

UMass Chan Medical School

eScholarship@UMassChan

University of Massachusetts Medical School Faculty Publications

2021-10-26

Comparative analysis reveals the long-term co-evolutionary history of parvoviruses and vertebrates [preprint]

Matthew A. Campbell
University of Alaska

Et al.

Let us know how access to this document benefits you.

Follow this and additional works at: https://escholarship.umassmed.edu/faculty_pubs



Part of the [Ecology and Evolutionary Biology Commons](#)

Repository Citation

Campbell MA, Loncar S, Kotin RM, Gifford RJ. (2021). Comparative analysis reveals the long-term co-evolutionary history of parvoviruses and vertebrates [preprint]. University of Massachusetts Medical School Faculty Publications. <https://doi.org/10.1101/2021.10.25.465781>. Retrieved from https://escholarship.umassmed.edu/faculty_pubs/2100

Creative Commons License



This work is licensed under a [Creative Commons Attribution-NonCommercial 4.0 International License](#).

This material is brought to you by eScholarship@UMassChan. It has been accepted for inclusion in University of Massachusetts Medical School Faculty Publications by an authorized administrator of eScholarship@UMassChan. For more information, please contact Lisa.Palmer@umassmed.edu.

1 Long title (250 characters):

2 **Comparative analysis reveals the long-term co-evolutionary history of parvoviruses and**
3 **vertebrates.**

4
5 Short title (70 characters):

6 **Origin and evolution of vertebrate parvoviruses**

7
8 Matthew A. Campbell*, Shannon Loncar*, Robert Kotin and Robert J. Gifford

9
10 *Equal contributions

11
12
13 **Matthew A. Campbell:** *University of Alaska Museum of the North, Fishes and Marine*
14 *Invertebrates, 1962 Yukon Drive, Fairbanks, AK 99775 USA*

15
16 **Shannon Loncar:** *University of Massachusetts Medical School, Department of Microbiology*
17 *and Physiological Systems, Gene Therapy Center, 55 Lake Ave. North, Worcester, MA 01655,*
18 *USA*

19 *Current address: Sanofi Genzyme*

20
21 **Robert Kotin:** *University of Massachusetts Medical School, Department of Microbiology and*
22 *Physiological Systems, Gene Therapy Center, 55 Lake Ave. North, Worcester, MA 01655, USA*
23 *And Synteny Therapeutics, Cambridge, MA 02138*

24
25 **Robert J. Gifford:** *MRC-University of Glasgow Centre for Virus Research, 464 Bearsden Rd,*
26 *Bearsden, Glasgow, UK, G61 1QH*

27
28 **Corresponding author:** Robert J. Gifford: robert.gifford@glasgow.ac.uk

29
30

1 SUMMARY

2 Parvoviruses (family *Parvoviridae*) are small, non-enveloped DNA viruses that infect a
3 broad range of animal species. Comparative studies, supported by experimental evidence,
4 show that many vertebrate species contain sequences derived from ancient parvoviruses
5 embedded in their genomes. These 'endogenous parvoviral elements' (EPVs), which arose via
6 recombination-based mechanisms in infected germline cells of ancestral organisms, constitute a
7 form of 'molecular fossil record' that can be used to investigate the origin and evolution of the
8 parvovirus family. Here, we use comparative approaches to investigate 198 EPV loci,
9 represented by 470 EPV sequences identified in a comprehensive *in silico* screen of 752
10 published vertebrate genomes. We investigated EPV loci by constructing an open resource that
11 contains all of the data items required for comparative sequence analysis of parvoviruses and
12 uses a relational database to represent the complex semantic relationships between them. We
13 used this standardised framework to implement reproducible comparative phylogenetic analysis
14 of combined EPV and virus data. Our analysis reveals that viruses closely related to
15 contemporary parvoviruses have circulated among vertebrates since the Late Cretaceous
16 epoch (100-66 million years ago). We present evidence that the subfamily *Parvovirinae*, which
17 includes ten vertebrate-specific genera, has evolved in broad congruence with the emergence
18 and diversification of major vertebrate groups. Furthermore, we infer defining aspects of
19 evolution within individual parvovirus genera - mammalian vicariance for protoparvoviruses
20 (genus *Protoparvovirus*), and inter-class transmission for dependoparvoviruses (genus
21 *Dependoparvovirus*) - thereby establishing an ecological and evolutionary perspective through
22 which to approach analysis of these virus groups. We also identify evidence of EPV expression
23 at RNA level and show that EPV coding sequences have frequently been maintained during
24 evolution, adding to a growing body of evidence that EPV loci have been co-opted or exapted
25 by vertebrate species, and especially by mammals. Our findings offer fundamental insights into
26 parvovirus evolution. In addition, we establish novel genomic resources that can advance the
27 development of parvovirus-related research - including both therapeutics and disease
28 prevention efforts - by enabling more efficient dissemination and utilisation of relevant,
29 evolution-related domain knowledge.

1 INTRODUCTION

2 Parvoviruses (family *Parvoviridae*) are a diverse group of small, non-enveloped DNA
3 viruses that infect a broad and phylogenetically diverse range of animal species [1, 2]. The
4 family includes numerous important pathogens of humans and domesticated species, including
5 erythroparvovirus B19 (fifth disease), carnivore protoparvovirus 1 (canine parvovirus) and
6 carnivore amdoparvovirus 1 (Aleutian mink disease). Parvoviruses are also being developed as
7 next-generation therapeutic tools - rodent protoparvoviruses (RoPVs) are promising anticancer
8 agents that show natural oncotropism and oncolytic properties [3, 4], while adeno-associated
9 virus (AAV), a non-autonomously replicating dependoparvovirus, has been successfully adapted
10 as a gene therapy vector, and parvoviruses are leading candidates for the further development
11 of human gene therapy [5, 6].

12 Parvoviruses have highly robust, icosahedral capsids (T=1) that contain a linear, single-
13 stranded DNA genome typically ~5 kilobases (kb) in length. Parvovirus genomes are typically
14 very compact and generally exhibit the same basic genetic organization comprising two major
15 gene cassettes, one (Rep/NS) that encodes the non-structural proteins, and another (Cap/VP)
16 that encodes the structural coat proteins of the virion [2]. However, some genera contain
17 additional open reading frames (ORFs) adjacent to these genes or overlapping them in
18 alternative reading frames. The genome is flanked at the 3' and 5' ends by palindromic inverted
19 terminal repeat (ITR) sequences that are the only *cis* elements required for replication.

20 Recent years have seen many important advances in understanding of parvovirus
21 evolution and diversity, driven primarily by dramatic increases in the availability of DNA
22 sequence data and investments in deriving novel adeno-associated virus capsid for gene
23 therapy applications. Metagenomic sequencing has enabled the discovery of numerous novel
24 parvovirus species, which in turn has led to the taxonomic re-organization of the *Parvoviridae* to
25 include additional subfamilies and genera [1]. In addition, whole genome DNA sequencing has
26 revealed that DNA sequences derived from parvoviruses are widespread in animal genomes [7-
27 13]. These endogenous parvoviral elements (EPVs) are thought to have arisen when parvovirus
28 infection of germline cells (i.e., gametes, gamete producing cells, or early-stage embryos) led to
29 integration of parvovirus-derived DNA into chromosomal DNA so that it was subsequently
30 inherited as a newly acquired allele. Integration of parvovirus DNA can occur via cell-mediated,
31 non-homologous recombination but may also be mediated by the activities of virus-encoded
32 proteins [14, 15]. Comparative genomic studies have shown that EPV sequences often occur as
33 orthologous loci in multiple related host species, demonstrating that they were incorporated into
34 the germline of a common ancestor. Thus, species divergence times – which are in part based

1 on evidence from the fossil record - provide a robust method of deriving minimum age estimates
2 for EVE insertions. Many EVEs represent virus lineages that have not been described
3 previously and may be extinct [7, 16]. However, others clearly represent members of
4 contemporary virus groups, and age estimates obtained for these EVEs provide insights into
5 their long-term evolutionary history [7, 17-20].

6 The extent to which EPVs have reached fixation through positive selection - as opposed
7 to incidental factors such as founder effects, population bottlenecks, and genetic hitchhiking -
8 remains unclear. Potentially, EPV genetic information might sometimes be co-opted or
9 “exapted” as has been reported for EVEs derived from other virus groups, including retroviruses
10 (family *Retroviridae*) [21, 22] and polydnaviruses (family *Polydnaviridae*) [23]. Recent studies
11 have revealed that two distinct, fixed EPVs in the germline of (i) the degu (*Octodon degus*) and
12 (ii) family Elephantidae (elephants) – both of which encode an intact Rep protein ORF - exhibit
13 similar patterns of tissue-specific expression in the liver [24, 25]. These observations suggest
14 that expression of Rep protein or mRNA might – in some way - be physiologically relevant. More
15 broadly, incorporation of parvovirus-derived DNA into animal germlines may provide a novel
16 DNA substrate for the evolution of new genes – for example, guinea pigs (*Cavia porcellus*)
17 encode a predicted polypeptide gene product comprising a partial myosin9-like (M9I) gene
18 fused to a 3' truncated, EPV-encoded replicase [26].

19 Comparative genomic analysis can reveal key insights into the biology and evolution of
20 viral species. EVE data are critically important components of these studies as they provide
21 calibrations in geologic time. Unfortunately, making effective use of these data is challenging for
22 a variety of reasons. This reflects a general lack of reproducibility and reusability in
23 computational genomics [27], particularly where rapidly evolving and highly divergent
24 sequences are involved [28, 29]. To address these issues we previously developed GLUE
25 (Genes Linked by Underlying Evolution), a sequence data-centric bioinformatics environment
26 computational genomics, with a focus on variation, evolution, and sequence interpretation [30].
27 Here, we used GLUE to create ‘Parvovirus-GLUE’, an extensible, open resource for
28 comparative genomic analysis of parvovirus and EPV sequence data. We catalogue hundreds
29 of EPV sequences in published whole genome sequence (WGS) data using a standardized
30 nomenclature system to systematize the parvovirus fossil record. We capture these data in
31 Parvovirus-GLUE and use reproducible approaches to examine their genomic and phylogenetic
32 characteristics, revealing new insights into parvovirus ecology and evolution.

1 RESULTS

2 Creation of open resources for reproducible genomic analysis of parvoviruses

3 Comparative genomic analyses generally entail the construction of complex data sets
4 comprising molecular sequence data linked to other kinds of information. Usually these include
5 genome feature annotations and multiple sequence alignments (MSAs) as well as other diverse
6 kinds of data. We used GLUE to create Parvovirus-GLUE [31], an open accessible online
7 resource for comparative genomic analysis of parvoviruses and EPVs that preserves the 'state'
8 of our data so that our analyses can be precisely and widely replicated (**Fig. S1a-b**).
9 Furthermore, hosting of the Parvovirus-GLUE project in an openly accessible online version
10 control system (GitHub) provides a platform for ongoing development of this resource by
11 multiple collaborators, following practices established in the software industry (**Fig. S1c**).

12 The Parvovirus-GLUE project incorporates all of the data items required for broad
13 comparative sequence analysis of parvoviruses, including: (i) a set of reference sequences
14 representing all known parvovirus species (**Supplementary**); (ii) sequence and isolate-specific
15 information (e.g. host species, vector species) in tabular form; (iii) a standardized set of
16 parvovirus genome features and their coordinates within selected reference genome
17 sequences; (iv) a set of MSAs incorporating all sequences in the project. Loading the project
18 into the GLUE 'engine' generates a relational database that not only contains the data items
19 associated with our analysis, but also represents the complex semantic links between them
20 (**Fig. S1a**). Reproducible comparative genomic analysis of parvoviruses can be implemented by
21 using GLUE's command layer to coordinate interactions between the Parvovirus-GLUE
22 database and bioinformatics software tools. The resource can be installed on all commonly-
23 used computing platforms, and is also fully containerised via Docker [32]. It can be used as a
24 local, stand-alone tool or as a robust foundation for the development of genome analysis-based
25 reporting tools for potential use in human and animal health (**Fig. S1b**) - e.g. see HCV-GLUE
26 [30], RABV-GLUE [33].

27

28 Systematic recovery of the parvovirus 'fossil record'

29 We screened *in silico* WGS data of 738 vertebrate species and recovered a total of 595
30 EPV sequences. EPV sequences ranged from near full-length virus genomes to fragments
31 ~150-300 nucleotides (nt) in length (**Fig. 1**). We classified EPVs into taxonomic groups based
32 on their phylogenetic relationships to contemporary viruses (**Fig. 2**). Among EPVs that were
33 >200nt in length, the majority (192/198=0.97%) were either unambiguous members of
34 contemporary genera or members of groups that emerge as sister clades to these genera. We

1 used comparative approaches to resolve these sequences into sets of orthologs, revealing that
2 they represent at least 198 distinct germline incorporation events (**Table 1, Tables S1-S6,**
3 **Supplementary**). As part of these efforts, we investigated the genes flanking each putatively
4 novel locus. We identified common sets of genes flanking most EPV loci (**Fig. 2**). We applied
5 unique identifiers to all EPVs identified in our study using a systematic nomenclature system
6 that was originally developed for endogenous retroviruses (ERVs) but has more recently been
7 applied to EVEs [19, 20]. This nomenclature captures information about orthology, enforcing a
8 higher level of order on our data set by unambiguously associating EPV sequences with
9 genomic loci. The semantic links between these data items are recorded in the Parvovirus-
10 GLUE database, so that they are available for retrieval and manipulation in computational
11 analyses.

12 EPVs were identified in the genomes of all major groups of terrestrial vertebrates (**Table**
13 **1**) except agnathans, crocodiles, or amphibians. However, they occur much more frequently in
14 mammals than in other groups. The majority of the EPVs identified in vertebrate WGS data
15 derived from genera within the subfamily *Parvovirinae*, but we also identified rare examples of
16 EPVs derived from subfamily *Hamaparvovirinae* (genus *Ichthamaparvovirus*) (**Table 2, Fig. S5**).
17 The identification of ichthamaparvovirus-derived EPVs in snakes provides the first evidence that
18 the host range of this viral genus extends to reptiles. Furthermore, orthologous copies of this
19 EPV were identified in multiple snake species, providing a minimum age of 62 My for the
20 *Ichthamaparvovirus* genus, and by extension the *Hamaparvovirinae* subfamily. Among
21 *Parvovirinae*-derived EPVs, those derived from proto- and dependoparvoviruses predominate.
22 Other genera (*Erythroparvovirus*, *Amdoparvovirus* [34]) are also represented in the mammalian
23 germline, but are relatively rare (<1% of species examined).

24 Previous studies have shown that some EPV loci express RNA with the potential to
25 encode polypeptide gene products, either as unspliced viral RNA [12, 24, 25], or as fusion
26 genes comprising RNA sequences derived from both host and viral sources [26]. We examined
27 coding potential in EPVs and identified numerous sequences capable of encoding uninterrupted
28 polypeptide sequences of 300 amino acids (aa) or more, with some ranging up to 722aa (**Table**
29 **S7**). Furthermore, screening of RNA databases revealed evidence for expression of EPV RNA
30 in several previously unreported EPVs) (**Supp. Doc 1**).

31

32 *EPVs reveal the deep evolutionary origins of the subfamily Parvovirinae*

33 We performed a comprehensive phylogenetic analysis of the *Parvoviridae*, using all
34 available information, including EPVs. Phylogenies revealed three robustly supported sub-

1 lineages within subfamily *Parvovirinae*, each encompassing multiple genera, as follows: (i)
2 “Amdo-Proto”: *Amdo-* and *Protoparvovirus*; (ii) “Ave-Boca”: *Ave-* and *Bocaparvovirus*; (iii)
3 “ETDC”: *Erythro-*, *Tetra-*, *Dependo-* and *Copiparvovirus* (**Fig. 2**).

4 Next, we compiled estimates of EPV integration dates (**Table S1-S6**) to provide an
5 overview of parvovirus and vertebrate interactions over the past 100 My (**Fig. 3**). This revealed
6 that germline incorporation of parvovirus DNA occurred throughout the Cenozoic Era in a broad
7 range of vertebrate species (**Fig. 3**). Our results reveal that mammals acquired EPVs at a much
8 higher frequency than other vertebrate groups (**Table 2**). In addition to the previously reported
9 dependoparvovirus-derived elements in “whippomorphs” (cetaceans and hippopotamus) [13],
10 lagomorphs, Old World rodents [7], New World rodents [24, 26], elephants [25], and macropoids
11 [12] we identified numerous other ancient EPV loci diverse range of animal species (**Table S1**,
12 **Fig 3**). Orthologous sets of EPV sequences demonstrating were also identified in passerine
13 birds (order Passeriformes), establishing the ancestral presence of these viruses among
14 ancestral members of clade Neoaves >85 Mya [35]. In addition, two ancient EPV loci were
15 identified in snakes - an EPV in pit vipers provides a minimum age of 30 My for the Amdo-Proto
16 lineage.

17 The parvovirus family likely has extremely ancient origins, perhaps dating back to the
18 origin of animal species [11], and the independent formation and fixation of EPVs in such a
19 diverse range of vertebrate groups demonstrates that *Parvovirinae* genera circulated widely
20 among vertebrate fauna throughout the Cenozoic Era. Furthermore, since we know that
21 transmission between distantly related host groups is rare, we can tentatively estimate the age
22 of sublineages within the *Parvovirinae*, based on the assumption that they reflect broad
23 codivergence of viruses and hosts – at least at higher taxonomic levels. For example, the “Ave-
24 Boca” lineage comprises clearly distinct avian (Ave-) and mammalian (Boca-) lineages,
25 suggesting that ancestral members of this virus lineage circulated among the common
26 ancestors of birds and mammals >300 Mya. Furthermore, we identified EPVs in the genomes of
27 cartilaginous and ray-finned fish, suggesting that the subfamily *Parvovirinae* may be as old, if
28 not older, than the vertebrate lineage itself (**Fig. 2**).

29 Interestingly, roseoloviruses (genus *Roseolovirus*) – the group of betaherpesviruses
30 (subfamily *Betaherpesvirinae*) that includes human herpesvirus 6 (HHV6) – have acquired a
31 homolog of the parvovirus *rep* gene, called U94, presumably when infection of the same host
32 cell led to parvovirus DNA being incorporated the ‘germline’ of an ancient betaherpesvirus [36].
33 The presence of U94 in an orthologous position in rodent and bat betaherpesviruses, as well as
34 within a betaherpesvirus-derived EVE in the tarsier genome [37], demonstrates that it arose

1 through a horizontal gene transfer event that occurred ancestrally, likely before the divergence
2 of (i) eutherian mammal orders and (ii) betaherpesvirus genera [38]. Phylogenetic
3 reconstructions show that this gene derives from the EDTC sub-lineage within the subfamily
4 *Parvovirinae* (**Fig. 2**).

5 Parvovirus genomes have palindromic ITR sequences at both the 3' and 5' ends which
6 can fold back on themselves to form "hairpin" structures that are stabilized by intramolecular
7 base-pairing. These "hairpin" structures are critical for genome replication in all parvoviruses,
8 however, whereas they are heterotelomeric (asymmetrical) in some genera (*Amdo*-, *Proto*-,
9 *Boca*-, and *Aveparvovirus*) they are homotelomeric (symmetrical) in others [39]. Interestingly,
10 the distribution of this trait (where it has been described) across sub-lineages within the
11 subfamily *Parvovirinae* suggests that - under the principles of maximum parsimony - the
12 asymmetrical form (which is found across the "Amdo-Proto" and "Ave-Boca" sublineages) would
13 be the ancestral form. Within the "ETDC" lineage, ITRs have only been described for the
14 *Dependoparvovirus* and *Erythroparvovirus* genera, both of which have homotelomeric ITRs (**Fig**
15 **2c**), suggesting this that the presence of homotelomeric ITRs is a derived characteristic in
16 subfamily *Parvovirinae*. Similarly, in all *Parvovirinae* groups except genus *Amdoparvovirus*, the
17 N-terminal region of VP1 (the largest of the capsid) contains a phospholipase A2 (PLA2)
18 enzymatic domain that becomes exposed at the particle surface during cell entry and is required
19 for escape from the endosomal compartments. Phylogenetic reconstructions indicate that loss
20 of PLA2 is an acquired characteristic of amdoparvoviruses (**Fig. 2**).

21 Another variable characteristic found in the *Parvovirinae* is the regulation of gene
22 expression strategies, with members of the *Proto*- and *Dependoparvovirus* genera using two to
23 three separate transcriptional promoters, whereas in the *Amdo*-, *Erythro*-, and *Boca*- genera all
24 genes are expressed from a single promoter and genus-specific read-through mechanisms are
25 used to produce alternative transcripts [2]. The presence of multiple separate promoters in the
26 distantly related *Proto*- and *Dependoparvovirus* genera indicates that this expression strategy is
27 probably ancestral, although the possibility that it evolved convergently in each lineage cannot
28 be formally ruled out.

29

30 *Mammalian vicariance shaped the evolution of protoparvoviruses*

31 We identified 121 protoparvovirus-related EPV sequences in mammals, which we
32 estimate to represent at least 105 distinct germline incorporation events (**Table S1**). Several
33 near full-length genomes were identified, and many elements spanned >50% of the genome
34 (**Fig 1a**). We reconstructed the evolutionary relationships between protoparvovirus-related

1 EPVs and contemporary protoparvoviruses, revealing three major subclades within the
2 *Protoparvovirus* genus, which we labelled “Archaeo-”, “Meso-” and “Neo-” protoparvovirus) (**Fig**
3 **4a**). The Archaeoprotoparvovirus (ApPV) clade is comprised exclusively of EPVs and is highly
4 represented in the genomes of Australian marsupials (Australidelphia), American marsupials
5 (Ameridelphia) and New World rodents. It includes numerous elements that are near full-length,
6 but none encoding intact open reading frames (ORFs).

7 The Mesoprotoparvovirus (MpPV) clade is also comprised exclusively of EPVs, and was
8 sparsely represented in the EPV fossil record, being detected in the Southern tamandua
9 (*Tamandua tetradactyla*) – a xenarthran – and the armadillo (*Oryzomys azer*). The EPV locus
10 found in armadillos is relatively degraded, but the tamandua EPV sequence is nearly full-length
11 and relatively intact (**Table S7**) (**Fig. 1a**). Finally, the Neoprotoparvovirus (NpPV) clade contains
12 EPVs along with all known contemporary protoparvoviruses (**Fig. 4a**). Of the four NpPV-derived
13 EPVs we identified here, three have been reported previously [9, 40], and all were identified in
14 rodents. The novel representative was identified in the steppe mouse (*Mus spicilegus*) and
15 comprises a near complete genome (**Fig. 1**). Notably, the VP gene of this element groups
16 robustly with a bat-derived virus [41] in phylogenetic trees (**Fig. S6a**), whereas those encoded
17 by other NpPV-derived EPVs group separately, in an entirely different subclade, together with
18 VP sequences derived from carnivore and porcine protoparvoviruses. Notably, phylogenetic
19 reconstructions show that none of the rodent EPVs in the NpPV clade groups with
20 contemporary RoPVs, but instead cluster robustly with pPVs found in other mammalian host
21 groups (e.g., carnivores, artiodactyls). This suggests that horizontal transfer from rodents to
22 other mammalian orders may have been a common feature of parvovirus evolution. Tusavirus,
23 a divergent protoparvovirus of uncertain host origin [42] groups basally in the NpPV clade (**Fig.**
24 **4a**), but could potentially represent an entirely distinct, under-sampled pPV lineage.

25 Continental drift over the past 150-200 My is widely accepted to have had a dramatic
26 impact on mammalian evolution [43]. Around 200 Mya, all continents were part of an
27 interconnected landmass (Pangaea) that later separated into two subcomponents (**Fig. 4b**).
28 One (Laurasia) comprised Europe and most of Asia, while the second (Gondwanaland)
29 comprised Africa, South America, Australia, India and Madagascar). Mammalian subpopulations
30 were fragmented by these events, and then fragmented further as Gondwanaland separated
31 into its component continents. The associated genetic isolation due to geographic separation
32 (vicariance) drove the early diversification of major subgroups, including indigenous mammalian
33 lineages in South America (xenarthans and marsupials), Australia (marsupials), and Africa
34 (afrotherians). At points throughout the Cenozoic Era, placental mammal groups that evolved in

1 Laurasia (Boreoeutherians) expanded into other continental regions. For example, the
2 ancestors of contemporary New World rodents (which include capybaras, chinchillas, and
3 guinea pigs among many other, highly diversified species), are thought to have reached the
4 South American continent ~35 Mya [44].

5 The reconstructed evolutionary relationships between protoparvoviruses and
6 protoparvovirus-derived EPVs strikingly reveal the impact of mammalian vicariance – and later
7 migration – on the emergence and spread of novel protoparvovirus sublineages (**Fig. 4**). The
8 protoparvovirus phylogeny can readily be mapped onto the phylogeny of mammalian host
9 species so that the three major protoparvovirus lineages emerge in concert with major groups of
10 mammalian hosts. These evolutionary relationships, which are supported by numerous,
11 independently acquired EPV loci and ortholog sets (**Table S1**), are consistent with a
12 parsimonious evolutionary scenario under which: (i) the ancestors of the contemporary
13 protoparvovirus species were present in the ancient supercontinent of Pangaea prior to its
14 breakup; (ii) the vicariance-driven, deep divergences in the mammalian phylogeny drove the
15 emergence of distinct protoparvovirus lineages in distinct biogeographic regions throughout the
16 course of the Cenozoic Era (from 65 Mya to present); (iii) the founder event associated with
17 migration of rodents into the New World allowed this group to escape infection with NpPVs, but
18 presumably, following their colonisation of the South American continent (estimated to have
19 occurred ~50-30 Mya [44]), they were then exposed to infection with ApPVs, to the extent that
20 numerous ApPV-derived EPVs were independently fixed in the germline. A previously reported
21 ApPV-derived EPV in the common opossum (*Monodelphis domestica*) [10] groups intermediate
22 between clades comprised of Australian marsupials EPVs and NW rodent EPVs, consistent with
23 this hypothesis. Biogeographic analysis of host species distributions and ancestral range
24 reconstruction support these findings (**Fig. S6b**).

25

26 *Ancient origins and inter-order transmission of erythroparvoviruses*

27 Our comprehensive *in silico* screen of vertebrate genomes identified the first reported
28 examples of EPVs derived from genus *Erythroparvovirus*. One was identified in the genome of
29 the Patagonian mara (*Dolichotis patagonum*) - a New World rodent – and another was identified
30 in the genome of the Indri (*Indri indri*), a Malagasy primate. The mara element spans a complete
31 NS gene, whereas the indri element encodes a complete viral genome with intact NS and VP
32 genes and incorporating putative ITR sequences (**Fig. 2**), suggesting it integrated relatively
33 recently. As reported in other erPVs, the viral protein 1 unique region (VP1u) of Erythro.1-Indri
34 is relatively long. Neither contained obvious homologs of the accessory proteins reported in

1 contemporary erPV genomes. Both erPV-derived EPVs grouped with erPVs derived isolated
2 from rodents in phylogenetic trees, indicating inter-order transmission from rodents to
3 lemuriforme primates (**Fig. S7**). Furthermore, when examined in relation to the biogeographic
4 distribution of host species, these phylogenetic relationships provide tentative age calibrations
5 for the *Erythroparvovirus* genus. based on the parsimonious assumption that the presence of
6 the EPVs derived from rodent erPVs in Madagascar and South America reflects their spread
7 into these isolated geographic regions during the Cenozoic Era (**Table 3**).

8

9 *Inter-class transmission and the evolution of non-autonomous dependoparvoviruses*

10 We identified 213 dependoparvovirus-related EPV sequences in mammals, which we
11 estimate to represent at least 80 distinct germline incorporation events (**Table S3**). A small
12 number of near full-length genomes were identified, but a large share of these elements
13 spanned only small fragments (i.e. >40%) of the dependoparvovirus (dPV) genome (**Fig 1b**).
14 We reconstructed the evolutionary relationships between dPV-related EPVs and contemporary
15 dPVs (**Fig 5, Fig. S9**). In trees that included EPVs, support for internal branching order was
16 typically quite low. This reflects the short length of many dPV-related EPV sequences, and the
17 fact that many parts of the viral genome sequence are relatively degraded [13]. However, when
18 only viruses and longer EPV sequences are included, phylogenies based on *rep* gene
19 sequences disclose several robustly supported subclades within the *Dependoparvovirus* genus
20 (**Fig. 5**). They include clades exclusive to reptilian species (Sauria-), Australian marsupials
21 (Oceania-), and Boreoeutherian mammals (Neo-). A fourth clade, which we named
22 “Shirdaldependoparvovirus” (ShdPV), contains dPV taxa derived from both avian and
23 mammalian hosts. Both the composition of this clade in terms of hosts, and its phylogenetic
24 position relative to other dPV groups, implies a role for interclass transmission between
25 mammals and birds in dPV evolution (**Fig. 5b**). Firstly, the avian viruses in this clade group
26 basally, forming a paraphyletic group relative to a derived subclade - here referred to as
27 ‘Lemuria-’ - of ancient EPVs obtained from a diverse range of mammalian hosts. This topology
28 suggests that clade Lemuria- may have originated via transfer from birds to mammals.
29 Furthermore, in both midpoint-rooted phylogenies, and in phylogenies rooted on the saurian
30 dependoparvoviruses (SdPVs) (as proposed by Penzes *et al* [45]), the ShdPVs as a whole fall
31 intermediate between two exclusively mammalian groups – the neodependoparvoviruses
32 (NdPVs) found in placental mammals, and the Oceaniadependoparvoviruses (OdPVs) found in
33 Australian marsupials (**Fig. 5a**). This suggests the ShPVs originated via transmission from
34 mammals to birds.

1 The NdPVs include the non-autonomous parvoviruses (AAVs), which require a helper
2 virus for replication typically a nuclear DNA virus (e.g. herpesvirus, adenovirus [1]). In
3 phylogenies rooted on the RdPVs, the NdPVs emerge as a derived clade with the autonomously
4 replicating avian viruses grouping basal. Fragmentary EPVs found in Cercopithecine primate
5 genomes arose between 23-16 Mya and appear to represent the ancient progenitors of
6 contemporary primate AAVs (**Fig. 5a**).

7

8 **DISCUSSION**

9 In this study we recovered the complete repertoire of EPV sequences in WGS data
10 representing 738 vertebrate species. While previous studies have reported a sampling of EPV
11 diversity in vertebrates [7-13, 24, 26, 34, 35, 40, 46], the present study is an order of magnitude
12 larger in scale – we identify 595 sequences representing nearly 200 discrete germline
13 incorporation events (**Table S1, Fig 1**). Furthermore, we introduced a higher level of order to
14 these data by: (i) discriminating between unique loci and orthologous copies; (ii) hierarchically
15 arranging MSAs so that phylogenetic analysis (and taxonomic classification) of individual EPVs
16 could utilise the maximum amount of available data; (iii) applying to EPVs a standardised
17 nomenclature that captures information about orthology and taxonomy; (iv) inferring ancestral
18 reference sequences for EPV coding domains.

19 The EPVs reported here are derived from a diverse array of distinct parvovirus groups
20 (**Fig. 2**). The majority grouped within subfamily *Parvovirinae*, but we also identify rare examples
21 of EPVs derived from *Icthamaparvovirus* (a genus in subfamily *Hamaparvovirinae*) in snakes.
22 Orthologous copies of this element demonstrate that the association between
23 hamaparvoviruses and vertebrates extends to the late Mesozoic Era >100 Mya, reinforces the
24 view that this recently described subfamily is ancient and broadly distributed [46, 47]. Among
25 EPVs derived from subfamily *Parvovirinae*, the majority derived from two genera –
26 *Protoparvovirus* and *Dependoparvovirus*. However, we also identified representatives of other
27 genera (*Amdoparvovirus*, *Erythroparvovirus*) as well as several highly divergent EPVs that likely
28 represent novel genera. Among these sequences, those that were identified in mammals may
29 simply represent mutationally degraded members of the established genera, since they are
30 ancient and relatively short (**Table 3, Fig 1f**). However, those identified in basal vertebrate
31 lineages such as cartilaginous fish (class Chondrichthyes) and lobe-finned fish (clade
32 Sarcopterygii) are likely to represent novel groups. These EPV sequences also demonstrate
33 that the host range of the subfamily *Parvovirinae* extends to basal vertebrates.

1 We obtained robust minimum age calibrations based on the identification of orthologous
2 genomic flanking sequences for all parvovirus genera represented in the viral fossil record
3 except *Erythroparvovirus*. However, for erythroparvoviruses (erPVs) and many other
4 *Parvovirinae* genera (including those that are not represent in the molecular fossil record) we
5 could infer more tentative calibrations based on the distribution of EPVs and viruses across host
6 groups (**Table 3**). The most striking example of this occurs in the *Protoparvovirus* genus, in
7 which the impact of biogeographic isolation throughout the Cenozoic Era is strongly reflected in
8 the phylogenetic relationships between virus subgroups and the distribution of virus subgroups
9 across host taxonomic groups and biogeographic host ranges. A simple, parsimonious
10 explanation of these relationships is presented in **Fig. 4**, wherein ancestral protoparvoviruses
11 (pPVs) were present in mammalian ancestors inhabiting Pangea ~200 Mya, and distinct pPV
12 lineages emerged as mammalian species were compartmentalised into distinct biogeographic
13 regions by continental drift. Later, the migration of mammalian groups into previously isolated
14 continental regions provided the opportunity for these pPV subgroups to infect new host groups.
15 Thus, the ancient AdPV lineage, which evolved primarily in marsupials spread into placental
16 mammal group (New World rodents) during the later Cenozoic Era (see **Fig. 4**). This extended
17 evolutionary timeline for pPVs is supported by evidence from orthology (**Table 3**), lending
18 credibility to similar, biogeography and distribution-based age estimated inferred for viral
19 lineages in which we did not obtain minimum age estimates based on orthologous EPVs (e.g.,
20 the Ave-Boca lineage).

21 Whereas some genera, such as *Protoparvovirus* and *Dependoparvovirus*, are highly
22 represented in the genomic 'fossil record', others are conspicuously absent. For example, no
23 EPVs derived from the 'Ave-Boca' lineage, or from the *Tetraparvovirus* and *Copiparvovirus*
24 genera, were identified. However, the ancient calibrations obtained for dPVs and pPVs imply
25 that other *Parvovirinae* genera have similarly ancient origins, and thus are consistent with the
26 avian and mammalian components of the Ave-Boca lineage emerging via broad codivergence
27 with vertebrate hosts ~400-300 Mya (**Table 3**). Extending this logic, the identification of
28 *Parvovirinae* lineages in basal vertebrate lineages such as fish suggests that the subfamily
29 *Parvovirinae* may have primordial origins within vertebrates.

30 While inter-class transmission of parvoviruses appears to be rare overall, we obtained
31 compelling evidence that it has occurred in the *Dependoparvovirus* genus, specifically in the
32 evolution of a lineage that contains parvoviruses and EPVs derived from both avian and
33 mammalian hosts, and which we named "Shirdaldependoparvovirus" (ShDPV). This robustly
34 supported clade contains both the avian autonomous dependoparvoviruses (dPVs) and the

1 lemuriadependoparvoviruses (LdPVs) – a clade of mammalian dPVs that existed >80 million
2 years ago (**Table 3**) and is so far only represented by EPVs. The topology of NS/Rep
3 phylogenies cannot be reconciled with codivergence and instead implies that both ShDPV and
4 the LdPV subclade it contains arose in separate inter-class transmission events involving
5 mammals and birds (**Fig. 5**).

6 The non-autonomous dependoparvoviruses – often referred to as “adeno-associated
7 viruses” (AAVs) – are characterised by the requirement for a helper virus to replicate. All of
8 these viruses group within the ‘neodependoparvovirus’ (NdPV) clade in our analysis. Most
9 recently described AAVs – such as those identified in bats, rodents and carnivores - have only
10 been characterised at sequence-level, and little is known about their phenotypic properties.
11 However, most of these viruses fall within the range of diversity defined by two AAV groups
12 (Dependo-A and Dependo-B) indicating that the requirement for a helper virus (“dependency”) is
13 an ancestral characteristic of the NdPV and likely to be shared among most if not all AAV
14 species. Furthermore, the EPV fossil record supports the view that the host range of NdPVs
15 encompasses all placental mammals. Dependency may have evolved as a means of timing
16 replication - some large DNA viruses, such as herpesviruses, are able to establish latent,
17 persistent infections within which they can ‘sense’ the cellular environment and switch to
18 replicative mode when conditions are optimal [48]. The success of this strategy is reflected in
19 the extremely high prevalence of herpesviruses in mammalian populations (often close to
20 100%). Possibly, NdPVs can optimise their transmission by tethering their replication cycle to
21 that of these ubiquitous, sophisticated DNA viruses.

22 It seems extraordinary that so many EPVs have been fixed in the mammalian germline,
23 since the formation of a novel EPVs is almost certainly a rare event - *in utero* virus infections are
24 often lethal and virus-infected gamete cells are unlikely to be viable under most circumstances.
25 Furthermore, the neutrality principle of population genetics predicts the loss of new alleles
26 occurring at low frequency (unless there are selective advantages from the genotype). In most
27 EPVs, ORFs have been disrupted by indels and contain multiple nonsense mutations rendering
28 the ancestral viral ORFs non-translatable, but some retain long regions of intact coding
29 sequence – both NS/Rep and VP/Capsid sequences are among the longest open regions
30 (**Table S7**). The Rep protein is structurally and functionally related to the rolling circle replication
31 (RCR) proteins that are among the oldest replicator proteins known [49, 50]. The RCR proteins
32 play a pivotal role for replication of both circular and linear genomes, and therefore, are
33 inextricably linked with single-stranded DNA viruses. With the exception of the mitochondrial
34 DNA polymerase, RCR proteins are restricted to microbial and viral species. Experimental

1 studies have shown that dPV Rep protein (over) expression affects healthy cells through a
2 variety of activities including DNA binding, constitutive ATPase, and inhibiting the (cyclic) cAMP
3 -activated protein kinase A (PKA) and protein kinase X (PrKX) [51, 52]. Rep-mediated inhibition
4 of these kinases not only affects the infected cell, but also diminishes the proliferation of
5 adenovirus helper virus, perhaps attenuating the virulence and virus-induced pathogenesis [53,
6 54]. Conceivably, it could be these properties that have favoured their capture by herpesviruses
7 and by host species genomes. The selective forces that have favoured the retention of open
8 VP/capsid genes in some EPVs (see **Table S7**) are unclear.

9 The extended evolutionary timescale implied by our analysis raises interesting questions
10 about parvovirus evolution. For example, all members of the subfamily *Parvovirinae* use similar
11 basic mechanisms to achieve specific steps in infection, but the specific details of these
12 processes differ between genera. Our study suggests that these differences could have evolved
13 gradually as distinct parvovirus lineages adapted to distinct ecological niches. It is clear from
14 phylogenetic and genomic analysis that most vertebrate species are infected with multiple
15 distinct parvovirus groups - for example, at least seven distinct genera circulate in mammals.
16 Has each parvovirus genus developed specializations that allow it to occupy a unique ecological
17 niche, or are some or all parvoviruses generalists? Other questions concern the current
18 distribution of parvoviruses – e.g., to what extent does it reflect long-term evolutionary
19 processes versus (possibly) more recent, anthropogenic influences? Also, which parvovirus
20 groups currently only known via EPVs e.g. (**Fig. 4, Fig. 5**) are still prevalent, and which are
21 extinct? Our study shows that parvovirus host-associations are relatively stable over time,
22 implying that further sampling of parvovirus and EPV diversity will help address these questions.

23 EPVs allow the present-day biological properties of parvoviruses to be examined in the
24 light of a time-calibrated evolutionary history. They can also be used to investigate the structural
25 properties of ancient capsids based on molecular modelling [34] and even to reconstruct viable
26 versions of capsids from extinct paleoviruses [40]. However, making effective use of EPV data
27 is often challenging, since high levels of sequence divergence preclude straightforward analysis.
28 In this study we introduced a template for computational genomics studies of viruses that
29 focusses on facilitating the reproduction of comparative analyses and re-use of the complex
30 datasets that underpin them (e.g., MSAs) (**Fig. S1**). This approach can not only scale to
31 accommodate greatly increased quantities of virus species and sequences, but also introduces
32 new levels of reproducibility and re-usability so that researchers working in different areas of
33 parvovirus genomics - but utilizing related data and domain knowledge - can benefit from one
34 another's work. The resources and approaches developed in this study can thus facilitate the

1 development of a broader understanding of parvovirus biology, covering multiple biological
2 scales, which can be used to mitigate their harmful impacts and inform their development as
3 therapeutic tools.

4

5

6 **METHODS**

7 *Creation of resources for reproducible comparative analysis of parvovirus genomes*

8 We used the GLUE software environment [30] to create a sequence-oriented resource
9 for comparative genomic analysis of parvoviruses (including extinct paleoviruses). This
10 resource, called ‘Parvovirus-GLUE’, not only contains the data items required for comparative
11 analysis (i.e., virus genome sequences, multiple sequence alignments (MSAs), genome feature
12 annotations, and other sequence-associated data), it also represents the semantic relationships
13 between these data items via a relational database (**Fig. S1**). A library of parvovirus reference
14 sequences was compiled from the virus reference genomes resource [55] and from recent
15 papers describing novel parvoviruses and parvovirus-derived EVEs (identified via PubMed
16 search). We obtained reference genome sequences for all Parvovirus species recognised by
17 the International Committee for Taxonomy of Viruses (ICTV), as well as a representative set of
18 recently identified parvovirus-related viruses that have not yet been incorporated into official
19 taxonomy. The database is constructed using GLUE’s native command layer (**Fig. S3a**). The
20 command layer can be used to interact with the database and with bioinformatics software tools
21 required for comparative sequence analysis, thus establishing a platform for implementing
22 comparative analyses in a reproducible, standardised way.

23 Parvovirus-GLUE incorporates MSAs representing distinct taxonomic levels within the
24 family *Parvoviridae* and uses a ‘constrained alignment tree’ data structure to represent the
25 hierarchical relationships between them (**Fig. S2, Table 1**). The root alignment contains all
26 reference sequences, whereas all children of the root inherit at least one reference from their
27 immediate parent. Thus, all alignments are linked to one another via our chosen set of
28 references. This allows a single unified alignment to be used for phylogenetic analyses across
29 distinct taxonomic levels and enables standardised sequence comparisons across the entire
30 parvovirus family. A set of ‘master’ reference sequences - each representing a distinct clade in
31 the parvovirus phylogeny – was defined. Reference sequences were used to define
32 ‘constrained’ alignments (i.e., alignments in which the genomic coordinate spaces are
33 constrained to the chosen reference sequence). For the lower taxonomic levels (i.e., genus and
34 below) we aligned complete coding sequences. For the root we aligned a conserved region of

1 NS protein consistent with the approach used by ICTV [1]. We used GLUE's 'constrained
2 alignment tree' data structure [30] to link MSAs constructed at distinct taxonomic levels, via a
3 set of common reference sequences.

4

5 Genome screening in silico

6 We used the Database-Integrated Genome Screening (DIGS) tool [56] to derive a non-
7 redundant database of EPV loci within published WGS assemblies. The DIGS tool is a Perl-
8 based framework that uses the Basic Local Alignment Search Tool (BLAST) program suite [57]
9 to perform similarity searches and the MySQL relational database management system to
10 coordinate screening and record output data. A user-defined reference sequence library
11 provides (i) a source of 'probes' for searching WGS data using the tBLASTn program [57, and
12 (ii) a means of classifying DNA sequences recovered via screening **Fig. S4**. For the purposes of
13 the present project, we collated a reference library composed of polypeptide sequences derived
14 from representative parvovirus species and previously characterised EPVs. Whole genome
15 sequence (WGS) data of animal species were obtained from the National Center for
16 Biotechnology Information (NCBI) genome database [58]. We obtained all animal genomes
17 available as of March 2020. We extended the core schema of the screening database to
18 incorporate additional tables representing the taxonomic classifications of viruses, EPVs and
19 host species included in our study. This allowed us to interrogate the database by filtering
20 sequences based on properties such as similarity to reference sequences, taxonomy of the
21 closest related reference sequence, and taxonomic distribution of related sequences across
22 hosts. Using this approach, we categorised sequences into: (i) putatively novel EPV elements;
23 (ii) orthologs of previously characterised EPVs (e.g., copies containing large indels); (iii) non-
24 viral sequences that cross-matched to parvovirus probes (e.g., retrotransposons). Sequences
25 that did not match to previously reported EPVs were further investigated by incorporating them
26 into genus-level, genome-length MSAs (see **Table 1**) with representative parvovirus genomes
27 and reconstructing maximum likelihood phylogenies using RAxML (version 8) [59].

28 Where phylogenetic analysis supported the existence of a novel EPV insertion, we also
29 attempted to: (i) determine its genomic location relative to annotated genes in reference
30 genomes; and (ii) identify and align EPV-host genome junctions and pre-integration insertion
31 sites (see below). Where these investigations revealed new information (e.g., by confirming the
32 presence of a previously uncharacterised EPV insertion) we updated our reference library
33 accordingly. This in turn allowed us to reclassify putative EPV loci in our database and group
34 sequences more accurately into categories. By iterating this procedure, we progressively

1 resolved the majority of EPV sequences identified in our screen into groups of orthologous
2 sequences derived from the same initial germline incorporation event (**Table S1-S6**).

3 We applied standard identifiers (IDs) to all EPV loci, following a convention established
4 for endogenous retroviruses [60] that has more recently been applied to EVEs [19, 20]. Each
5 EVE is assigned a unique identifier (ID) constructed from two components. The first component
6 is the classifier ‘EPV’. The second component is itself a composite of two distinct
7 subcomponents separated by a period; (i) the name of the lowest level taxonomic group (i.e.,
8 species, genus, subfamily, or other clade) into which the element can be confidently placed by
9 phylogenetic analysis; (ii) a numeric ID that uniquely identifies the insertion.

10

11 *Phylogenetic and Phylogeographic analysis*

12 A process for reconstructing evolutionary relationships across the entire *Parvoviridae*
13 was implemented using GLUE. We used a data structure called a ‘constrained MSA tree’ to
14 coordinate genomic analyses across the large phylogenetic distances found in parvoviruses
15 **Fig. S7d**. This approach addresses the issue that MSAs constructed at higher taxonomic levels
16 (e.g., above genus-level in the *Parvoviridae*) typically contain far fewer columns than those
17 constructed at low taxonomic levels (e.g., genus, subgenus), meaning that to fully investigate
18 phylogenetic relationships using all available information it is often necessary to construct
19 several separate MSAs each representing a distinct taxonomic grouping (e.g., see **Table 2**).
20 The difficulties encountered in maintaining these MSAs in sync with one another while avoiding
21 irreversible data loss are a key factor underlying the low levels of reproducibility and re-use in
22 comparative genomic analyses [27], particularly those that examine more distant evolutionary
23 relationships genomic [29]. To address these issues the ‘constrained MSA tree’ data structure
24 represents the hierarchical links between MSAs constructed at distinct taxonomic levels,
25 creating in effect a single, unified MSA that can be used to reconstruct both shallow and deep
26 evolutionary relationships while making use of the maximum amount of available information at
27 each level (**Table 2, Fig. S1d**). This approach also has the effect of standardising the genomic
28 coordinate space to the constraining reference sequences selected for each MSA without
29 imposing any limitations on which references are used (e.g., laboratory strains versus wild-type
30 references), since additional, alternative constraining references can be incorporated, and
31 contingencies such as insertions relative to the constraining reference are dealt with in a
32 standardised way [30].

33 MSAs partitions derived from the constrained MSA tree were used as input for
34 phylogenetic reconstructions. Nucleotide and protein phylogenies were reconstructed using

1 maximum likelihood (ML) as implemented in RAxML (version 8.2.12) [59]. To handle coverage-
2 related issues we generated gene coverage data prior to phylogenetic analysis and used this
3 information to condition the way in which taxa are selected into MSA partitions. Protein
4 substitution models were selected via hierarchical maximum likelihood ratio test using the
5 PROTAUTOGAMMA option in RAxML. For multicopy EPV lineages we constructed MSAs and
6 phylogenetic trees to confirm that branching relationships follow those of host species (**Fig S4b**,
7 [31]).

8 Time-calibrated vertebrate phylogenies were obtained via TimeTree [61]. We used a
9 time-calibrated phylogeny of protoparvovirus host species and the present continental
10 distribution of host organisms to model ancestral biogeographical range of protoparvovirus host
11 [62] (**Fig. 6b**). Country-level distribution information for each host species was obtained via the
12 `occ_search` function of the `rgbif` library in R [63]. Country records were consolidated in continent
13 entries with the `continents` function of the `countrycode` library and manually curated to ensure
14 accuracy. Within continents, North Africa and Sub-Saharan Africa were considered distinct
15 distributions and coded separately. The Dispersal-Extinction-Cladogenesis (DEC) model
16 implemented in the program Lagrange C++ was applied without constraining the number of
17 ancestral states nor limiting connectivity between biogeographic units [64]. Ancestral states at
18 all nodes in the tree were inferred and the tree visualized in R with the `ggtree` and `ggplot`
19 libraries [65, 66]. Input data and configuration files for Lagrange along with the time tree and
20 Lagrange output are provided in the Data Supplement [31].

21

22 Genomic analysis of EPVs

23 ORFs were inferred by manual comparison of sequences to those of reference viruses.
24 For phylogenetic analysis, the putative peptide sequences of EVEs (i.e., the virtually translated
25 sequences of EVE ORFs, repaired to remove frameshifting indels) were aligned with
26 polypeptide sequences encoded by reference genomes. We used PAL2NAL [67] to generate in-
27 frame, DNA alignments of virus coding domains from alignments of polypeptide gene products.
28 Phylogenies were reconstructed using maximum likelihood (ML) as implemented in RAxML [59]
29 and GTR model of nucleotide selection as selected using the likelihood ratio test. The putative
30 peptide sequences of EPVs were aligned with NS and VP polypeptides of representative
31 exogenous parvoviruses using MUSCLE.

32

33 Expression and intactness of EPVs

1 We identified open coding regions of coding sequence in EPVs by using PERL scripts,
2 (included with Parvovirus-GLUE [31]) to process EPV sequence data. To determine if there was
3 evidence of expression of EPVs in host species, we searched the NCBI Reference RNA
4 Sequences (refseq_rna) with Dependoparvovirus VP and Rep sequences (NC_002077). We
5 searched a translated nucleotide query and a translated database using tBLASTx [57] and
6 evaluated alignments found between refseq_rna sequences and Dependoparvovirus VP and
7 Rep sequences. To further verify expression, we determined if the annotations were solely
8 based on computational prediction or if there is RNAseq data alignment to the annotation in
9 support of the feature. For those host species with evidence of expression, we conducted blastn
10 searches within refseq_rna to identify expressed EPVs.

11

Table 1. Multiple sequence alignments included in Parvoviridae-GLUE

#	Scope/Name	Parent	Children	Constraining reference	Coverage*	Viruses	EVE loci
Root							
1	<i>Parvoviridae</i>	none	3	CPV	NS (13%)	3	0
Subfamily							
2	<i>Parvovirinae</i>	<i>Parvoviridae</i>	2	CPV	NS (63%)	13	4
3	<i>Hamaparvovirinae</i>	<i>Parvoviridae</i>	2	PPV7	NS	5	0
4	<i>Densoparvovirinae</i>	<i>Parvoviridae</i>	0	JcDENV	NS	9	1*
Cross-genus							
5	Boca-Ave	<i>Parvovirinae</i>	2	ChPV	Genome (70%)	2	0
6	Amdo-Proto	<i>Parvovirinae</i>	2	CPV	Genome (77%)	2	0
7	EDCT	<i>Parvovirinae</i>	3	HPV4	Genome (57%)	4	0
8	Chaphama-Icthama	<i>Hamaparvovirinae</i>	2	PPV7	Genome	2	0
Genus							
9	<i>Aveparvovirus</i>	Boca-Ave	0	ChPV	Genome (88%)	4	0
10	<i>Bocaparvovirus</i>	Boca-Ave	0	BPV	Genome (75%)	7	0
11	<i>Erythroparvovirus</i>	EDCT	0	B19	Genome (80%)	9	2*
12	<i>Tetraparvovirus</i>	EDCT	0	HPV4	Genome (80%)	10	0
13	<i>Dependoparvovirus</i>	EDCT	0	AAV2	Genome (84%)	27	81
14	<i>Copiparvovirus</i>	EDCT	0	BPV2	Genome (62%)	2	0
15	<i>Amdoparvovirus</i>	Amdo-Proto	0	AMDV	Genome (85%)	6	6
16	<i>Protoparvovirus</i>	Amdo-Proto	0	CPV	Genome (90%)	18	106
17	<i>Chaphamaparvovirus</i>	<i>Hamaparvovirinae</i>	0	PPV7	Genome (85%)	12	4**
18	<i>Icthamaparvovirus</i>	<i>Hamaparvovirinae</i>	0	SyIPV	Genome (62%)	2	2
Totals						137	201

Footnote: Linking alignments that represent internal nodes contain only the reference sequences for their 'child' alignments in the constrained alignment tree. *Putatively exogenous parvoviruses identified in this study. **Putatively exogenous parvoviruses identified in previous studies. Abbreviations: EDCT=Erythro-Dependo-Copi-Tetra group; CPV=canine parvovirus; JcDENV=*Junonia coenia* densovirus; PPV7=porcine parvovirus 7; ChPV=chicken parvovirus; HPV4=human parvovirus 4; B19=Human erythroparvovirus B19; BPV=bovine parvovirus; AMDV=Aleutian mink disease virus; *Syngnathus scovelli* ichthamaparvovirus.

Table 2. Incorporation of parvovirus DNA into the vertebrate germline

Parvovirus genus	Host species group									
	Chondrichthyes <i>species=5</i>		Actinopterygii <i>species =175</i>		Sauria <i>species =200</i>		Mammalia <i>species=353</i>		Vertebrata <i>species =752*</i>	
	<i>loci</i>	<i>ratio*</i>	<i>loci</i>	<i>ratio</i>	<i>loci</i>	<i>ratio</i>	<i>loci</i>	<i>ratio</i>	<i>loci</i>	<i>ratio</i>
<i>Ichthamaparvovirus</i>	0	-	1	0.01	2	0.01	0	-	2	0.003
<i>Erythroparvovirus</i>	0	-	0	-	0	-	2	0.01	2	0.003
<i>Amdoparvovirus</i>	0	-	0	-	1	0.01	2	0.01	3	0.004
<i>Dependoparvovirus</i>	0	-	0	-	13	0.07	66	0.19	80	0.108
<i>Protoparvovirus</i>	0	-	0	-	0	-	105	0.3	105	0.142
New clades	1	0.2	2	0.01	0	-	3	0.01	6	0.008
Totals	<u>1</u>	<u>0.2</u>	<u>3</u>	<u>0.02</u>	<u>16</u>	<u>0.08</u>	<u>178</u>	<u>0.5</u>	<u>198</u>	<u>0.263</u>

Footnote: *Agnathans (n=3), Amphibians (n=15), and Lungfish (*Latimeria chalumnae*) were also screened but results are not shown since no species in these groups were found to have EPVs. * Ratio = unique loci / genomes screened.

Table 3. Calibrations of parvovirus evolution

Parvovirus lineage	Host species lineage(s)	High	Low
Ortholog-based*			
Primate AAVs (Dependo-A)	OW Primates	23	16
Neodependo-	Glires	88	76
Neodependo-	Lagomorpha	77	23
Neodependo-	Vespertilionidae	49	38
Neodependo-	Elephantidae	23	9
Neodependo-	Eulemur	9	6
Neodependo-	Hyracoidea	14	7
Lemuriadependo-	Whippomorpha	56	52
Lemuriadependo-	Rhinocerotidae	51	15
Lemuriadependo-	Phyllostomidae	39	35
Oceaniadependo-	Macropus	45	27
Amdo-	Serpentes	111	100
Amdo-	Hyracoidea	14	7
Proto-	Afrotheria		
Proto-	Marsupials		
EDTC	Laurasitheria	84	73
Ichthama	Serpentes	74	49
Archaeo-Proto-	Ctenomyidae-Octodontidae	24	16
Codivergence based			
Proto-Amdo lineage	Sharks/Bony fish		
Neoprotoparvovovirus	Eutherian mammals	200	150
Boca-Ave lineage	Aves/Mammalia	393	294
Copiparvovirus	Eutherian mammals	111	100
Tetraparvovirus	Eutherian mammals	111	100
Erythroparvovirus	Eutherian mammals	111	100
Dependoparvovirus	Euteleostii	446	425
Biogeography-linked			
Archeoproto-	Pangea	200	180
Archeoproto- NW rodent clade	NW Rodents S. America	50	30
Erythroparvo- Rodent clade	Malagasy rodents Madagascar	30	20
U94 gene transfer			
EDTC lineage	Betaherpesviruses	200	80

Footnote: *Not all ortholog-based calibrations are shown, only the oldest for each virus lineage in which we identified orthologous sets of EPV sequences. Complete records of ortholog-based dates can be found in **tables S1-S6**.

Figure 1. Genomic structures of unique EPV loci.

(a) Protoparvovirus-derived EPV loci shown relative to the canine parvovirus (CPV) genome; (b) Dependoparvovirus-derived EPVs loci shown relative to the adeno-associated virus 2 (AAV-2) genome; (c) EPV loci derived from Amdoparvovirus-like viruses shown relative to the Aleutian mink disease (AMDV) genome; (d) Erythroparvovirus-derived loci shown relative to the parvovirus B19 genome; (e) EPVs derived from unclassified parvoviruses shown relative to a generic parvovirus genome. (f) Icthamaparvovirus-derived loci shown relative to *Syngnathus scovelli* parvovirus (SscPV); EPV locus identifiers are shown on the left. Solid bars to the right of each EPV set show taxonomic subgroupings below genus level. Where numbers are shown to the immediate right, the sequence shown is a consensus and numbers indicate how many individual orthologs sequences were used to create the consensus. Boxes bounding EPV elements indicate either (i) the presence of an identified gene (see **Tables S1-S6**), (ii) an uncharacterised genomic flanking region, or (iii) a truncated contig sequence (see key). EPV locus identifiers use six letter abbreviations to indicate host species (**Table S8**). **Abbreviations:** NS = non-structural protein; VP = capsid protein; ORF = open reading frame. ITR=Inverted terminal repeat; PLA2 = phospholipase A2 motif.

Figure 2. Evolution of subfamily *Parvoviridae*.

A maximum likelihood phylogeny showing the reconstructed evolutionary relationships between contemporary parvoviruses and the ancient parvovirus species represented by endogenous parvoviral elements (EPVs). Panels (A) and (B) show a more detailed view of subclades (labelled I and II) within the phylogeny shown in panel (C). The complete phylogeny, which is midpoint rooted for display purposes, was reconstructed using a multiple sequence alignment spanning 270 amino acid residues positions of the Rep protein and the LG likelihood substitution model. Coloured brackets indicate the established parvovirus genera recognised by the International Committee for the Taxonomy of Viruses. Bootstrap support values (1000 replicates) are shown for deeper internal nodes only. Scale bars show evolutionary distance in substitutions per site. Taxa labels are coloured based on taxonomic grouping as indicated by brackets, unclassified taxa are shown in black. Viral taxa are shown in bold, while EPV taxa are shown in regular text. **Abbreviations:** PV=Parvovirus; HHV=Human herpesvirus; AAV=Adeno-associated virus; AMDV=Aleutian mink disease; CPV=canine parvovirus; BPV=bovine parvovirus; BrdPV=Bearded dragon parvovirus; MdPV=Muscovy duck parvovirus; SIPV=slow loris parvovirus. TS=transcription strategy; MTSP=Multiple transcriptional start

positions; STSP+=single transcription start position plus additional strategies;
HOMO=homoteleomeric; HETERO=heteroteleomeric.

Figure 3. Incorporation of EPVs into the vertebrate germline.

A time-calibrated evolutionary tree of vertebrate species examined in this study, illustrating the distribution of germline incorporation events over time. Colours indicate parvovirus genera as shown in the key. Diamonds on internal nodes indicate minimum age estimates for EPV loci endogenization (calculated for EPV loci found in >1 host species). Coloured circles adjacent to tree tips indicate the presence of EPVs in host taxa, with the diameter of the circle reflecting the number of EPVs identified. Brackets to the left show taxonomic groups within vertebrates.

Figure 4. Protoparvovirus evolution has been shaped by mammalian vicariance.

(A) Maximum likelihood-based phylogenetic reconstructions of evolutionary relationships between contemporary parvovirus species and the ancient parvovirus species represented by endogenous parvoviral elements (EPVs). The phylogeny was constructed from a multiple sequence alignment spanning 712 amino acid residues in the Rep protein (substitution model=LG likelihood). The tree is midpoint rooted for display purposes. Asterisks indicate nodes with bootstrap support >70% (1000 replicates). The scale bar shows evolutionary distance in substitutions per site. Coloured brackets to the right indicate (i) subgroups within the *Protoparvovirus* genus (outer set of brackets) and (ii) the host range of each subgroup (inner set of brackets). Terminal nodes are represented by squares (EPVs) and circles (viruses) and are coloured based on the biogeographic distribution of the host species in which they were identified. Coloured diamonds on internal nodes show the inferred ancestral distribution of parvovirus ancestors, using colours that reflect the patterns of continental drift and associated mammalian vicariance shown in the maps in panel (B). **Evidence for the presence of the “Mesoprotoparvovirus” group in Afrotherians is presented in **Fig. 2.** **(B)** Mollweide projection maps showing how patterns of continental drift from 200-35 led to periods of biogeographic isolation for terrestrial mammals in Laurasia (Europe and Asia), South America, Australia Africa and Madagascar. The resulting vicariance is thought have contributed to the diversification of mammals, reflected in the mammalian phylogeny as shown in in panel (c). The majority of placental mammals (including rodents, primates, ungulates and bats) evolved in Laurasia. However, these groups later expanded into other continents, and fossil evidence indicates that the ancestors of today’s “New World rodents” had arrived on the South American continent by~35 million years ago (Mya), if not earlier **(C)** A time-calibrated phylogeny of mammals with

annotations indicating the biogeographic associations of the major taxonomic groups of contemporary mammals and ancestral mammalian groups, following panel (b) and key 1. **(D)** A time-calibrated phylogeny of mammals annotated to indicate the distribution of protoparvovirus subgroups among mammalian groups, following key 2. Question marks indicate where it is unknown whether viral counterparts of the lineages represented by EPVs still circulate among contemporary members of the host species groups in which they are found. **Abbreviations:** Mya = millions of years ago; NW=New World; (OW); CPV=carnivore parvovirus type 1; PPV=porcine parvovirus; HV=Hamster parvovirus; TuV=Tusavirus.

Figure 5. Dependoparvovirus evolution and the influence of inter-class transmission.

(A) A maximum likelihood phylogeny showing the reconstructed evolutionary relationships between contemporary dependoparvovirus species and the ancient dependoparvovirus species represented by EPVs. Virus taxa names are shown in bold, EPVs are shown in regular text. The phylogeny was constructed from a multiple sequence alignment spanning (MSA) 330 amino acid residues of the Rep protein and the LG likelihood substitution model and is rooted on the reptilian lineage as proposed by Penzes *et al* [45]. Brackets to the right indicate proposed taxonomic groupings. Shapes on leaf nodes indicate full-length EPVs and EPVs containing intact/expressed genes. Numbers next to leaf nodes indicate minimum age calibrations for EPV orthologs. Shapes on branches and internal nodes indicate different kinds of minimum age estimates for parvovirus lineages, as shown in the key. Numbers adjacent to node shapes show minimum age estimates in millions of years before present. For taxa that are not associated with mammals, organism silhouettes indicate species associations, following the key. The scale bar shows evolutionary distance in substitutions per site. Asterisks in circles indicate nodes with bootstrap support >70% (1000 replicates) in the tree shown. Plain asterisks next to internal nodes indicate nodes that are not supported in the tree shown here but do have bootstrap support >70% (1000 replicates) in phylogenies based on longer MSA partitions within Rep (but including less taxa). *Age calibrations based on data obtained in separated publications – see references [13] and [25]. **A contemporary virus derived from the marsupial clade has been reported in marsupials, but only transcriptome-based evidence is available [12]. **(B)** A time-calibrated phylogeny of vertebrate lineages showing proposed patterns of inter-class transmission in the Shirdaldependoparvovirus lineage. **Abbreviations:** PV=Parvovirus; AAV=Adeno-associated virus; BrdPV=Bearded dragon parvovirus; MdPV=Muscovy duck parvovirus.

ACKNOWLEDGEMENTS

This work was supported by funding from the Association Monégasque Contre les Myopathies, and the Bill & Melinda Gates Foundation (OPP1202116).

COMPETING INTERESTS

R.M.K. is a co-founder of Synteny Therapeutics, Inc., which is a co-assignee of a patent application filed on behalf of University of Massachusetts Medical School and Synteny Therapeutics, Inc.

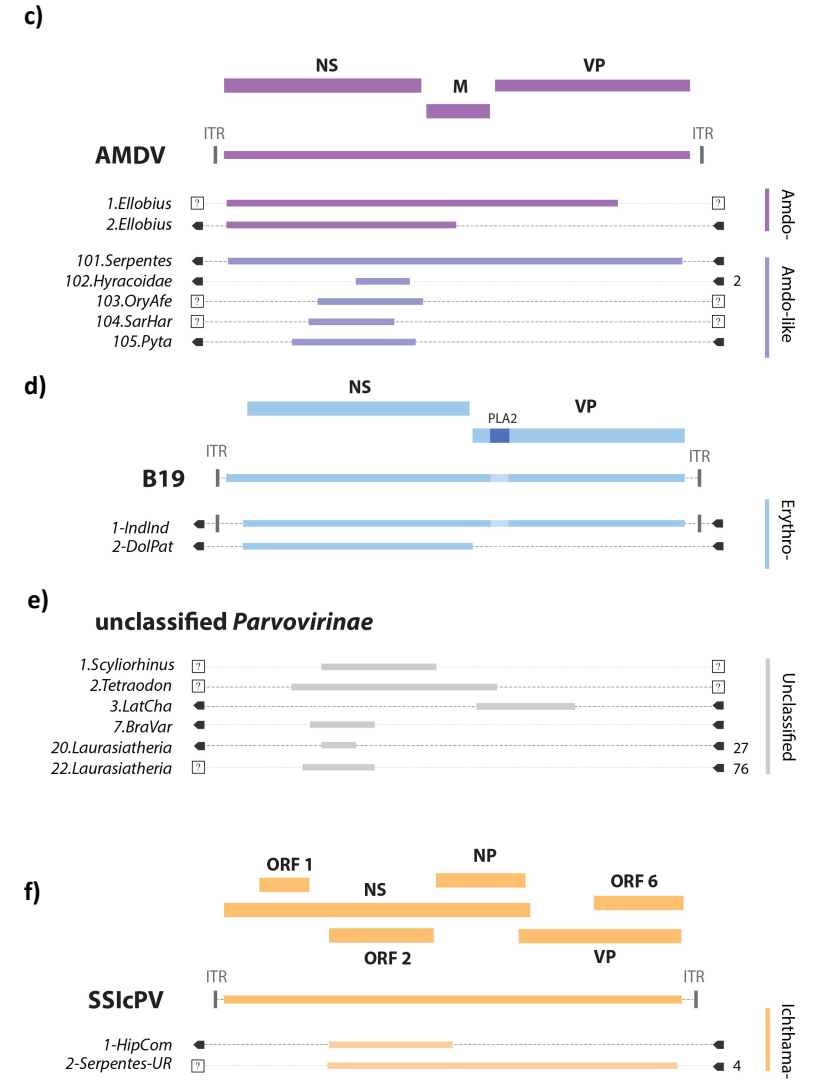
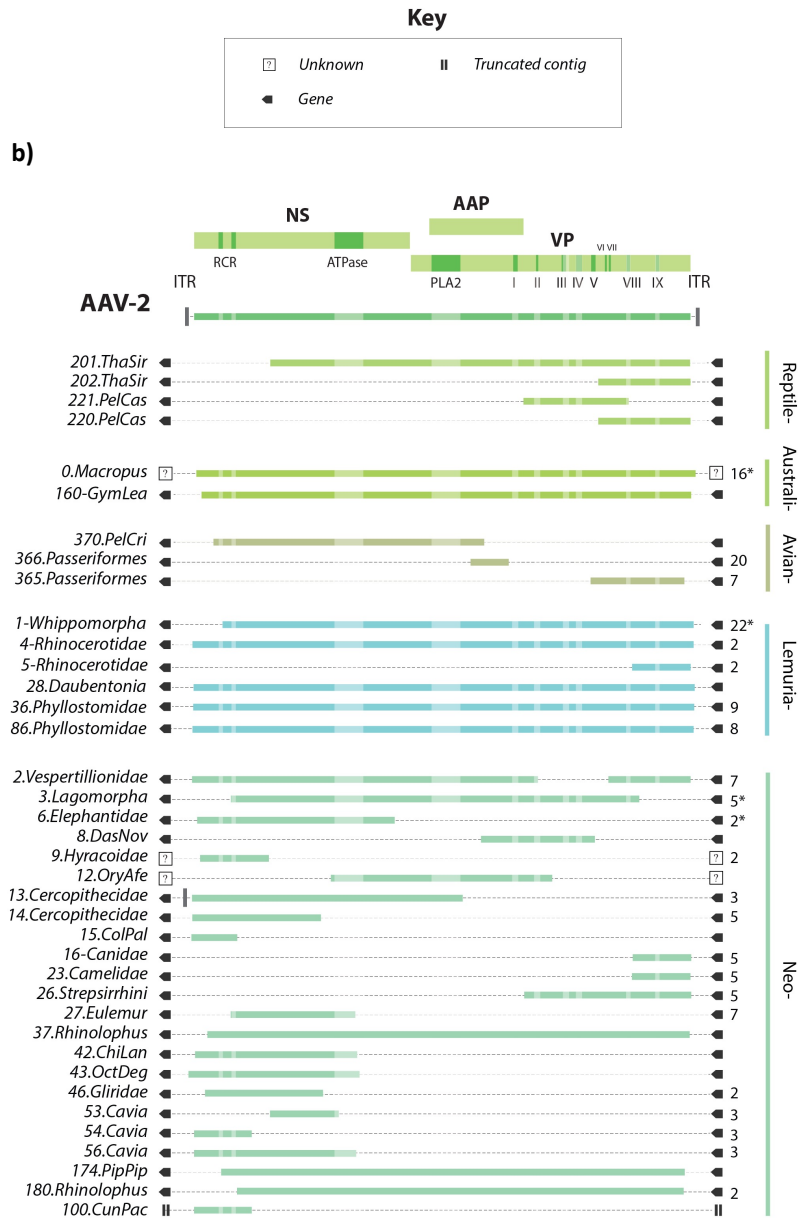
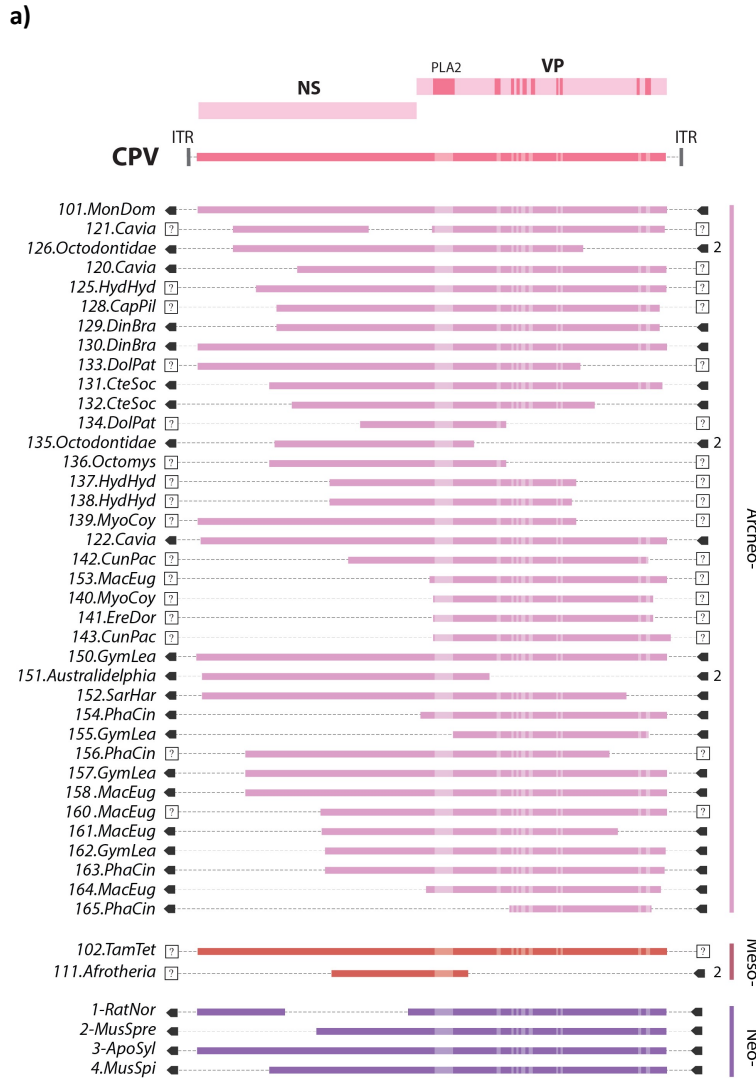
References

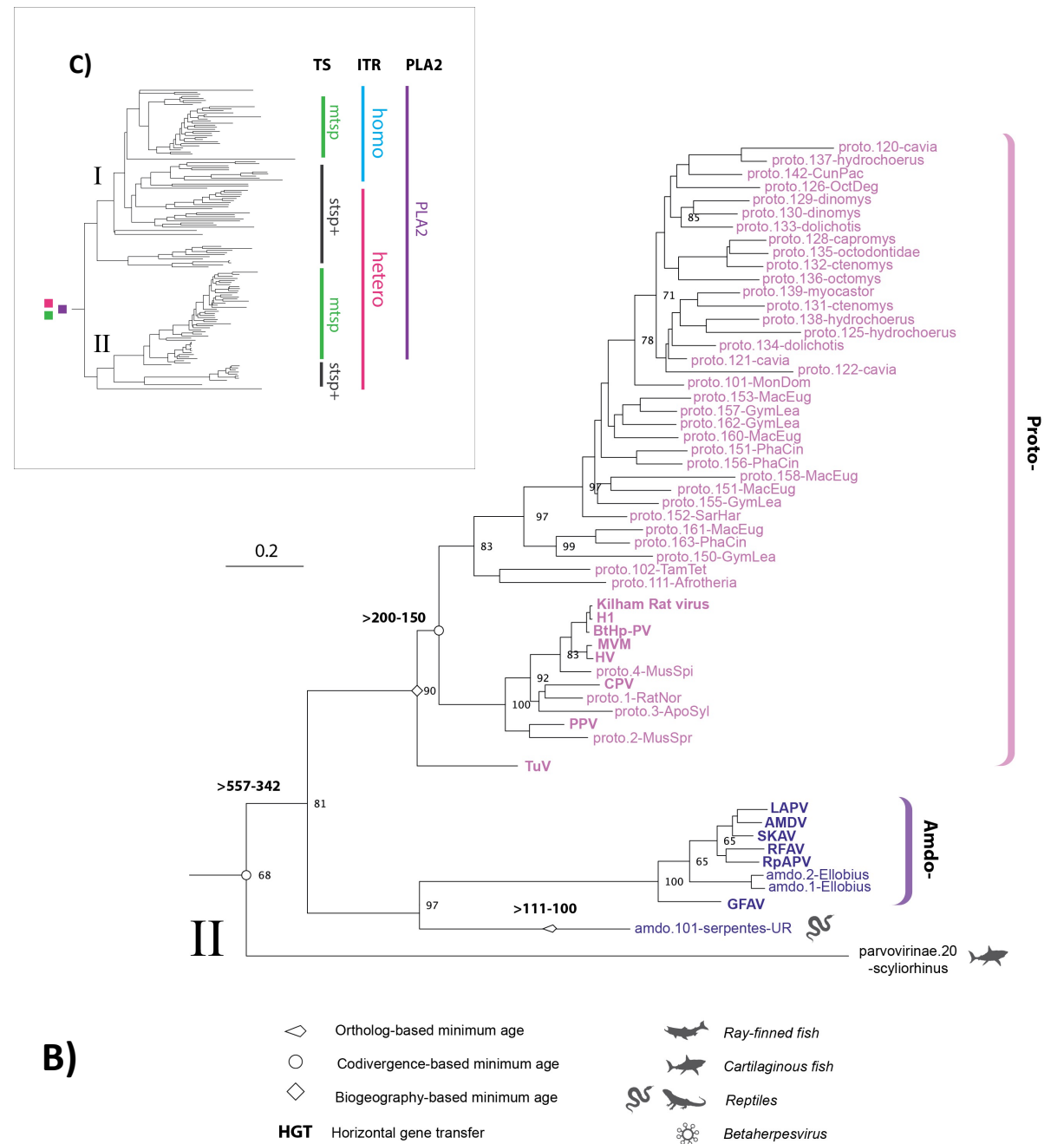
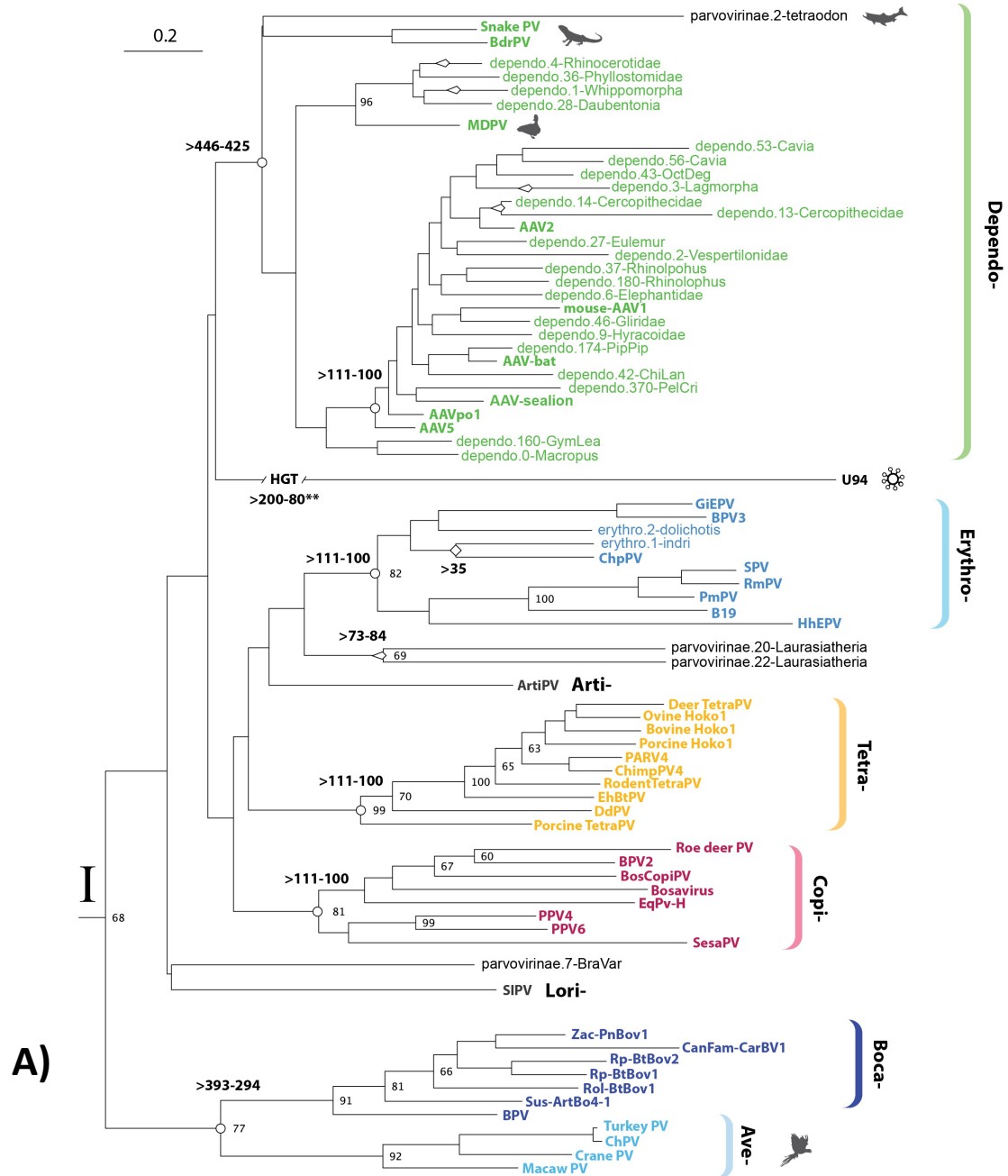
1. Cotmore, S.F., et al., *ICTV Virus Taxonomy Profile: Parvoviridae*. J Gen Virol, 2019. **100**(3): p. 367-368.
2. Cotmore, S.F. and P. Tattersall, *Parvoviruses: Small Does Not Mean Simple*. Annu Rev Virol, 2014. **1**(1): p. 517-37.
3. Nüesch, J.P., et al., *Molecular pathways: rodent parvoviruses--mechanisms of oncolysis and prospects for clinical cancer treatment*. Clin Cancer Res, 2012. **18**(13): p. 3516-23.
4. Hartley, A., et al., *A Roadmap for the Success of Oncolytic Parvovirus-Based Anticancer Therapies*. Annu Rev Virol, 2020. **7**(1): p. 537-557.
5. Fakhiri, J. and D. Grimm, *Best of most possible worlds: Hybrid gene therapy vectors based on parvoviruses and heterologous viruses*. Mol Ther, 2021.
6. Naso, M.F., et al., *Adeno-Associated Virus (AAV) as a Vector for Gene Therapy*. BioDrugs, 2017. **31**(4): p. 317-334.
7. Katzourakis, A. and R.J. Gifford, *Endogenous viral elements in animal genomes*. PLoS Genet, 2010. **6**(11): p. e1001191.
8. Belyi, V.A., A.J. Levine, and A.M. Skalka, *Sequences from ancestral single-stranded DNA viruses in vertebrate genomes: the parvoviridae and circoviridae are more than 40 to 50 million years old*. J Virol, 2010. **84**(23): p. 12458-62.
9. Kapoor, A., P. Simmonds, and W.I. Lipkin, *Discovery and characterization of mammalian endogenous parvoviruses*. J Virol, 2010. **84**(24): p. 12628-35.
10. Liu, H., et al., *Widespread endogenization of densoviruses and parvoviruses in animal and human genomes*. J Virol, 2011. **85**(19): p. 9863-76.
11. Francois, S., et al., *Discovery of parvovirus-related sequences in an unexpected broad range of animals*. Sci Rep, 2016. **6**: p. 30880.
12. Smith, R.H., et al., *Germline viral "fossils" guide in silico reconstruction of a mid-Cenozoic era marsupial adeno-associated virus*. Sci Rep, 2016. **6**: p. 28965.
13. Hildebrandt, E., et al., *Evolution of dependoparvoviruses across geological timescales – implications for design of AAV-based gene therapy vectors*. Virus Evolution, 2020.
14. Kotin, R.M., R.M. Linden, and K.I. Berns, *Characterization of a preferred site on human chromosome 19q for integration of adeno-associated virus DNA by non-homologous recombination*. Embo j, 1992. **11**(13): p. 5071-8.
15. Weitzman, M.D., et al., *Adeno-associated virus (AAV) Rep proteins mediate complex formation between AAV DNA and its integration site in human DNA*. Proc Natl Acad Sci U S A, 1994. **91**(13): p. 5808-12.

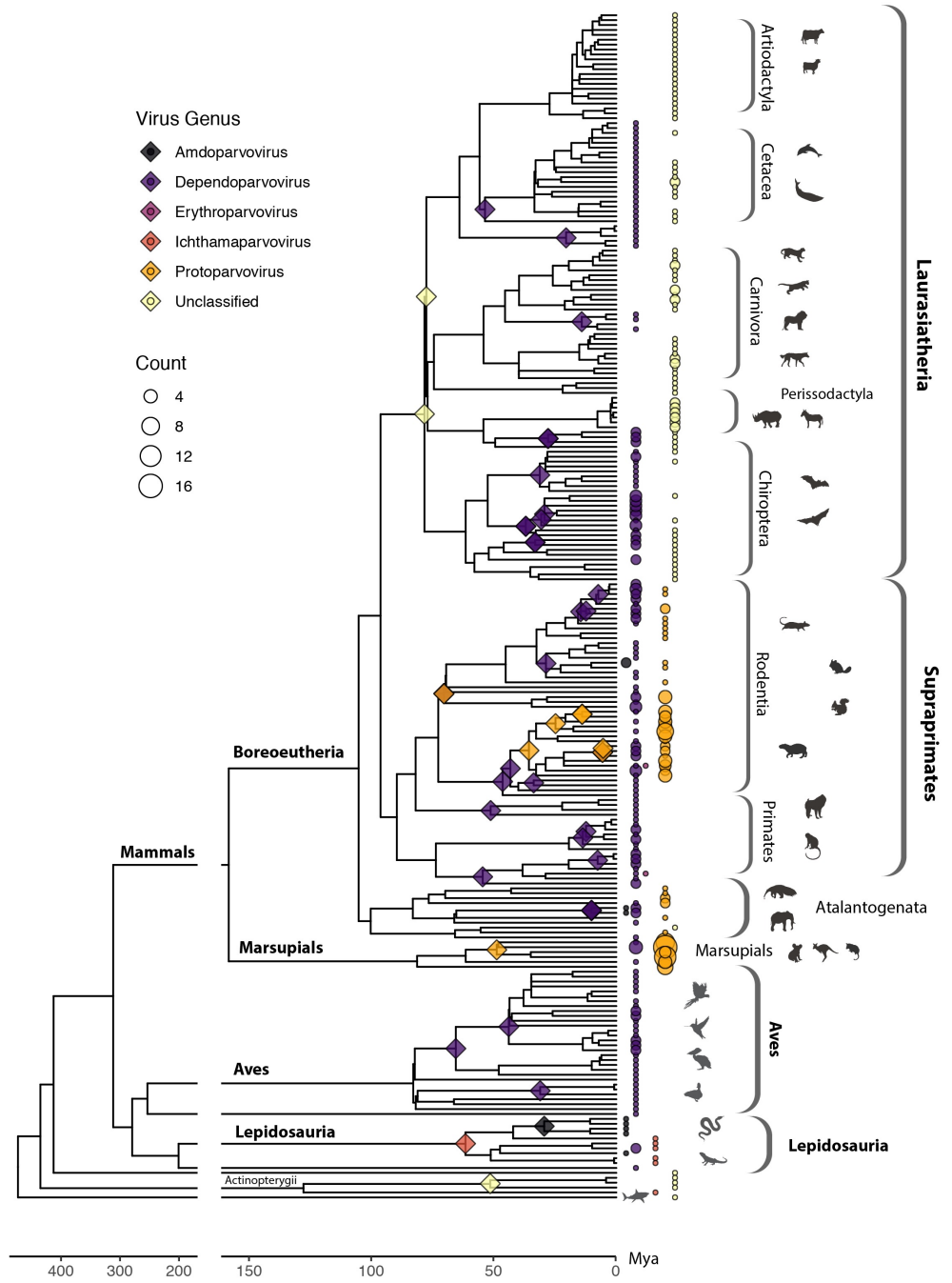
16. Kawasaki, J., et al., *One hundred million years history of bornavirus infections hidden in vertebrate genomes*. Proc Natl Acad Sci U S A, 2021: p. 2020.12.02.408005.
17. Katzourakis, A., et al., *Discovery and analysis of the first endogenous lentivirus*. Proc Natl Acad Sci U S A, 2007. **104**(15): p. 6261-5.
18. Feschotte, C. and C. Gilbert, *Endogenous viruses: insights into viral evolution and impact on host biology*. Nat Rev Genet, 2012. **13**(4): p. 283-96.
19. Lytras, S., G. Arriagada, and R.J. Gifford, *Ancient evolution of hepadnaviral paleoviruses and their impact on host genomes*. Virus Evol, 2021. **7**(1): p. veab012.
20. Kawasaki, J., et al., *100-My history of bornavirus infections hidden in vertebrate genomes*. Proc Natl Acad Sci U S A, 2021. **118**(20).
21. Cornelis, G., et al., *An endogenous retroviral envelope syncytin and its cognate receptor identified in the viviparous placental Mabuya lizard*. Proc Natl Acad Sci U S A, 2017. **114**(51): p. E10991-e11000.
22. Pastuzyn, E.D., et al., *The Neuronal Gene Arc Encodes a Repurposed Retrotransposon Gag Protein that Mediates Intercellular RNA Transfer*. Cell, 2018. **173**(1): p. 275.
23. Herniou, E.A., et al., *When parasitic wasps hijacked viruses: genomic and functional evolution of polydnviruses*. Philos Trans R Soc Lond B Biol Sci, 2013. **368**(1626): p. 20130051.
24. Arriagada, G. and R.J. Gifford, *Parvovirus-derived endogenous viral elements in two South American rodent genomes*. J Virol, 2014. **88**(20): p. 12158-62.
25. Kobayashi, Y., et al., *An endogenous adeno-associated virus element in elephants*. Virus Res, 2018.
26. Valencia-Herrera, I., et al., *Molecular Properties and Evolutionary Origins of a Parvovirus-Derived Myosin Fusion Gene in Guinea Pigs*. J Virol, 2019. **93**(17).
27. Grüning, B., et al., *Practical Computational Reproducibility in the Life Sciences*. Cell Syst, 2018. **6**(6): p. 631-635.
28. Ali, R.H., M. Bogusz, and S. Whelan, *Identifying Clusters of High Confidence Homologies in Multiple Sequence Alignments*. Mol Biol Evol, 2019. **36**(10): p. 2340-2351.
29. Holmes, E.C. and S. Duchêne, *Can Sequence Phylogenies Safely Infer the Origin of the Global Virome?* mBio, 2019. **10**(2).
30. Singer, J.B., et al., *GLUE: a flexible software system for virus sequence data*. BMC Bioinformatics, 2018. **19**(1): p. 532.
31. Gifford, R.J. *Parvovirus-GLUE*. 2021; Available from: <https://giffordlabcvr.github.io/Parvovirus-GLUE/>.
32. Merkel, D., *Docker: lightweight linux containers for consistent development and deployment*. Linux Journal, 2014. **239**(2).
33. Campbell, K., et al., *Making Genomic Surveillance Deliver: A Lineage Classification and Nomenclature System to Inform Rabies Elimination*. bioRxiv, 2021: p. 2021.10.13.464180.
34. Péntzes, J.J., et al., *Endogenous amdoparvovirus-related elements reveal insights into the biology and evolution of vertebrate parvoviruses*. Virus evolution, 2018. **4**(2): p. vey026-vey026.
35. Cui, J., et al., *Low frequency of paleoviral infiltration across the avian phylogeny*. Genome Biol, 2014. **15**(12): p. 539.
36. Gompels, U.A., et al., *The DNA sequence of human herpesvirus-6: structure, coding content, and genome evolution*. Virology, 1995. **209**(1): p. 29-51.
37. Aswad, A. and A. Katzourakis, *The first endogenous herpesvirus, identified in the tarsier genome, and novel sequences from primate rhadinoviruses and lymphocryptoviruses*. PLoS genetics, 2014. **10**(6): p. e1004332-e1004332.

38. McGeoch, D.J., F.J. Rixon, and A.J. Davison, *Topics in herpesvirus genomics and evolution*. Virus Res, 2006. **117**(1): p. 90-104.
39. Cotmore, S.F. and P. Tattersall, *Parvovirus diversity and DNA damage responses*. Cold Spring Harb Perspect Biol, 2013. **5**(2).
40. Callaway, H.M., et al., *Examination and Reconstruction of Three Ancient Endogenous Parvovirus Capsid Protein Gene Remnants Found in Rodent Genomes*. J Virol, 2019. **93**(6).
41. Wu, Z., et al., *Deciphering the bat virome catalog to better understand the ecological diversity of bat viruses and the bat origin of emerging infectious diseases*. Isme j, 2016. **10**(3): p. 609-20.
42. Väisänen, E., et al., *Human Protoperoviruses*. Viruses, 2017. **9**(11).
43. Springer, M.S., et al., *The historical biogeography of Mammalia*. Philos Trans R Soc Lond B Biol Sci, 2011. **366**(1577): p. 2478-502.
44. Poux, C., et al., *Arrival and diversification of caviomorph rodents and platyrrhine primates in South America*. Syst Biol, 2006. **55**(2): p. 228-44.
45. Péntzes, J.J., et al., *Novel parvoviruses in reptiles and genome sequence of a lizard parvovirus shed light on Dependoparvovirus genus evolution*. J Gen Virol, 2015. **96**(9): p. 2769-2779.
46. Souza, W.M., et al., *Chapparvoviruses occur in at least three vertebrate classes and have a broad biogeographic distribution*. J Gen Virol, 2017. **98**(2): p. 225-229.
47. Péntzes, J.J., et al., *An Ancient Lineage of Highly Divergent Parvoviruses Infects both Vertebrate and Invertebrate Hosts*. Viruses, 2019. **11**(6).
48. Imperiale, M.J. and M. Jiang, *What DNA viral genomic rearrangements tell us about persistence*. J Virol, 2015. **89**(4): p. 1948-50.
49. Hickman, A.B., et al., *Structural unity among viral origin binding proteins: crystal structure of the nuclease domain of adeno-associated virus Rep*. Mol Cell, 2002. **10**(2): p. 327-37.
50. Koonin, E.V., *Temporal order of evolution of DNA replication systems inferred by comparison of cellular and viral DNA polymerases*. Biol Direct, 2006. **1**: p. 39.
51. Schmidt, M., S. Afione, and R.M. Kotin, *Adeno-associated virus type 2 Rep78 induces apoptosis through caspase activation independently of p53*. J Virol, 2000. **74**(20): p. 9441-50.
52. Zimmermann, B., et al., *PrKX is a novel catalytic subunit of the cAMP-dependent protein kinase regulated by the regulatory subunit type I*. J Biol Chem, 1999. **274**(9): p. 5370-8.
53. Di Pasquale, G. and J.A. Chiorini, *PKA/PrKX activity is a modulator of AAV/adenovirus interaction*. Embo j, 2003. **22**(7): p. 1716-24.
54. Hermonat, P.L., *The adeno-associated virus Rep78 gene inhibits cellular transformation induced by bovine papillomavirus*. Virology, 1989. **172**(1): p. 253-61.
55. Brister, J.R., et al., *NCBI viral genomes resource*. Nucleic Acids Res, 2015. **43**(Database issue): p. D571-7.
56. Zhu, H., et al., *Database-integrated genome screening (DIGS): exploring genomes heuristically using sequence similarity search tools and a relational database*. bioRxiv, 2018.
57. Altschul, S.F., et al., *Gapped BLAST and PSI-BLAST: a new generation of protein database search programs*. Nuc. Acids Res., 1997. **25**: p. 3389-3402.
58. Kitts, P.A., et al., *Assembly: a resource for assembled genomes at NCBI*. Nucleic Acids Res, 2016. **44**(D1): p. D73-80.
59. Stamatakis, A., *RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies*. Bioinformatics, 2014. **30**(9): p. 1312-3.
60. Gifford, R.J., et al., *Nomenclature for endogenous retrovirus (ERV) loci*. Retrovirology, 2018. **15**(1): p. 59.

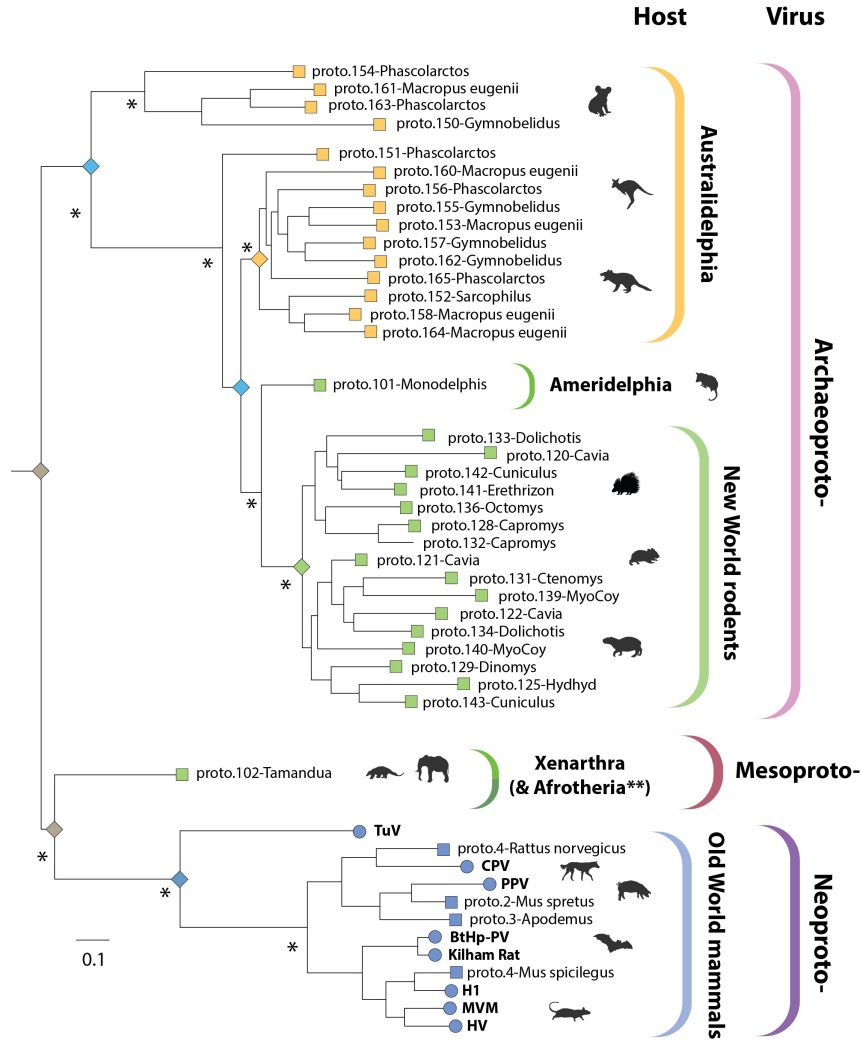
61. Kumar, S., et al., *TimeTree: A Resource for Timelines, Timetrees, and Divergence Times*. Mol Biol Evol, 2017. **34**(7): p. 1812-1819.
62. Clark, J.R., et al., *A comparative study in ancestral range reconstruction methods: retracing the uncertain histories of insular lineages*. Syst Biol, 2008. **57**(5): p. 693-707.
63. Chamberlain, S., et al., *rgbif: Interface to the Global Biodiversity Information Facility API*. 2021.
64. Ree, R.H. and S.A. Smith, *Maximum likelihood inference of geographic range evolution by dispersal, local extinction, and cladogenesis*. Syst Biol, 2008. **57**(1): p. 4-14.
65. Wickham, H., *ggplot2: Elegant Graphics for Data Analysis*. 2016.
66. Yu, G., *Using ggtree to Visualize Data on Tree-Like Structures*. Curr Protoc Bioinformatics, 2020. **69**(1): p. e96.
67. Suyama, M., D. Torrents, and P. Bork, *PAL2NAL: robust conversion of protein sequence alignments into the corresponding codon alignments*. Nucleic Acids Res, 2006. **34**(Web Server issue): p. W609-12.



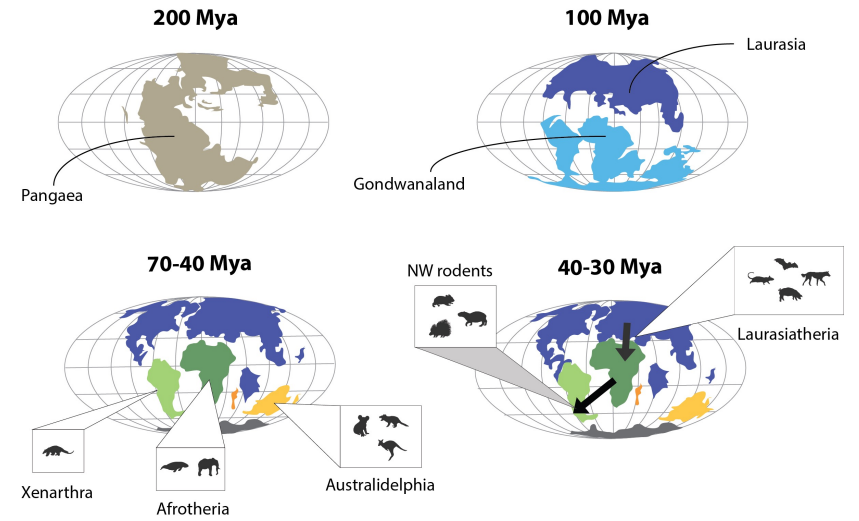




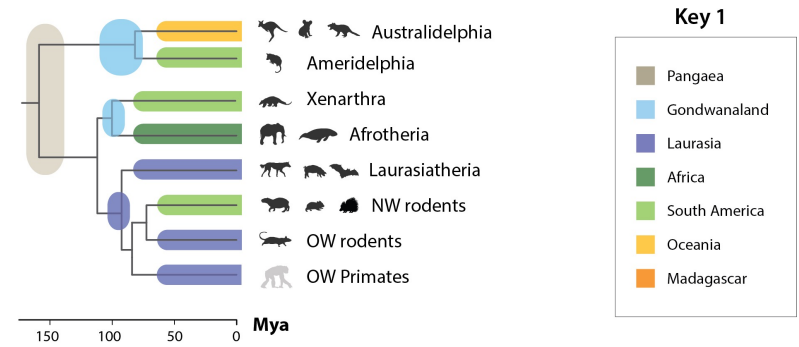
A)



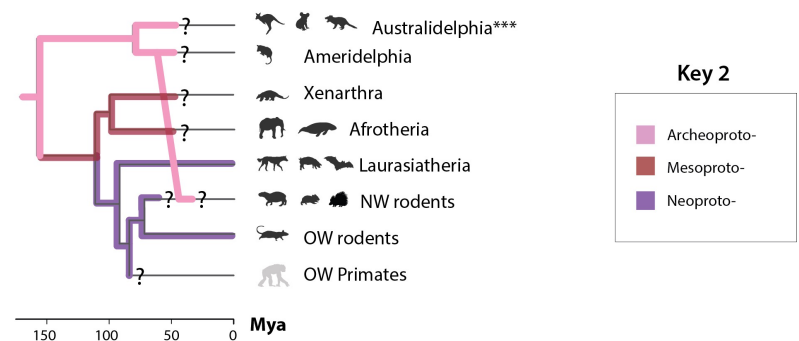
B)



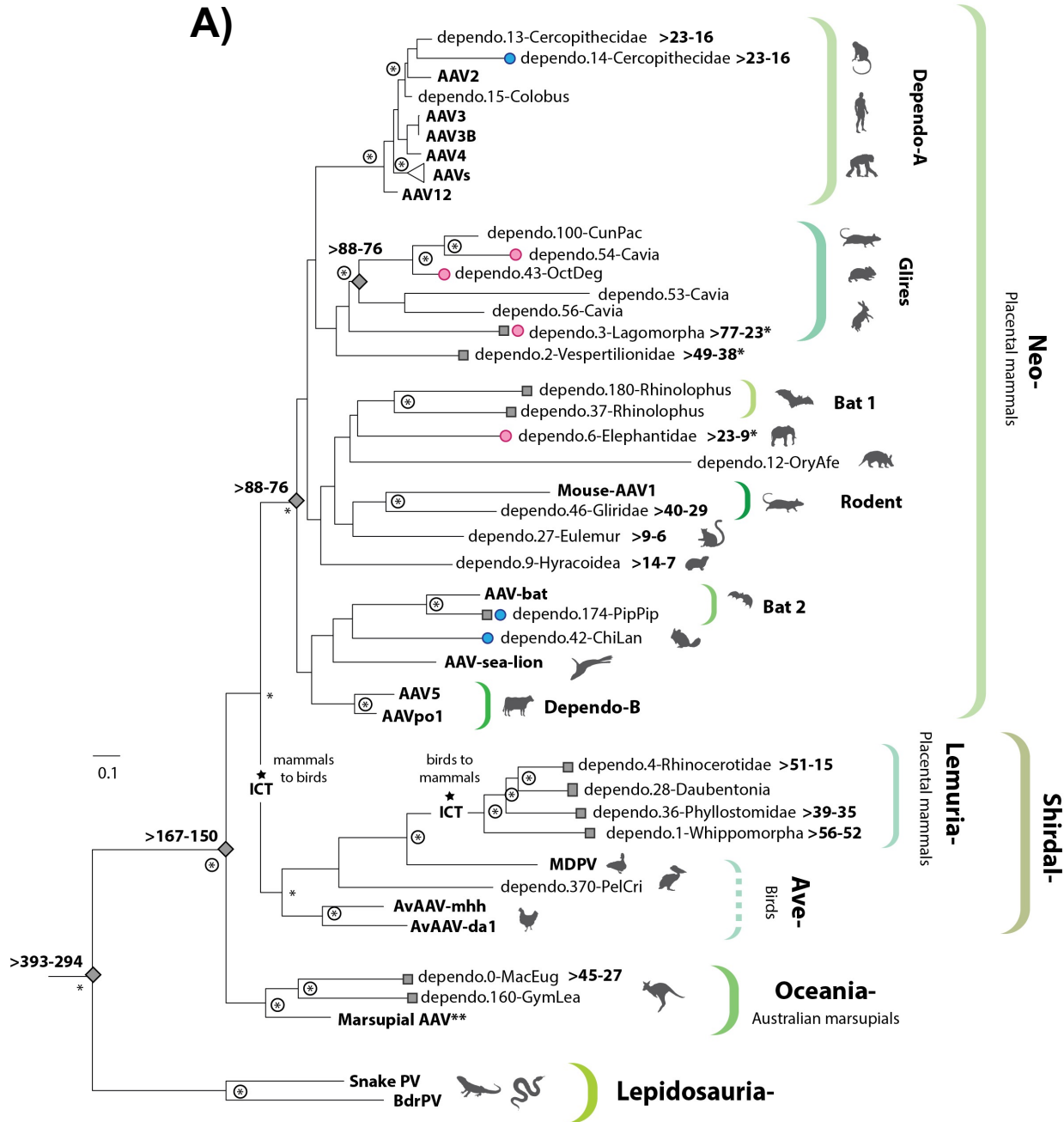
C)



D)

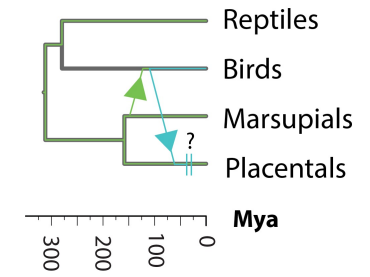


Virus phylogeny



B)

Host phylogeny



Key

- partially intact
- expressed & intact
- full-length
- ◆ internal node calibration
- ⊗ Bootstrap >70 in H.T. tree (shown)
- * Bootstrap >70 in L.T. tree
- ★ ICT Inter-class transmission