2011. The evolution of endogenous viral elements. *Cell Host*
11 *Microbe* **10**:368-377.
- 12 19. **Baker KS, Leggett RM, Bexfield NH, Alston M, Daly G, Todd S,**
13 **Tachedjian M, Holmes CE, Crameri S, Wang LF, Heeney JL, Suu-Ire R,**
14 **Kellam P, Cunningham AA, Wood JL, Caccamo M, Murcia PR.** 2013.
15 Metagenomic study of the viruses of African straw-coloured fruit bats:
16 detection of a chiropteran poxvirus and isolation of a novel adenovirus.
17 *Virology* **441**:95-106.
- 18 20. **Palinski RM, Mitra N, Hause BM.** 2016. Discovery of a novel Parvovirinae
19 virus, porcine parvovirus 7, by metagenomic sequencing of porcine rectal
20 swabs. *Virus Genes* **52**:564-567.
- 21 21. **Roediger B, Lee Q, Tikoo S, Cobbin JCA, Henderson JM, Jormakka M,**
22 **O'Rourke MB, Padula MP, Pinello N, Henry M, Wynne M, Santagostino**
23 **SF, Brayton CF, Rasmussen L, Lisowski L, Tay SS, Harris DC, Bertram**
24 **JF, Dowling JP, Bertolino P, Lai JH, Wu W, Bachovchin WW, Wong JJ,**
25 **Gorrell MD, Shaban B, Holmes EC, Jolly CJ, Monette S, Weninger W.**
26 2018. An Atypical Parvovirus Drives Chronic Tubulointerstitial Nephropathy
27 and Kidney Fibrosis. *Cell* **175**:530-543.e524.
- 28 22. **Schoonvaere K, Smaghe G, Francis F, de Graaf DC.** 2018. Study of the
29 Metatranscriptome of Eight Social and Solitary Wild Bee Species Reveals
30 Novel Viruses and Bee Parasites. *Front Microbiol* **9**:177.
- 31 23. **Chapman MS, Agbandje-McKenna M.** 2006. Atomic structure of viral
32 particles, p 107–123. *In* Kerr JR, Cotmore SF, Bloom ME, Linden RM, Parrish
33 CR (ed), *Parvoviruses*. Hodder Arnold, Ltd;.
- 34 24. **Kailasan S, Agbandje-McKenna M, Parrish CR.** 2015. Parvovirus Family
35 Conundrum: What Makes a Killer? *Annu Rev Virol* **2**:425-450.
- 36 25. **Penzes JJ, Marsile-Medun S, Agbandje-McKenna M, Gifford RJ.** 2018.
37 Endogenous amdoparvovirus-related elements reveal insights into the biology
38 and evolution of vertebrate parvoviruses. *Virus Evol* **4**:vey026.
- 39 26. **Arriagada G, Gifford RJ.** 2014. Parvovirus-derived endogenous viral
40 elements in two South American rodent genomes. *J Virol* **88**:12158-12162.
- 41 27. **Deyle DR, Russell DW.** 2009. Adeno-associated virus vector integration.
42 *Curr Opin Mol Ther* **11**:442-447.
- 43 28. **Majumder K, Etingov I, Pintel DJ.** 2017. Protoparvovirus Interactions with
44 the Cellular DNA Damage Response. *Viruses* **9**.
- 45 29. **Willmann R.** 2004. Phylogenetic relationships and evolution of insects., p
46 330–344. *In* Cracraft J, Donoghue MJ (ed), *Assembling the Tree of Life*.
47 Oxford University Press.
- 48 30. **Cohen S, Pante N.** 2005. Pushing the envelope: microinjection of Minute
49 virus of mice into *Xenopus* oocytes causes damage to the nuclear envelope.
50 *J Gen Virol* **86**:3243-3252.
- 51 31. **Sonntag F, Bleker S, Leuchs B, Fischer R, Kleinschmidt JA.** 2006.
52 Adeno-associated virus type 2 capsids with externalized VP1/VP2 trafficking
53

- 1 domains are generated prior to passage through the cytoplasm and are
2 maintained until uncoating occurs in the nucleus. *J Virol* **80**:11040-11054.
- 3 32. **Kailasan S, Halder S, Gurda B, Bladek H, Chipman PR, McKenna R,**
4 **Brown K, Agbandje-McKenna M.** 2015. Structure of an enteric pathogen,
5 bovine parvovirus. *J Virol* **89**:2603-2614.
- 6 33. **Kaufmann B, El-Far M, Plevka P, Bowman VD, Li Y, Tijssen P, Rossmann**
7 **MG.** 2011. Structure of Bombyx mori densovirus 1, a silkworm pathogen. *J*
8 *Virol* **85**:4691-4697.
- 9 34. **Williams SH, Che X, Garcia JA, Klena JD, Lee B, Muller D, Ulrich W,**
10 **Corrigan RM, Nichol S, Jain K, Lipkin WI.** 2018. Viral Diversity of House
11 Mice in New York City. *MBio* **9**.
- 12 35. **Drouin LM, Lins B, Janssen M, Bennett A, Chipman P, McKenna R, Chen**
13 **W, Muzyczka N, Cardone G, Baker TS, Agbandje-McKenna M.** 2016.
14 Cryo-electron Microscopy Reconstruction and Stability Studies of the Wild
15 Type and the R432A Variant of Adeno-associated Virus Type 2 Reveal that
16 Capsid Structural Stability Is a Major Factor in Genome Packaging. *J Virol*
17 **90**:8542-8551.
- 18 36. **Simpson AA, Hebert B, Sullivan GM, Parrish CR, Zadori Z, Tijssen P,**
19 **Rossmann MG.** 2002. The structure of porcine parvovirus: comparison with
20 related viruses. *J Mol Biol* **315**:1189-1198.
- 21 37. **Xie Q, Bu W, Bhatia S, Hare J, Somasundaram T, Azzi A, Chapman MS.**
22 2002. The atomic structure of adeno-associated virus (AAV-2), a vector for
23 human gene therapy. *Proc Natl Acad Sci U S A* **99**:10405-10410.
- 24 38. **Simpson AA, Chipman PR, Baker TS, Tijssen P, Rossmann MG.** 1998.
25 The structure of an insect parvovirus (*Galleria mellonella* densovirus) at 3.7 Å
26 resolution. *Structure* **6**:1355-1367.
- 27 39. **Kaufmann B, Bowman VD, Li Y, Szelei J, Waddell PJ, Tijssen P,**
28 **Rossmann MG.** 2010. Structure of *Penaeus stylirostris* densovirus, a shrimp
29 pathogen. *J Virol* **84**:11289-11296.
- 30 40. **Zhu H, Dennis T, Hughes J, Gifford RJ.** 2018. Database-integrated genome
31 screening (DIGS): exploring genomes heuristically using sequence similarity
32 search tools and a relational database. bioRxiv doi:10.1101/246835.
- 33 41. **Carver T, Harris SR, Berriman M, Parkhill J, McQuillan JA.** 2012. Artemis:
34 an integrated platform for visualization and analysis of high-throughput
35 sequence-based experimental data. *Bioinformatics* **28**:464-469.
- 36 42. **Knudsen S.** 1999. Promoter2.0: for the recognition of PolIII promoter
37 sequences. *Bioinformatics* **15**:356-361.
- 38 43. **Reese MG, Eeckman FH, Kulp D, Haussler D.** 1997. Improved splice site
39 detection in Genie. *J Comput Biol* **4**:311-323.
- 40 44. **Dogan RI, Getoor L, Fau - Wilbur WJ, Wilbur Wj Fau - Mount SM, Mount**
41 **SM.** SplicePort--an interactive splice-site analysis tool.
- 42 45. **Salamov AA, Solovyev VV.** 1997. Recognition of 3' -processing sites of
43 human mRNA precursors. *Bioinformatics* **13**:23-28.
- 44 46. **Armougom F, Moretti S, Poirot O, Audic S, Dumas P, Schaeli B, Keduas**
45 **V, Notredame C.** 2006. Espresso: automatic incorporation of structural
46 information in multiple sequence alignments using 3D-Coffee. *Nucleic Acids*
47 *Res* **34**:W604-608.
- 48 47. **Edgar RC.** 2004. MUSCLE: multiple sequence alignment with high accuracy
49 and high throughput. *Nucleic Acids Res* **32**:1792-1797.
- 50 48. **Wallace IM, O'Sullivan O, Higgins DG, Notredame C.** 2006. M-Coffee:
51 combining multiple sequence alignment methods with T-Coffee. *Nucleic Acids*
52 *Res* **34**:1692-1699.

- 1 49. **Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O.**
2 2010. New algorithms and methods to estimate maximum-likelihood
3 phylogenies: assessing the performance of PhyML 3.0. *Syst Biol* **59**:307-321.
4 50. **Felsenstein J.** 2005. PHYLIP (Phylogeny Inference Package) version 3.6. ,
5 Department of Genome Sciences, University of Washington, Seattle.
6 51. **Lobley A, Sadowski MI, Jones DT.** 2009. pGenTHREADER and
7 pDomTHREADER: new methods for improved protein fold recognition and
8 superfamily discrimination. *Bioinformatics* **25**:1761-1767.
9 52. **Yang J, Yan R, Roy A, Xu D, Poisson J, Zhang Y.** 2015. The I-TASSER
10 Suite: protein structure and function prediction. *Nat Methods* **12**:7-8.
11 53. **Carrillo-Tripp M, Shepherd CM, Borelli IA, Venkataraman S, Lander G,**
12 **Natarajan P, Johnson JE, Brooks CL, 3rd, Reddy VS.** 2009. VIPERdb2: an
13 enhanced and web API enabled relational database for structural virology.
14 *Nucleic Acids Res* **37**:D436-442.
15 54. **Schrödinger L.** The PyMOL Molecular Graphics System, Version 2.0.
16 55. **Pettersen EF, Goddard TD, Huang CC, Couch GS, Greenblatt DM, Meng**
17 **EC, Ferrin TE.** 2004. UCSF Chimera--a visualization system for exploratory
18 research and analysis. *J Comput Chem* **25**:1605-1612.
19 56. **Zuker M.** 2003. Mfold web server for nucleic acid folding and hybridization
20 prediction. *Nucleic Acids Res* **31**:3406-3415.
21 57. **Gifford RJ, Blomberg J, Coffin JM, Fan H, Heidmann T, Mayer J, Stoye J,**
22 **Tristem M, Johnson WE.** 2018. Nomenclature for endogenous retrovirus
23 (ERV) loci. *Retrovirology* **15**:59.
24
25

1 **Supporting Information Legends**

2

3 **Figure S1**

4 Secondary structure predictions of the *Syngnathus scovelli* chapparovirus
5 genome termini.

6

7 **Figure S2**

8 (A) Secondary structure and disordered region predictions of the
9 nucleoprotein (NP) ORF from an amniote chapparovirus (mouse kidney
10 virus) and an endogenous chapparoviral element from an invertebrate
11 genome (ChPV.3-Latrodectus hesperus). (B) Predictions of potential
12 phosphorylation sites in case of the same amniote chapparovirus and an
13 invertebrate endogenous chapparoviral element NPs.

14

15 **Figure S3**

16 (A) Alignment of amniote chapparovirus capsid protein ORF (VP) derived
17 amino acid sequences, containing both isolates of murine origin. Variable
18 regions (VRs) are marked by the black bars and coloring is based on
19 sequence similarity (red = highly similar, blue = not similar). The conserved
20 loops making up the fivefold symmetry axes of the capsid are highlighted in
21 bold. (B) The same alignment incorporating the complete VP protein
22 sequences of all novel chapparovirus sequences reported in this study.

23

24

Figure Legends

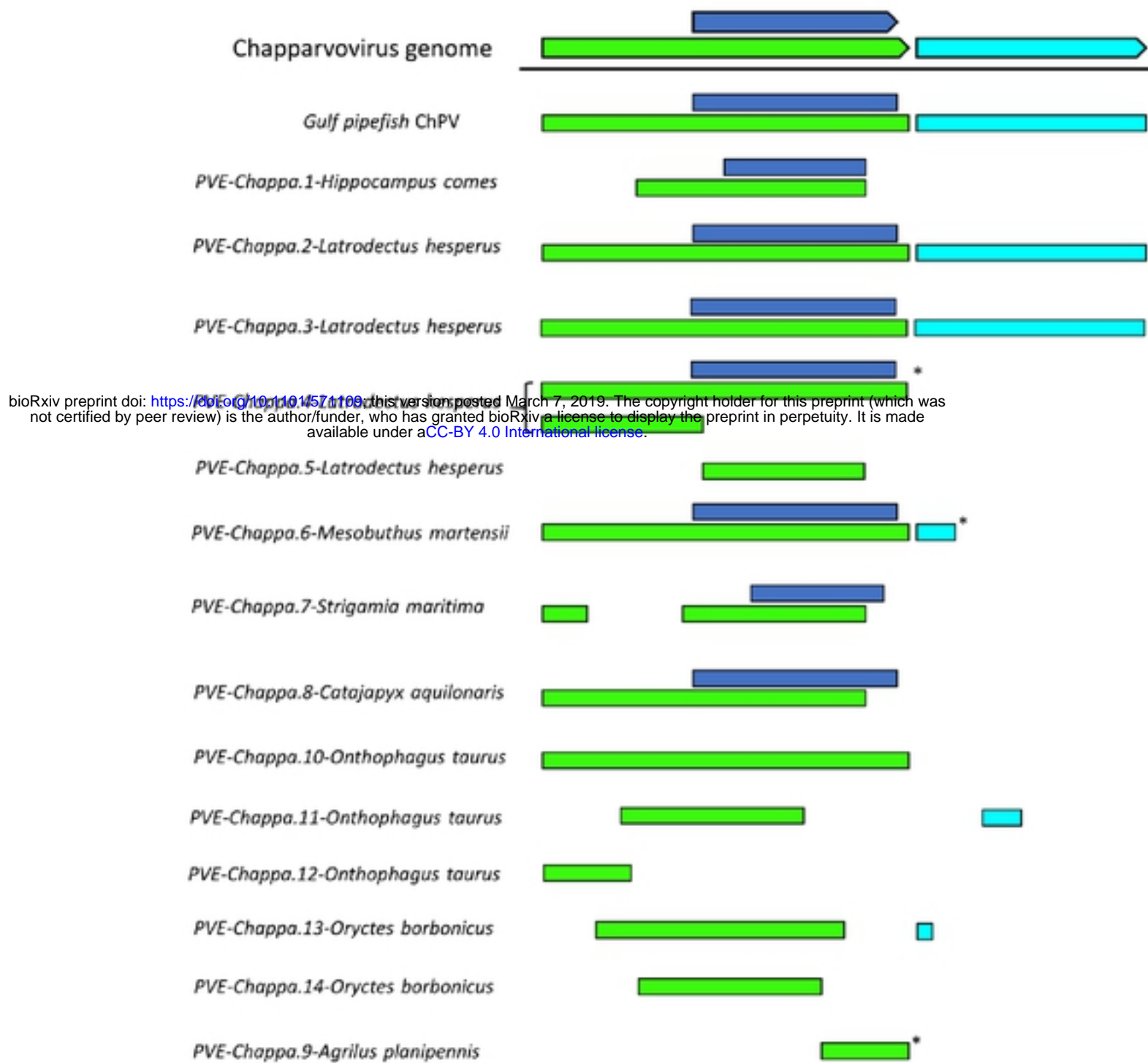
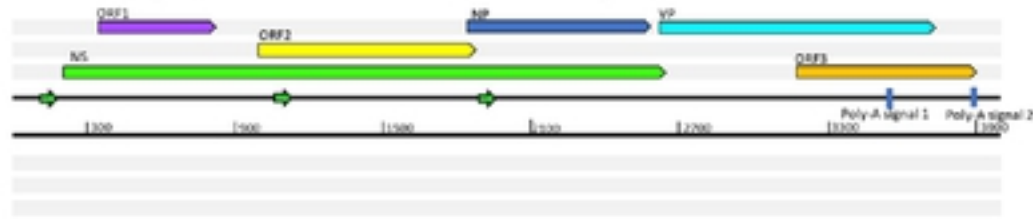


Figure 1. Basic gene content of novel chapparvoviruses and chapparvovirus EPV, shown in relation to a representative chapparvovirus genome (mouse kidney parvovirus). Abbreviations: non-structural protein (NS); capsid protein (VP); nucleoprotein (NP).

Figure 2a

Vertebrates

Gulf pipefish (*Syngnathus scovelli*) chapparvovirus (exogenous)



Invertebrates

ChPV.2-*Latrodectus hesperus*



bioRxiv preprint doi: <https://doi.org/10.1101/571109>; this version posted March 7, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

ChPV.3-*Latrodectus hesperus*



Taurus scarab (*Onthophagus taurus*) ambidensoviral PVE and ChPV.10-*Onthophagus taurus*

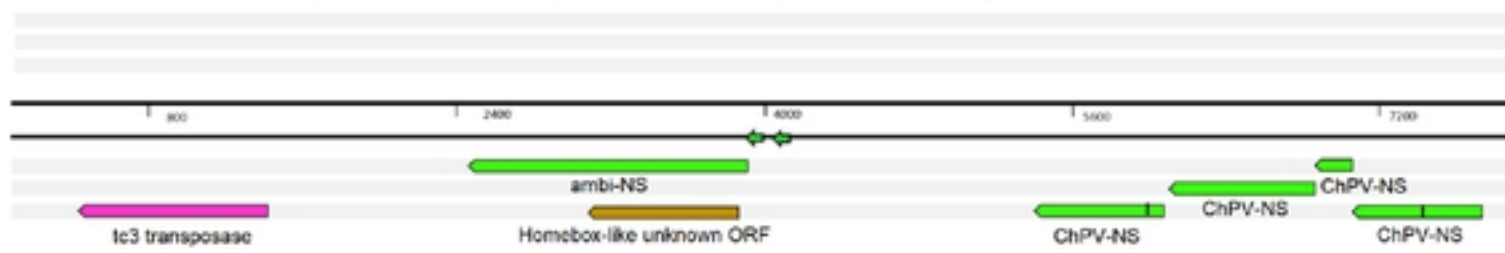
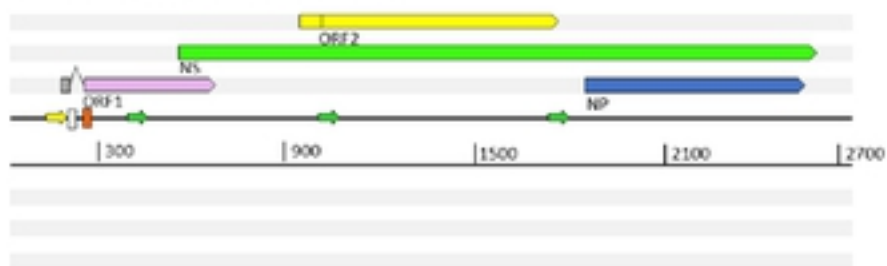
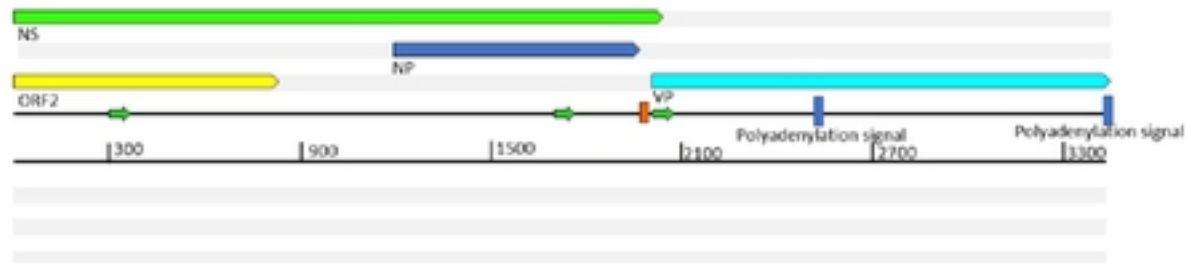


Figure 2b

Simian parvo-like virus 3

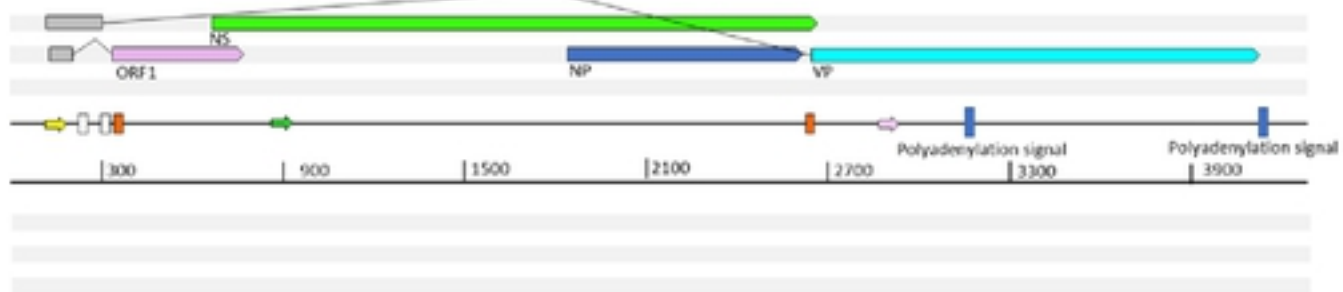


Porcine parvovirus 7



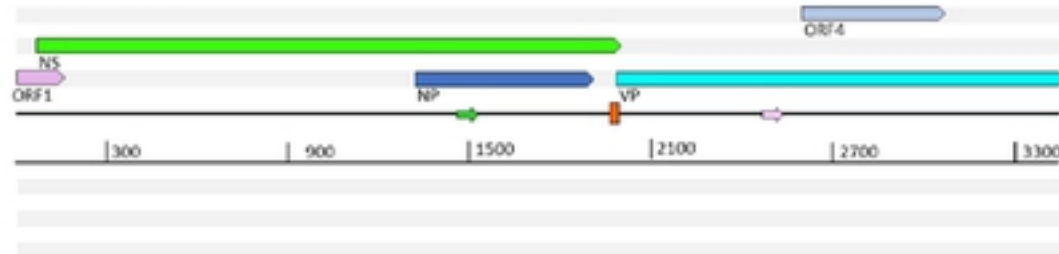
Type 1

bioRxiv preprint doi: <https://doi.org/10.1101/571109>; this version posted March 7, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.



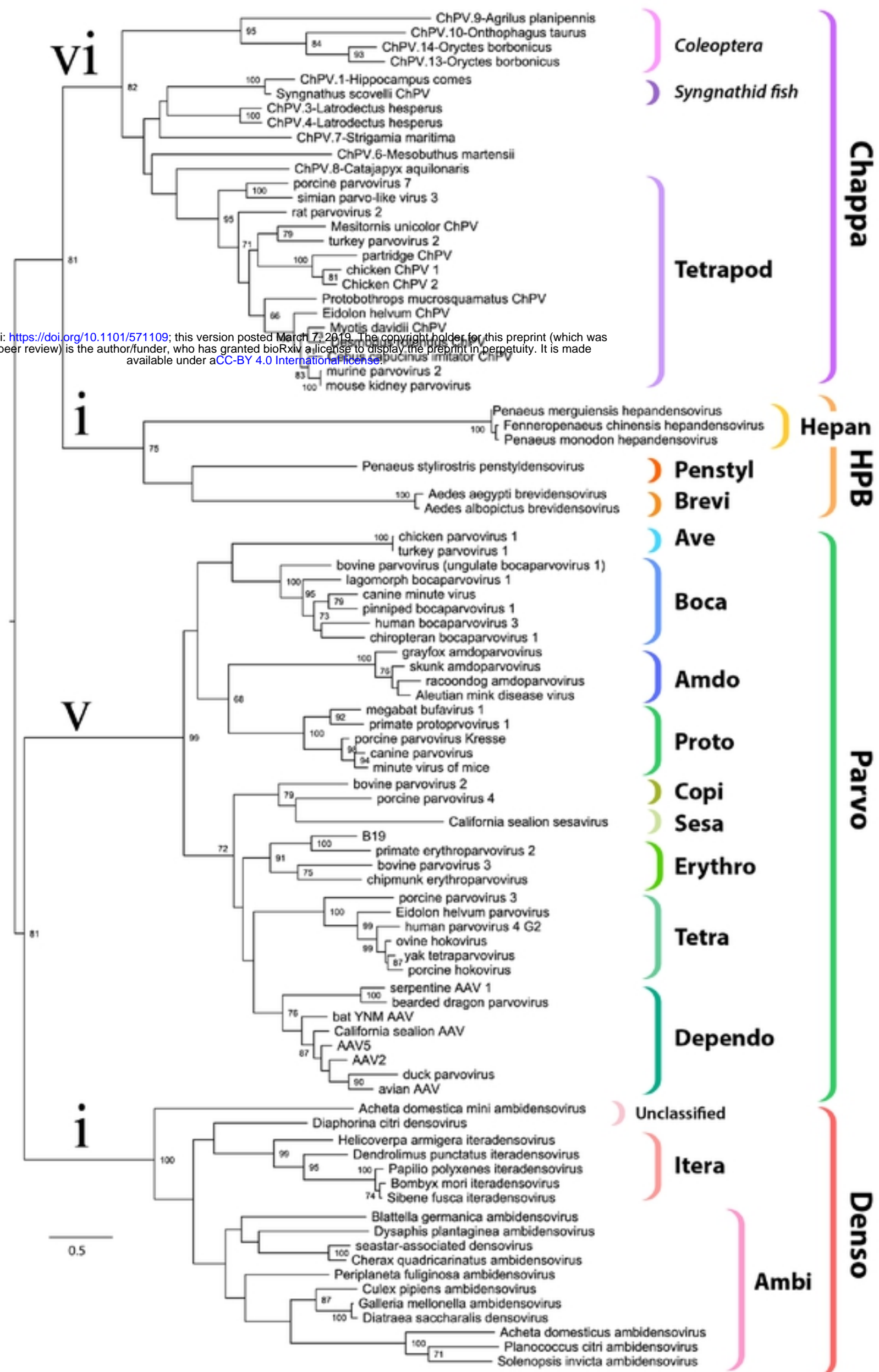
Type 2

Turkey parvovirus 2



Type 3

Fig.2 legend: Genome organization, position of open reading frames (ORFs) and predicted cis transcription elements of chapparvoviruses and endogenous chapparvoviral elements shown in six reading frames. ORFs are represented by arrows, colored according to homology. In-frame stop codons are shown as vertical lines. Splice donor sites are marked by white-colored, acceptor sites by orange-colored bars. Blue colored bars show predicted polyadenylation signals. Promoters are presented as small green arrows if predicted with a higher score than 0.95 of 1 and pink if the score is between 0.9 and 0.95. Sequences marked by grey boxes are assessed to be transcribed but not translated. (A) Near complete entries detected include the exogenous, potentially circulating chapparvovirus of the gulf pipefish (*Syngnathus scovelli*) as well as two, non-identical, endogenous chapparvoviruses of the Western black widow spider (*Latrodectus hesperus*). ChPV.2-*Latrodectus hesperus* contains a previously unidentified repetitive element, present as multiple copies scattered in the *Latrodectus* genome, marked by the white box within the VP gene. The element ChPV.10-*Onthophagus taurus* shares its integration site with another endogenous element disclosing similarity to ambidensoviruses. (B) Representatives of the three basic genome organization types of exogenous amniote chapparvoviruses. ORF4, however, does not occur in any of the other avian chapparvoviruses.



bioRxiv preprint doi: <https://doi.org/10.1101/571109>; this version posted March 7, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

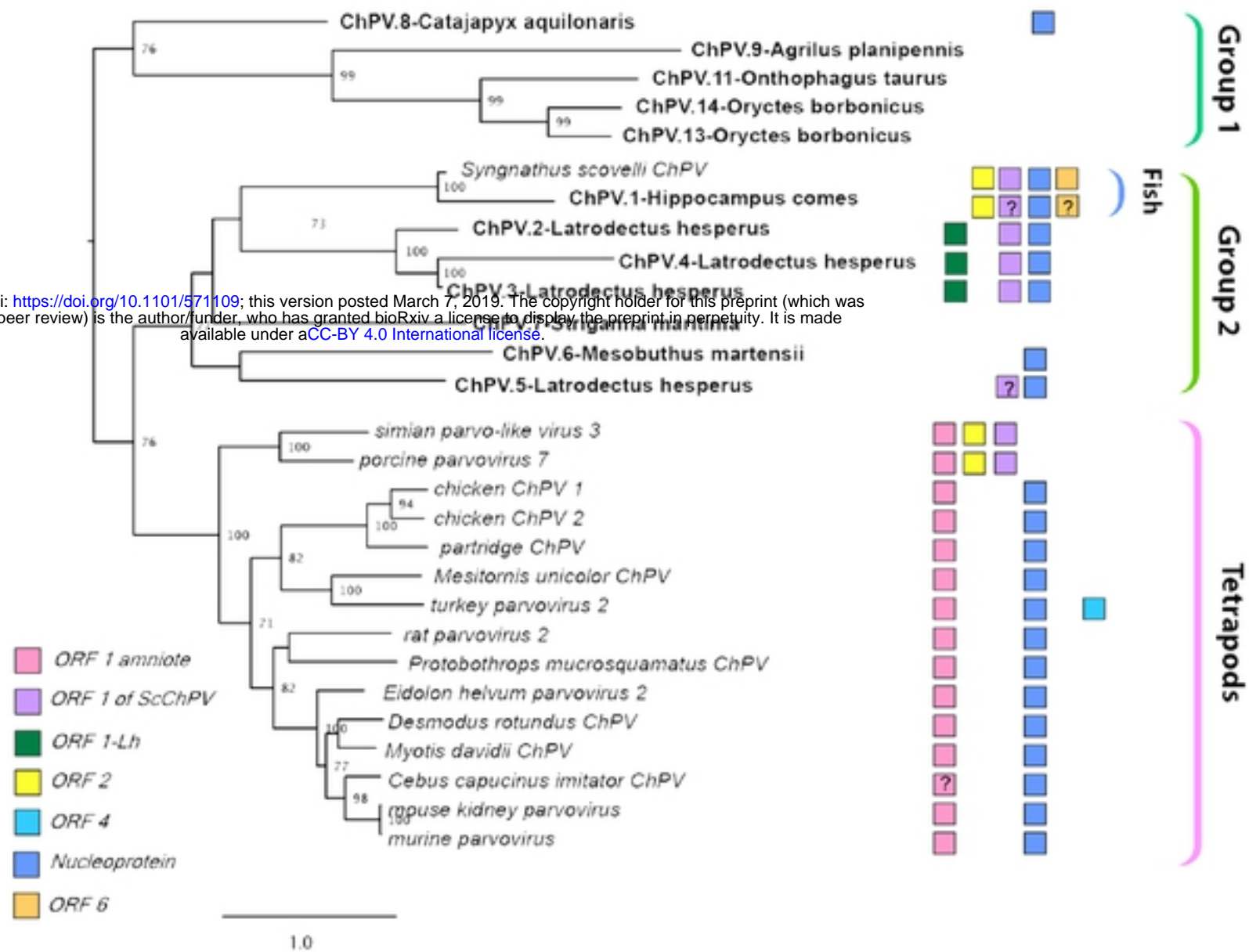


Figure 3b

Figure 3 legend: (A) Maximum likelihood phylogeny of the family *Parvoviridae* based on the tripartite helicase domain of the large non-structural protein NS1. (B) Maximum likelihood phylogenetic reconstructions of the Chapparvovirus clade based on the complete aligned amino acid sequences of the NS1. The presence or suspected presence (marked by a question mark) of predicted auxiliary protein encoding open reading frames (ORFs) are mapped as boxes of various colors next to each entry. Sequences of vertebrate origin are indicated by black, from members of the arthropod subphyla Chelicerata by blue, Myriapoda by purple and Hexapoda by green text. Bold taxa labels indicate endogenous sequences, whereas taxa labels in italics indicate sequences known or believed to derive from viruses.

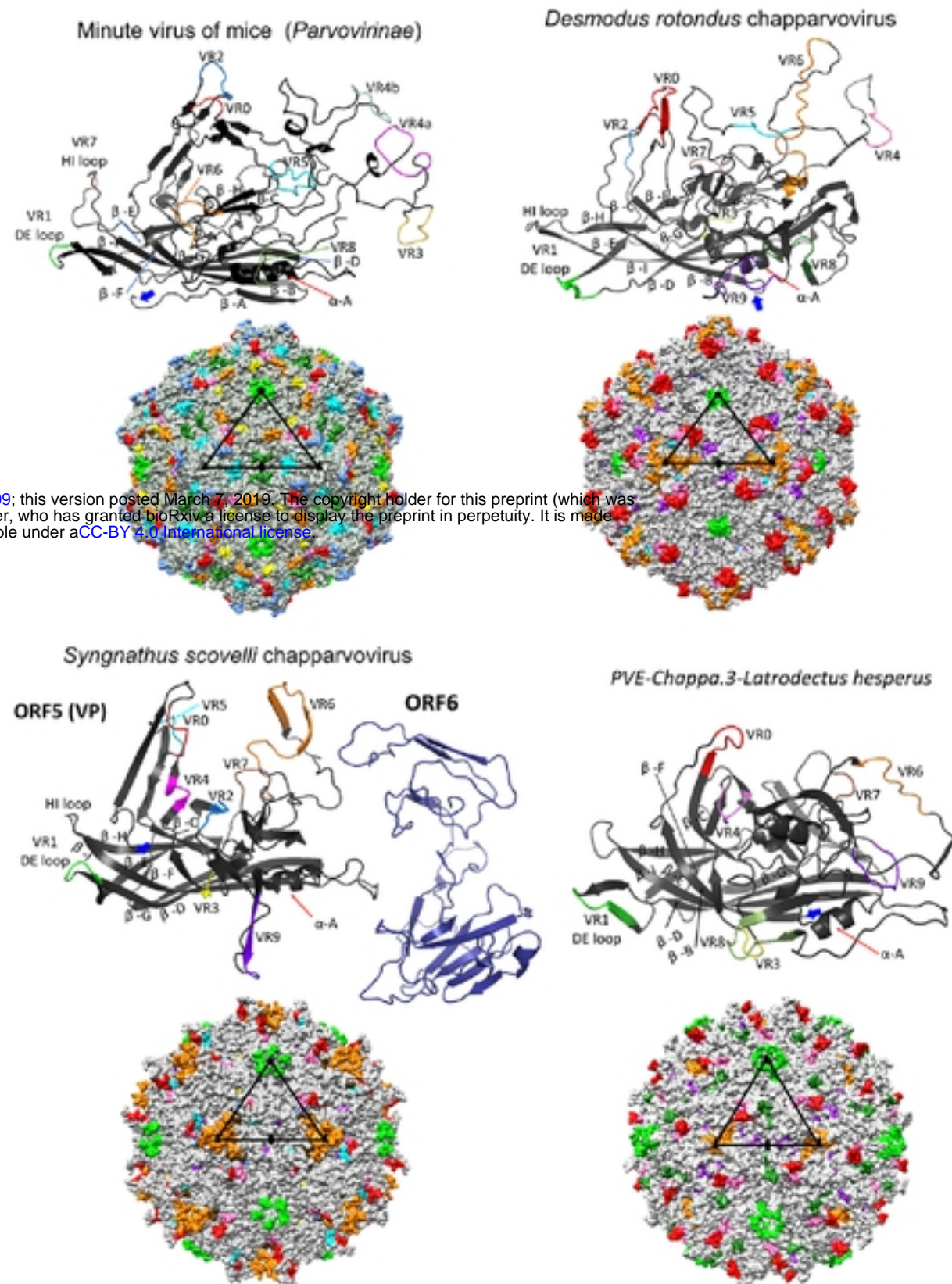


Figure 4

Comparison of VP monomer ribbon diagrams of the protoparvovirus minute virus of mice (PDB ID: 1Z14) from subfamily *Parvovirinae* to homology models of an amniote, a fish and an endogenous arthropod chapparovirus. Variable regions (VRs) of the same number are marked by the same color and mapped to the surface area of the icosahedral capsid model constructed of 60 monomers. In case of the novel *Syngnathus scovelli* chapparovirus the hypothetical structural protein product of ORF6 was also modelled. Triangles mark the position of a single monomer within the capsid, the fivefold

symmetry axis marked by a pentagon, the threefolds with triangles and the twofold with an ellipsoid.

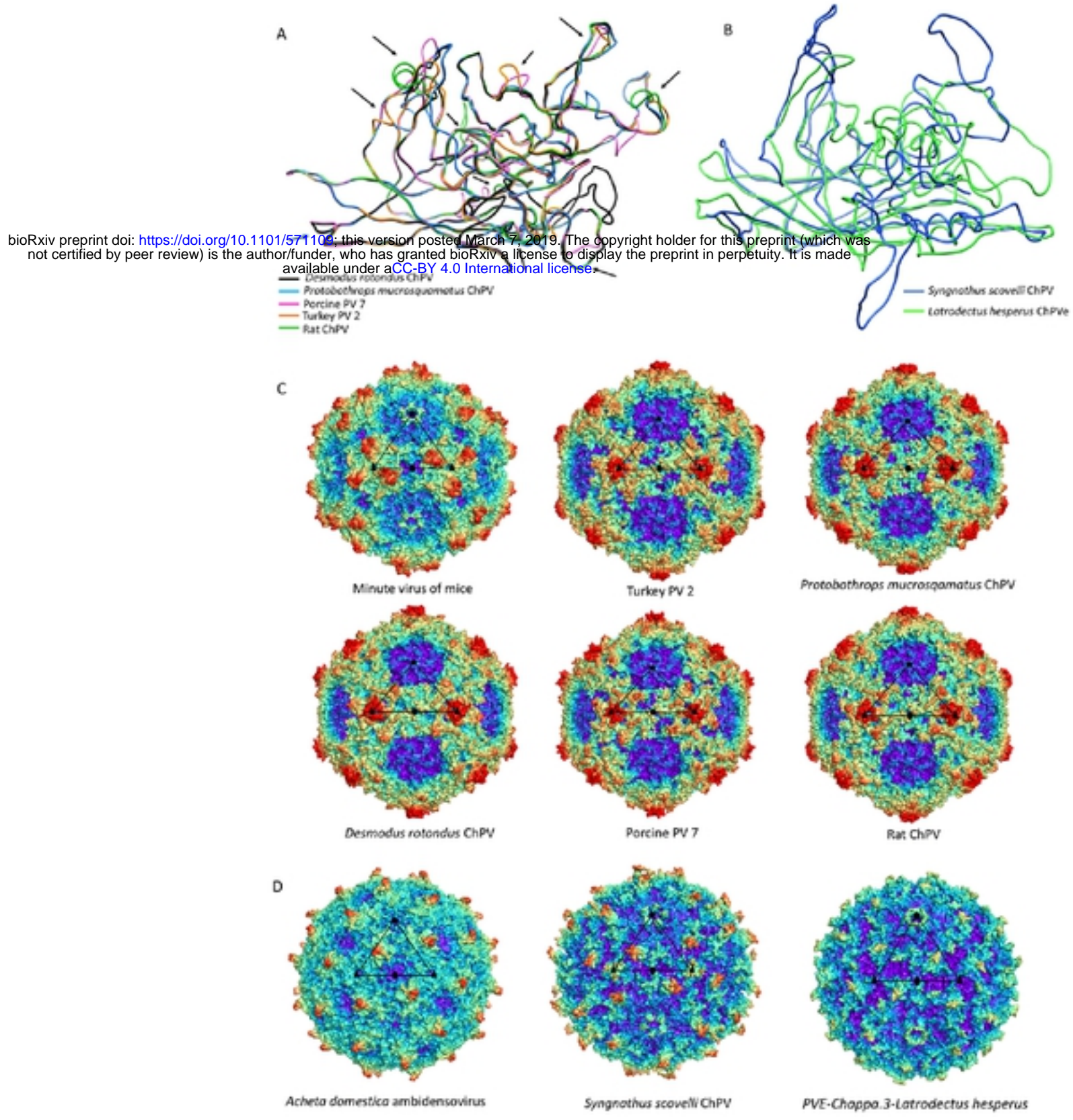


Figure 5

(A) Superposition of protein monomer homology models of amniote chapparovirus capsids, including reptilian, avian, rodent, chiropteran and ungulate representatives. Black arrows show variable regions (VRs) (B) Superposition of capsid protein monomer homology model of the novel exogenous fish chapparovirus of *Syngnathus scovelli* with that of an endogenous arthropod chapparovirus element. (C) Capsid surface morphology of amniote chapparovirus homology models compared to that obtained for a prototypic parvovirus, minute virus of mice (MVM) (PDB ID: 1Z14 at 3.25 Å resolution). Triangles mark the position of a single monomer within the capsid, the fivefold symmetry axis marked by a pentagon, the threefolds with triangles and the twofold with an ellipsoid. (D) Homology models of complete viral capsids to compare surface morphology of the newly-identified fish chapparovirus and arachnid endogenous chapparoviral element with that of the actual capsid structure of a densovirus of subfamily Densovirinae, genus Ambidensovirus (PDB ID: 4MGU at 3.5 Å resolution).