



1 **ABSTRACT**

2

3 **Background:**

4 Genomic regions that remain poorly understood, often referred to as the "dark genome,"  
5 contain a variety of functionally relevant and biologically informative genome features. These  
6 include endogenous viral elements (EVEs) - virus-derived sequences that can dramatically  
7 impact host biology and serve as a virus "fossil record". In this study, we introduce a  
8 database-integrated genome screening (DIGS) approach to investigating the dark genome  
9 *in silico*, focusing on EVEs found within vertebrate genomes.

10

11 **Results:**

12 Using DIGS on 874 vertebrate species genomes, we uncovered approximately 1.1 million  
13 EVE sequences, with over 99% originating from endogenous retroviruses or transposable  
14 elements that contain EVE DNA. We show that the remaining 6038 sequences represent  
15 over a thousand distinct horizontal gene transfer events across ten virus families, including  
16 some that have not previously been reported as EVEs. We explore the genomic and  
17 phylogenetic characteristics of non-retroviral EVEs and determine their rates of acquisition  
18 during vertebrate evolution. Our study uncovers novel virus diversity, broadens knowledge of  
19 virus distribution among vertebrate hosts, and provides new insights into the ecology and  
20 evolution of vertebrate viruses.

21

22 **Conclusions:**

23 We comprehensively catalogue and analyse EVEs within 874 vertebrate genomes, shedding  
24 light on the distribution, diversity and long-term evolution of viruses, and revealing their  
25 extensive impact on vertebrate genome evolution. Our results demonstrate the power of  
26 linking a relational database management system to a similarity search-based screening  
27 pipeline for *in silico* exploration of the dark genome.

## 1 INTRODUCTION

2 The availability of whole genome sequence (WGS) data from a broad range of species  
3 provides unprecedented scope for comparative genomic investigations [1-3]. However,  
4 these investigations rely to a large extent on *annotation* - the process of identifying and  
5 labelling genome features - which usually lags far behind the generation of sequence data.  
6 Consequently, most whole genome sequences are comprised of DNA that is incompletely  
7 understood in terms of its evolutionary origins and functional significance. The portion of  
8 sequenced genome space that lacks annotations is sometimes referred to as the 'dark  
9 genome' [4], and contains a wide variety of yet-to-be-characterized genome features. Some  
10 of these may have functional roles, such as encoding proteins [5] or regulating gene  
11 expression [6]. Others, such as non-expressed pseudogenes, may not, but can nonetheless  
12 provide valuable insights into genome biology and evolution.

13

14 Within the dark genome, endogenous viral elements (EVEs) constitute a particularly  
15 intriguing group of genome features. EVEs are virus-derived DNA sequences that become  
16 integrated into the germline genome of host species and are stably inherited as host alleles  
17 – a form of horizontal gene transfer [7-14]. While once considered genetic 'junk,' it has  
18 become evident over recent years that EVEs can profoundly impact host biology and  
19 genome evolution, with many now known to have physiologically relevant roles [15-19]. In  
20 addition, EVE sequences (whether functional or not) provide a rare source of retrospective  
21 information about ancient viruses, akin to a viral 'fossil record' [7, 20-22].

22

23 Identifying genome features contained within the dark genome, such as EVEs, often relies  
24 on the use of sequence similarity searches, such as those implemented in the Basic Local  
25 Alignment Search Tool (BLAST) [23, 24], to search WGS databases. Because sequence  
26 similarity reflects homology (evolutionary relatedness), novel genome features can often be  
27 identified based on their resemblance to ones that have been described previously. One  
28 example of this approach is implemented in the PSI-BLAST [5] and HMMER [8] programs, in

1 which iterated search strategies are used to progressively increase sensitivity so that novel  
2 homologs of previously characterised genes may be detected. A related approach is  
3 ‘systematic *in silico* genome screening’ which extends the basic concept of a similarity  
4 search in two ways: (i) inclusion of multiple query sequences and/or target databases; (ii)  
5 similarity-based classification of matching sequences (‘hits’) via comparison to reference  
6 sequence library (**Fig. 1a**). Hits may also be further investigated using additional  
7 comparative or experimental approaches (**Fig. 1b, Table 1**). Thus, screening can provide  
8 one component of a broader analytical pipeline.

9

10 While straightforward in principle, *in silico* genome screening is computationally expensive  
11 and can be difficult to implement efficiently. Moreover, large-scale screens can produce  
12 copious output data that are difficult to manage and interpret without an appropriate  
13 analytical framework. To address these issues, we developed a database-oriented approach  
14 to *in silico* screening, called *database-integrated genome screening* (DIGS). To demonstrate  
15 the use of this approach, we first created an open software framework for performing it, then  
16 used this framework to search published vertebrate genomes for EVE loci. Besides  
17 demonstrating that DIGS provides a powerful, flexible approach for exploring the dark  
18 genome, our analysis provides a comprehensive and detailed overview of EVE diversity in  
19 vertebrate genomes and reveals new information about the long-term evolutionary  
20 relationships between viruses and vertebrate hosts.

21

## 22 **RESULTS**

### 23 **1. A database-integrated approach to exploring the dark genome**

24 We developed a robust, database-integrated approach to systematic *in silico* genome  
25 screening, referred to as database-integrated genome screening (DIGS). This approach  
26 integrates a similarity search-based screening pipeline with a relational database  
27 management system (RDBMS) to enable efficient exploration of the dark genome. The  
28 rationale for this integration is twofold: it not only provides a solid foundation for conducting

1 large-scale, automated screens in an efficient and non-redundant manner but also allows for  
2 the structured querying of screening output using SQL, a powerful and well-established tool  
3 for database interrogation [25]. Additionally, an RDBMS offers advantages such as data  
4 recoverability, multi-user support, and networked data access.

5

6 The DIGS process comprises three key input data components:

7 **Target Database (TDb):** A collection of whole genome sequence assemblies (or other large  
8 sequence datasets such as transcriptomes) that will serve as the target for sequence  
9 similarity searches.

10 **Query Sequences (Probes):** A set of sequences to be used as input for similarity searches  
11 of the TDb.

12 **Reference Sequence Library (RSL):** The RSL represents the broad range of genetic  
13 diversity associated with the genome feature(s) under investigation. Its composition varies  
14 according to the analysis context (see **Table 1**). It should always include sequences  
15 representing diversity within the genome feature under investigation. It may also include  
16 genetic marker sequences and potentially cross-matching genome features. Probes are  
17 typically a subset of sequences contained in the RSL.

18

19 As illustrated in **Fig. 2**, the DIGS process involves systematic searching of a user-defined  
20 TDb with user-defined probes, merging fragmented hits, and classifying merged sequences  
21 through BLAST-based comparison to the RSL. The output - a set of non-redundant,  
22 defragmented 'hits' - is captured in a project-specific relational database. Importantly, this  
23 integration allows database queries to be employed in real time, with SQL queries  
24 referencing any information captured by the database schema. SQL-based querying of  
25 screening databases facilitates the identification of loci of interest, which can then be  
26 explored further using comparative approaches (see **Fig. 1b**).

27

1 It is important to note that screening is usually an iterative discovery process, wherein initial  
2 results inform the development of subsequent screens. For instance, novel diversity  
3 detected by an initial screen can subsequently be incorporated into the RSL, and hits within  
4 the screening database can be reclassified using the updated library (**Fig. 2**). Additionally,  
5 probe sets used in initial searches can be expanded to incorporate sequences identified  
6 during screening, broadening the range of sequences detected in subsequent screens [26].  
7 However, care must be taken when using this approach, since it can potentially produce  
8 misleading results, or generate excessive hits (e.g. if highly repetitive sequences are  
9 contained within the new probes). Importantly, database integration allows results to be  
10 observed and interrogated in real time - as they are being generated. This means that  
11 configuration issues (e.g. badly composed RSL, inappropriate choice of probes) can be  
12 detected early on – potentially saving a significant amount of time and effort. Furthermore, it  
13 facilitates the implementation of agile, heuristic screening strategies, in which approaches  
14 are adjusted in line with results.

15

## 16 **2. An open software framework for implementing DIGS**

17 We constructed a software framework for implementing DIGS, called ‘the DIGS tool’. The  
18 DIGS tool is implemented using the PERL scripting language. It uses the BLAST+ program  
19 suite [24] to perform similarity searches, and the MySQL RDBMS (to capture their output).  
20 Accessible through a text-based console interface, it simplifies the complex process of large-  
21 scale genome screening, and provides a versatile basis for implementing screens.

22

23 To initiate screening using the DIGS tool, researchers provide a project-specific command  
24 file (**Fig. S1**) that serves as the blueprint for the screening process. This command file  
25 specifies parameters, including the user-defined name of the screening database, and file  
26 paths to the TDb, RSL, and probe sequences. When a screen is initiated a project-specific  
27 database is created with the schema shown in **Fig. S2**. This core schema can be extended  
28 to include any relevant "side data" – e.g., taxonomic information related to the species and

1 sequences included in the screen - increasing the power of SQL queries to reveal  
2 informative patterns (**Fig. S2, Fig. S3**).

3

4 Systematic screening proceeds automatically until all searches have been completed. If the  
5 process is interrupted at any point, or if novel probe/target sequences are incorporated into  
6 the project, screening will proceed in a non-redundant way on restarting. Thus, screening  
7 projects can readily be expanded to incorporate new TDb files (e.g. recently published WGS  
8 assemblies), or novel probe/reference sequences, as they become available. The DIGS tool  
9 console allows reclassification of sequences held in the results table (e.g. following an RSL  
10 update). To increase efficiency, this process can be tailored to specific subsets of database  
11 sequences by supplying SQL constraints via the DIGS tool console.

12

13 BLAST algorithms emphasise local similarity and consequently tend to fragment contiguous  
14 matches into several separate hits if similarity across some internal regions of the match is  
15 low. The DIGS tool allows screening pipelines to be configured with respect to how  
16 overlapping/adjacent hits are handled, so that the screening process can be tailored to the  
17 specific needs of diverse projects. The DIGS tool also provides a 'consolidation' function that  
18 concatenates, rather than merges, adjacent hits and stores concatenated results, along with  
19 information about their structure, in a new screening database table.

20

21 For program validation, we mined mammalian genomes for sequences disclosing similarity  
22 to the antiviral restriction factor tetherin [27, 28]. Tetherin provides a useful test case as it is  
23 a relatively unique gene and its evolution, distribution and diversity have previously been  
24 examined [27, 28]. Results were compared with those provided by two alternative, widely  
25 used genome mining pipelines: OrthoDB [29] and Ensembl [30] and found to overlap by  
26 >99% (**Fig. S4**).

27

1 The DIGS tool provides functionality for exporting FASTA-formatted sequences and  
2 managing screening database tables (e.g., add/drop tables, import table data). Further  
3 information regarding program installation and usage is provided online, in a repository  
4 associated website [31]. In the sections below we illustrate the application of the DIGS tool  
5 to cataloguing of EVEs in vertebrate genomes, focussing on both high and low copy number  
6 elements.

7

### 8 **3. Use of DIGS to catalogue RT-encoding endogenous retroviruses**

9 Unlike other vertebrate viruses, retroviruses (family *Retroviridae*) integrate their genome into  
10 the nuclear genome of infected cells as an obligate part of their life cycle. As a result,  
11 retroviruses gain more opportunities to become a permanent part of the host germline.  
12 Furthermore, the initial integrated form of a retrovirus genome, called a provirus, is typically  
13 replication competent and increases in germline copy number can thus occur through  
14 reinfection of germ line cells [32]. Accordingly, 'endogenous retroviruses' (ERVs) are by far  
15 the most common type of EVE found in vertebrate genomes [7, 33].

16

17 Retrovirus genomes contain a *pol* coding domain that encodes a reverse transcriptase (RT)  
18 gene. The RT gene can be used to reconstruct phylogenetic relationships across the entire  
19 *Retroviridae* and hence provides the lynchpin for unravelling the evolutionary history and  
20 origins of ERV loci [34, 35]. We therefore implemented a screening procedure to detect RT-  
21 encoding ERV loci, based on an RSL comprised of previously classified exogenous  
22 retrovirus and ERV RT sequences (**see Methods**). Screening involved more than 1.5 million  
23 discrete tBLASTn searches and resulted in the identification of 1,073,137 ERV RT hits. This  
24 set was filtered based on higher BLAST bitscore cut-off to obtain a high confidence set of  
25 500,701 loci (**Table 2**).

26

27 High confidence ERV RT hits were identified in all vertebrate classes. However, the  
28 frequency among classes was found to vary dramatically (**Fig. 3**). ERVs occur most



1 frequently in mammals (class Mammalia) and amphibians (class Amphibia), and at relatively  
2 similar, intermediate frequencies in the genomes of reptiles (class Squamata) and birds  
3 (class Aves). By contrast, RT-encoding ERVs are relatively rare in the genomes of fish,  
4 including ray-finned fish (class Actinopterygii) and jawless fish (class Agnatha). Cartilaginous  
5 fish (class Chondrichthyes) represent a possible exception, although only a few genomes  
6 were available for this group (**Fig. 3**). These findings are broadly consistent with previous  
7 studies, conducted using a smaller amount of species genomes [33, 36-38].

8  
9 ERVs have been taxonomically grouped into three clades (I, II and III) based on their  
10 phylogenetic relatedness in the RT gene to the exogenous *Gammaretrovirus*, *Betaretrovirus*  
11 and *Spumavirus* genera respectively [1,2]. We incorporated into our RT screening database  
12 taxonomic information for (i) host species examined in our screen and (ii) RSL RT  
13 sequences. We then used SQL statements referencing these tables to summarise the  
14 frequency of clade I, II and III ERVs in distinct vertebrate classes. Whereas clade I and III  
15 ERVs are found in all vertebrate groups, clade II ERVs appear to have a more restricted  
16 distribution, occurring only at low frequency in amphibians, and being completely absent  
17 from agnathans and cartilaginous fish (**Table 2**). A few clade II ERVs were identified in ray-  
18 finned fish, but these were very closely related to mammalian ERVs and likely represent  
19 contamination of WGS builds with mammalian genomic DNA. While RT-encoding ERV copy  
20 number is relatively high in cartilaginous fish, RT diversity is relatively low, with the majority  
21 of ERV RT sequences belonging to clade III.

22

#### 23 **4. Use of DIGS to catalogue non-retroviral EVEs vertebrate genomes.**

24 To identify non-retroviral EVEs, we first obtained an RSL representing all known viruses [39].  
25 From this library, a set of representative probes was selected. Probes include representative  
26 proteomes of all known vertebrate viruses except retroviruses. Screening entailed >1.5  
27 million discrete tBLASTn searches, and initial results comprised 33,654 hits. However, many  
28 of these represented matches to host genes and TEs. We identified these spurious matches

1 by interrogating screening results with a combination of SQL queries and *ad hoc*  
2 phylogenetic analysis. We also identified and excluded hits that appeared likely to derived  
3 from exogenous viruses (see **Table S1**). For example, SQL-generated summaries of our  
4 initial screen results revealed several WGS sequences disclosing unexpected similarity to  
5 plant virus genomes (**Fig. S3a**). Among these, matches to geminiviruses (family  
6 *Geminiviridae*) and potyviruses (family *Potyviridae*) lack evidence for germline integration  
7 and likely to represent diet-related contamination. Other unusual matches were contained  
8 within large contigs and thus could not be dismissed as contaminating DNA but were  
9 revealed to be spurious by genomic and phylogenetic analysis. For example, a sequence  
10 identified in the genome of the pig-nosed turtle (*Carettochelys insculpta*) disclosed similarity  
11 to caulimoviruses (family Caulimoviruses) - but was revealed by closer analysis to represent  
12 an unusual ERV (**Fig. S3a, Fig. S5**).

13

14 Next we removed matches to transposons partly comprised of virus-derived DNA such as  
15 the adintovirus-derived mavericks [40] and alloherpesvirus-derived teratorns [41] (**Fig. S3a**).  
16 Once TEs had been removed, results comprised 6038 putative non-retroviral EVE  
17 sequences, representing 10 virus families (**Table 3**, [42]). We did not identify any EVEs  
18 derived from vertebrate viruses with genomes comprised of double-stranded RNA (e.g.,  
19 order Reovirales) or circular single-stranded RNA (e.g., genus *Deltavirus*). However, all  
20 other virus genome 'classes' were represented including reverse-transcribing DNA (DNART)  
21 viruses, double-stranded DNA (DNA ds) viruses, single-stranded DNA (DNA ss) viruses,  
22 single-stranded negative sense RNA (RNA ss-ve) viruses, and single-stranded positive  
23 sense RNA (RNA ss+ve) viruses. Plotting the distribution of EVEs and exogenous viruses  
24 from distinct virus families and genera across vertebrate phyla, revealed that many virus  
25 groups have had a broader distribution across vertebrate hosts than recognised on the basis  
26 of exogenous isolates (**Fig. 4**).

27

1 We examined all EVE loci identified in our study to determine their coding potential. We  
2 identified numerous EVE loci encoding open reading frames (ORFs) >300 amino acids (aa)  
3 in length (**Fig. S6**). Among these, four encoded ORFs longer than 1000 aa. One of these – a  
4 1718aa ORF encoded by an endogenous borna-like L-protein (EBLL) element in bats  
5 (EBLL-Cultervirus.29-EptFus) – has been reported previously [43]. However, we also  
6 identified an endogenous chuvirus-like L-protein (ECLL) element encoding an ~1400 aa  
7 ORF in livebearers (subfamily Poeciliinae). This element encodes long ORFs in two distinct  
8 livebearer species (*P. formosa* and *P. latapina*), indicating it's coding capacity has been  
9 conserved for >10 million years [44]. We also detected herpesvirus and alloherpesvirus  
10 EVEs encoding ORFs >1000 aa, but as discussed below, the integration status of these  
11 sequences remains unclear.

12

## 13 **5. Diversity of non-retroviral EVEs in vertebrate genomes**

### 14 5.1 EVEs derived from viruses with double-stranded DNA genomes

15 We detected DNA derived from herpesviruses (family *Herpesviridae*) in mammalian and  
16 reptilian genomes (**Fig. 4, Table 3**, [42]). DNA sequences derived from betaherpesviruses  
17 (subfamily *Betaherpesvirinae*) and gammaherpesviruses (subfamily *Gammaherpesvirinae*)  
18 have previously been reported in WGS assemblies of the tarsier (*Carlito syrichta*) and aye-  
19 aye (*Daubentonia madagascensis*), respectively [45]. In addition to these sequences, we  
20 detected gammaherpesvirus DNA in WGS data of red squirrels (*Sciurus vulgaris*) and the  
21 Amazon river dolphin (*Inia geoffrensis*), while betaherpesvirus DNA was detected in the  
22 stoat (*Mustela ermina*) WGS assembly, and DNA derived from an alphaherpesvirus  
23 (subfamily *Alphaherpesvirinae*) in the Oriximina lizard (*Tretioscincus oriximinensis*) WGS  
24 (**Fig. S7a-b**). Germline integration of human betaherpesviruses has been demonstrated [46,  
25 47], and the presence of a betaherpesvirus-derived EVE in the tarsier genome EVE has  
26 been established [45]. However, herpesviruses can also establish latent infections, and  
27 might be considered likely to occur as contaminants of DNA samples used to generate  
28 whole genome sequence assemblies. Due to the limitations of the WGS assemblies in which

1 they were identified, it was not possible to confirm that the novel herpesvirus DNA  
2 sequences detected here represent EVEs rather than DNA derived from contaminating  
3 exogenous viruses.

4  
5 DNA derived from alloherpesviruses (family *Alloherpesviridae*) was detected in fish and  
6 amphibians. In ray-finned fish, most of these sequences belonged to the 'teratorn' lineage of  
7 transposable elements, which have arisen via fusion of alloherpesvirus genomes and  
8 piggyBac transposons, and have been intragenomically amplified in the genomes of teleost  
9 fish (Infraclass Teleostei) [41]. Additional alloherpesvirus-related elements were identified in  
10 three amphibian species and five ray-finned fish species [42]. One of these elements,  
11 identified in the Asiatic toad (*Bufo gargarizans*) occurred within a contig that was significantly  
12 larger than a herpesvirus genome, demonstrating that it represents an EVE rather than an  
13 exogenous virus. Phylogenetic analysis revealed that alloherpesvirus-like sequences  
14 identified in amphibian genomes clustered robustly with amphibian alloherpesviruses, while  
15 those identified in fish genomes clustered with fish alloherpesviruses (**Fig. S7c**).

16

#### 17 5.2 EVEs derived from viruses with single-stranded DNA genomes

18 EVEs derived from parvoviruses (family *Parvoviridae*) and circoviruses (family *Circoviridae*)  
19 are widespread in vertebrate genomes, being found in the majority of vertebrate classes  
20 (**Fig. 4**). Both endogenous circoviral elements (ECVs) and endogenous parvoviral elements  
21 (EPVs) are only absent in major vertebrate groups represented by a relatively small number  
22 of sequenced species genomes (i.e. between 1 and 6). No ECVs or EPVs were identified in  
23 the tuatara (order Rhynchocephalia) or in crocodiles (order Crocodylia). EPVs were not  
24 identified in agnathans, while ECVs were not identified in cartilaginous fish.

25

26 We identified a total of 1192 ECVs, most of which derived from an element in carnivore  
27 (Class Mammalia: order Carnivora) genomes that is embedded within a non-LTR  
28 retrotransposon and has undergone intragenomic amplification (**Fig. S8**). While many of the

1 ECVs identified in our screen have been reported in previous publications [7, 26, 48-50], we  
2 also identified novel loci in mammals, reptiles, amphibians, and ray-finned fish [42].  
3 Phylogenetic analysis (see **Fig. S7d**) revealed that a novel ECV locus in turtles was found to  
4 group with avian circoviruses, while amphibian ECV elements grouped with fish circoviruses,  
5 though bootstrap support for this relationship was lacking. A circovirus-like sequence  
6 detected in the WGS data of Allen's wood mouse (*Hylomyscus alleni*) grouped robustly with  
7 the exogenous 'rodent circovirus', but integration of this sequence into the *H. alleni* genome  
8 could not be confirmed.

9

10 We identified 627 EPVs, representing two distinct subfamilies within the *Parvoviridae* and  
11 five distinct genera (see **Fig. 4**). The majority of these loci have been reported in a previous  
12 study of vertebrate genomes [50] or were related to these loci. However, we also identified  
13 novel EPVs in reptiles, amphibians and mammals (**Table 3**, [42]). In reptiles the novel  
14 elements derived from genus *Dependoparvovirus* while the amphibian elements were more  
15 closely related to viruses in genus *Protoparvovirus*. Interestingly, these EPVs clustered  
16 basally within a clade of protoparvovirus-related viruses in phylogenetic reconstructions (**Fig.**  
17 **S7e**), consistent with previous analyses indicating that this genus may have broadly co-  
18 diverged with vertebrate phyla [50].

19

### 20 5.3 EVEs derived from reverse-transcribing DNA viruses

21 EVEs derived from hepadnaviruses (family *Hepadnaviridae*), which are reverse-transcribing  
22 DNA viruses, were identified in reptiles, birds and amphibians (**Table 3**, [42]). Most of these  
23 EVEs, which are commonly referred to as 'endogenous hepatitis B viruses' (eHBVs), have  
24 been reported previously [51, 52]. However, we identified novel elements in the plateau  
25 fence lizard (*Sceloporus tristichus*), and also in vertebrate classes where they have not been  
26 reported previously. These include one eHBV identified in a cartilaginous fish, the Australian  
27 ghostshark (*Callorhynchus milii*), and another identified in an amphibian, the common coquí  
28 (*Eleutherodactylus coqui*).

1

2 Phylogenetic analysis (see **Fig. S7f**) revealed that novel eHBV elements identified in lizards  
3 (suborder Lacertilia) group robustly with the exogenous skink hepadnavirus (SkHBV), while  
4 the amphibian element groups with robustly within a clade comprised of the exogenous  
5 spiny lizard hepadnavirus (SIHBV), Tibetan frog hepadnavirus (TfHBV) and eHBV elements  
6 identified in crocodile genomes. The eHBV identified in sharks was relatively short and not  
7 amenable to phylogenetic analysis, but nonetheless provides the first evidence that  
8 hepadnaviruses infect this host group.

9

#### 10 5.4 EVEs derived from viruses with single-stranded, negative sense RNA genomes

11 Screening revealed that vertebrate genomes contain numerous EVEs derived from  
12 mononegaviruses (order *Mononegvirales*), which are characterised by non-segmented  
13 ssRNA-ve genomes. These EVEs derive from four mononegavirus families: bornaviruses  
14 (family *Bornaviridae*), filoviruses (family *Filoviridae*), paramyxoviruses (family  
15 *Paramyxoviridae*) and chuviruses (family *Chuviridae*) (**Fig. 4, Table 3**, [42]). We did not  
16 detect any EVEs derived from other mononegavirus families that infect vertebrates  
17 (*Pneumoviridae*, *Rhabdoviridae*, *Nyamiviridae*, *Sunviridae*), nor any EVEs derived from virus  
18 families with segmented, negative sense RNA genomes (e.g., *Peribunyaviridae*,  
19 *Orthomyxoviridae*).

20

21 The majority of mononegavirus EVEs identified in our screen were derived from  
22 bornaviruses and filoviruses and have been described in previous reports [7, 48, 50, 51, 53].  
23 However, we also identified novel EVEs derived from these groups, as well as previously  
24 unreported EVEs derived from paramyxoviruses and chuviruses (**Table 3**). Germline  
25 integration of DNA derived from mononegaviruses can occur if, in an infected germline cell,  
26 viral mRNA sequences are reverse transcribed and integrated into the nuclear genome by  
27 cellular retroelements [54]. EVE loci generated in this way preserve the sequences of  
28 individual genes of ancient mononegaviruses, but not entire viral genomes. Among

1 mononegavirus-derived EVEs, regardless of which family, EVEs derived from the  
2 nucleoprotein (NP) and large polymerase (L) genes predominate, but other genes are also  
3 represented, including the glycoprotein (GP) genes of filoviruses, bornaviruses, and  
4 chuviruses, the VP30 and VP35 genes of filoviruses, and the hemagglutinin-neuraminidase  
5 (HA-NM) gene of paramyxoviruses.

6

7 Paramyxovirus-like EVEs were identified in ray-finned fish, amphibians and sharks (**Fig. 4,**  
8 **Table 3**, [42]). Many of these EVEs were highly divergent and/or degenerated and  
9 consequently their evolutionary relationships to contemporary paramyxoviruses were poorly  
10 resolved in phylogenetic analysis. However, an L polymerase-derived sequence identified in  
11 the pobblebonk frog (*Limnodynastes dumerilii*) genome was found to group robustly with  
12 sunshine virus, a contemporary paramyxovirus of Australian pythons [55] in phylogenetic  
13 trees (**Fig. S7g**).

14

15 Chuvirus-like sequences were identified in agnathans, ray-finned fish, reptiles, and  
16 mammals (**Fig. 4, Table 3**, [42]). The majority of the mammalian elements were identified in  
17 marsupials, but we also identified a single chuvirus-like EVE in the genome of a  
18 laurasiatherian mammal – the bottlenose dolphin (*Tursiops truncatus*). Phylogenetic trees  
19 reconstructed using alignments of NP-derived chuvirus EVEs and NP genes of  
20 contemporary chuviruses revealed evidence for the existence of distinct clades specific to  
21 particular vertebrate classes (**Fig. S7h**). These included a clade including both a snake EVE  
22 and an exogenous chuviruses of snakes, and two clades comprised of EVEs and viruses of  
23 teleost fish. In addition, these phylogenies revealed a robustly supported relationship  
24 between chuvirus EVEs in the Tibetan frog (*Nanorana parkeri*) and zebrafish (*Danio rerio*)  
25 genomes. Taken together, these results provide evidence for the existence of numerous  
26 diverse lineages of chuviruses in vertebrates, adding to recent evidence for the presence of  
27 exogenous chuviruses in marsupials [56].

28

1 Filovirus-derived EVEs were mainly identified in mammals (**Fig. 4, Table 3, [42]**). However,  
2 we also identified one filovirus-derived EVE in an amphibian – the mimic poison frog  
3 (*Ranitomeya imitator*) - providing the first evidence that filoviruses infect this vertebrate  
4 group (**Table 1**). We identified novel, ancient filovirus EVEs in anteaters (family  
5 Myrmecophagidae) and spiny mice (genus *Acomys*).

6

7 Strikingly, the inclusion of Tapajos virus (TAPV), a snake filovirus, in phylogenetic  
8 reconstructions revealed evidence for the existence of two highly distinct filovirus lineages in  
9 mammals (**Fig. 5**). These two lineages, which are robustly separated from one another by  
10 TAPV, are evident in phylogenies constructed for both the NP and VP35 genes. One lineage  
11 (here labelled ‘Mammal-1’) is comprised of EVEs and all contemporary mammalian  
12 filoviruses, whereas the other (‘Mammal-2’) is comprised exclusively of EVEs. Notably,  
13 within the Mammal-1 group, EVEs identified in mammals indigenous to Southern  
14 Hemisphere continents (e.g. marsupials, xenarthrans) group basally, whereas EVEs and  
15 viruses isolated from ‘Old World’-associated placental mammals occupy a more derived  
16 position.

17

18 The ‘Mammal-2’ clade contains filovirus EVEs from rodents, primates and bats. Because  
19 EVEs belonging to this clade were obtained from several distinct lineages, and show  
20 conservation across these groups, we can be reasonably confident they represent a *bona*  
21 *fide* lineage within the *Filoviridae*, rather than just a set of highly degraded filo-like EVEs that  
22 group together due to long branch attraction [57]. One member of this group (eflp-filo.1-  
23 Myotis) encodes an intact VP35 protein, the properties of which have been experimentally  
24 investigated in recent studies [58, 59]. Interestingly, we found that spiny mice also harbour a  
25 filovirus EVE encoding an intact VP35 protein (eflp-filo.3-Acomys), however, this insertion  
26 belongs to the ‘Mammal 1’ clade and is relatively closely related to the VP35 proteins found  
27 in contemporary mammalian filoviruses (**Fig. 5b**).

28



1 Bornavirus-like EVEs were identified in all vertebrate classes except Chondrichthyes (**Fig. 4,**  
2 **Table 3,** [42]). The majority have been reported previously or are orthologs of previously  
3 reported EVEs. However, we identified novel bornavirus-like EVEs in the genomes of ray-  
4 finned fish and amphibians. The amphibian EVEs grouped robustly with culterviruses in  
5 phylogenetic reconstructions (**Fig. S7i-j**).

6

#### 7 5.5 EVEs derived from viruses with single-stranded, positive sense RNA genomes

8 EVEs derived from positive sense RNA viruses are rare in vertebrate genomes (**Fig. 4,**  
9 **Table 3,** [42]). The only examples we identified were a small number of sequences derived  
10 from flavivirids (family *Flaviviridae*). These include an EVE derived from the *Pestivirus* genus  
11 of flavivirids, the reference genome of the Indochinese shrew (*Crocidura indochinensis*), as  
12 reported previously [60], and EVEs identified in ray-finned fish, also reported previously [61].  
13 In fish genomes, flavivirid EVEs derive from the proposed ‘Tamanavirus’ genus, and a  
14 lineage labelled ‘X2’ that groups as a sister taxon to the proposed ‘Jingmenvirus’ genus.  
15 However, jingmenviruses are actually segmented, RNAss-ve viruses whose genomes  
16 include flavivirid-derived segments [62]. Since is possible that the X2 lineage shares a  
17 common RNAss-ve ancestor with jingmenviruses, EVEs belonging to this lineage may in fact  
18 be derived from viruses with ssRNA-ve genomes.

19

## 20 **6. Frequency of germline incorporation events across distinct vertebrate phyla**

21 We used the DIGS framework to dissect the history of horizontal gene transfer events  
22 involving germline incorporation of DNA derived from non-retroviral viruses. We excluded  
23 EVEs derived from DNAds viruses from this analysis because most of these are  
24 mavericks/polintons or teratorn elements that have undergone intragenomic amplification.  
25 For these groups, the large number of insertions, and the fact that amplified lineages appear  
26 to have been independently established on multiple occasions, meant that such an analysis  
27 would be beyond the scope of this study. Furthermore, for most of the few dsDNA-derived  
28 EVEs that did not belong to these groups, it was not possible to determine if they

1 represented germline-integrated elements, exogenous viruses, or integrations occurring in  
2 somatic cells.

3

4 To examine the rate of germline incorporation in the remaining groups of non-retroviral EVEs  
5 representing DNAss, DNArt, RNAss-ve and RNAss+ve viruses, we compiled an expanded  
6 RSL containing a single reference sequence for each putative (or previously confirmed)  
7 ortholog. By classifying our hits against this expanded RSL, we could discriminate novel  
8 EVE loci (paralogs) from orthologs of previously described EVE loci. Where novel paralogs  
9 were identified, we incorporated these into our RSL and then reclassified related sequences  
10 in our screening database against this updated library. By investigating loci in this way, and  
11 iteratively reclassifying database sequences, we progressively resolved the various non-  
12 retroviral EVEs identified in our screen into sets of putatively orthologous insertions. Via this  
13 analysis we estimated that the non-retroviral EVEs identified in our study (excluding those  
14 derived from dsDNA viruses) represent ~1137 distinct germline incorporation events (**Table**  
15 **3**). Using orthology information we calculated minimum age estimates for all non-retroviral  
16 EVEs identified in two or more species [42]. We applied standardised nomenclature to EVE  
17 loci (see **Methods**), capturing information about EVE orthology, taxonomy, and host  
18 distribution [42].

19

20 Next, we estimated the rate of germline incorporation for each endogenized virus family, in  
21 all vertebrate classes represented by at least ten species (**Fig. 6**). Rates were found vary  
22 dramatically across each of the vertebrate groups examined. Overall, rates were highest in  
23 mammals, and lowest in reptiles. Fish and amphibians disclosed similar rates with ssDNA  
24 and ssRNA-ve viruses being incorporated at similar, intermediate rates. Birds were generally  
25 similar to reptiles but show a higher rate of ssDNA virus incorporations and a markedly  
26 elevated rate of hepadnavirus incorporation. Rates of parvovirus, filovirus, and bornavirus  
27 infiltration were very high in mammals compared to other vertebrate classes, with  
28 bornaviruses being incorporated at a particularly high rate (>0.03 per million years of species

1 evolution). A relatively high rate of incorporation of RNAss+ve viruses was observed in ray-  
2 finned fish, but since the elements in question are closely related to jingmenviruses, as  
3 described above, they may in fact reflect incorporation of DNA derived from an RNAss-ve  
4 virus group [62].

5  
6 In addition to estimating the frequency of germline incorporation of non-retroviral viruses, we  
7 used our screening data to reconstruct a time-calibrated overview of virus integration  
8 throughout vertebrate evolutionary history (**Fig. 7, Table S2**). Among putatively orthologous  
9 groups of EVEs for which we were able to estimate minimum dates of integration, the  
10 majority were found to have been incorporated in the Cenozoic Era (1-66 Mya). So far, the  
11 oldest integration event identified involves a metahepadnavirus (genus *Metahepadnavirus*)-  
12 derived EVE that appears to be orthologous in tuataras and birds, indicating it was  
13 incorporated into the saurian germline >280-300 Mya (see [51]). Other ancient EVEs include  
14 circovirus and herpetohepadnavirus (genus *Herpetohepadnavirus*)-derived EVEs in turtles  
15 (order Testudines) (see [52]), a circovirus-derived EVE in frogs (order Anura), and  
16 bornavirus integrations in placental mammals (see [53]). Besides revealing the landscape of  
17 non-retroviral EVE integration throughout vertebrate history, plotting EVE distribution in this  
18 way clearly reveals the main differences in EVE distribution across host groups (**Fig. 7**).

19

## 20 **DISCUSSION**

21 Sequencing of genomes is advancing rapidly but deciphering the complex layers of  
22 information they contain is a challenging, long-term endeavour [58, 59]. Genomes are not  
23 only inherently complex, they also exhibit remarkable dynamism, with phenomena such as  
24 recombination, transposition and horizontal gene transfer contributing to the creation of  
25 genomic 'churn' that makes feature distribution difficult to map [60]. These issues, combined  
26 with rapid data accumulation, coverage limitations, and assembly errors – make generation  
27 of complete and accurate annotations difficult [62, 63]. Consequently, labour intensive

1 manual genome annotation remains important [58, 61], and most published whole genome  
2 sequences are comprised of genomic 'dark matter'.

3

4 An exciting aspect of these circumstances is that they provide immense scope to make  
5 interesting biological discoveries using low cost, *in silico* approaches. While experimental  
6 studies are generally required to characterise genome features at a functional level,  
7 approaches based solely on comparative sequence analysis (see **Fig 1b**) can often reveal  
8 useful insights into their biology and evolution [1, 63]. Furthermore, comparative  
9 investigations *in silico* can often be productively combined with functional genomics or  
10 experimental approaches (**Fig 1b, Table 1**).

11

12 Systematic *in silico* genome screening is computational approach that facilitates  
13 investigation of the dark genome (**Fig. 1**). However, it can be challenging to implement  
14 efficiently. Automated pipelines are generally required to implement large-scale screens [64],  
15 and these can produce copious output data that are difficult to manage and interpret without  
16 an appropriate analytical framework. In this report, we introduce DIGS – a robust analytical  
17 platform for conducting large-scale *in silico* screens - and describe an open software  
18 framework (the DIGS tool) for implementing it.

19

20 EVEs constitute one interesting and informative group of genome features that can be found  
21 within the dark genome [22]. They are poorly annotated for several reasons. Firstly, they  
22 arise sporadically via horizontal gene transfer, and consequently their distribution is  
23 unpredictable [7, 22]. Additionally, uncharacterised EVE loci may be hard to recognise due  
24 to their being highly degraded or fragmented, or because their exogenous virus counterparts  
25 are either unknown or extinct [65, 66]. Finally, there are numerous potential sources of  
26 confounding or artefactual results that can arise during EVE screening, including host genes  
27 that exhibit similarity to virus genes, and contamination of WGS assemblies with DNA  
28 derived from exogenous viruses.

1

2 To illustrate how DIGS facilitates identification and characterisation of features hidden within  
3 the dark genome, we use the DIGS tool to perform a broad-based investigation of EVE  
4 diversity in vertebrates. We first focussed on high-copy number EVEs - which in vertebrate  
5 genomes mainly comprise ERVs. We screened 874 vertebrate genomes for RT-encoding  
6 ERVs and identified 700,000 high confidence matches. This screen revealed marked  
7 differences in ERV RT copy number between vertebrate classes. An in-depth investigation  
8 of ERV diversity in vertebrates – for example, examining their composition in finer detail, or  
9 incorporating insertions that lack RT sequences, was considered beyond the scope of this  
10 study. However, the RT dataset generated here provides a robust foundation for further ERV  
11 studies that are underpinned by phylogenetic analysis. For example, we have previously  
12 used RT data in combination with other *in silico* approaches for in-depth, phylogenetical  
13 characterisation of ERVs within discrete mammalian subgroups (e.g. see [67]).

14

15 ERVs constitute a unique type of EVE, in that they can remain replication-competent  
16 following integration and may increase their germline copy number through continued virus  
17 replication. However, the germline copy number of any EVE can potentially increase  
18 through interactions with TEs - this has been described for ERVs [68-70], as well as for  
19 EVEs derived from dsDNA viruses [40, 41, 71]. In addition, data obtained here and in our  
20 previous investigations show that EVEs derived from hepadnaviruses have been amplified in  
21 cormorants [51], while circovirus-derived sequences have been amplified in carnivore  
22 genomes [48], apparently in association with LINE1 activity [42]. These findings underline  
23 the impact of fusion between EVEs and vertebrate transposons on vertebrate genome  
24 evolution. This phenomenon has occurred on multiple independent occasions and involved a  
25 diverse range of vertebrate viruses. Interestingly, we found that circovirus EVEs in carnivore  
26 genomes are associated with a retroelement lineage (LINE1) that has also inserted into  
27 gammaherpesvirus and Chikungunya virus genomes (**Fig. S8**). These findings suggest that

1 retroelement-mediated transposition can establish a complex network of horizontal gene  
2 transfer events linking host and virus genomes.

3

4 DIGS is not only well-suited to exploring the distribution and diversity of high copy number  
5 genome features such as ERVs and TEs, it can also be used in ‘beach combing’ searches  
6 that aim to identify rare and unusual genome features. These kinds of screens typically  
7 require a rigorous filtering process to distinguish genuine from spurious matches, and as  
8 shown here, this is facilitated by database integration. DIGS enabled the efficient  
9 identification of EVEs derived from non-retroviral viruses (which are relatively rare and  
10 diverse) and provided a powerful framework for filtering spurious results (**Fig. S3**).

11

12 Via DIGS, we established a broad overview of non-retroviral EVE diversity in vertebrate  
13 genomes (**Fig. 4, Fig. 6**), shedding new light on virus distribution and diversity in  
14 vertebrates. Notably, our findings extend the known host range of important virus families.  
15 For example, we identify a filovirus-derived EVE in a frog (order Anura), providing the first  
16 evidence for the existence of amphibian filoviruses. In addition, we provide the first evidence  
17 for the presence (at least historically) of hepadnaviruses in sharks, and chuviruses in  
18 mammals (**Fig. 4**). In addition, we reveal novel virus diversity. For example, we identify novel  
19 lineages of parvoviruses and circoviruses in amphibians (**Fig S7d-e**), as well as a novel  
20 circovirus lineage in turtles (**Fig S7d**) and a novel hepadnavirus lineage in frogs (**Fig S7f**).  
21 We also identify novel paramyxovirus, chuvirus and bornavirus lineages in fish and  
22 amphibians (**Fig. S7g-j**).

23

24 Mammalian filoviruses include some of the most lethal viruses in the world [72], and while  
25 the natural reservoirs of some are known, they remain unclear for the highly pathogenic  
26 ebolavirus (EBOV) and its closest relatives (**Fig. 5**). EBOV is assumed to have a zoonotic  
27 origin, but it has rarely been possible to formally link outbreaks to a given animal reservoir,  
28 limiting understanding of its emergence. So far, efforts to identify the true reservoirs of

1 ebolaviruses have tended to focus on bats [73]. However, the widespread presence of  
2 filovirus EVEs in rodents [42], including some groups that have not been examined as  
3 potential EBOV reservoirs, such as spiny mice, suggests that the potential of this group to  
4 serve as a reservoir should not be overlooked.

5

6 Previous studies have noted that filovirus EVE sequences in the genomes of cricetid rodents  
7 (family Cricetidae) robustly split the *Ebolavirus* and *Cuevavirus* genera from the  
8 *Marburgvirus* and *Dianlovirus* genera, demonstrating that these groups diverged >20 million  
9 years ago (Mya) [74], rather than within the past 10,000 years as suggested by molecular  
10 clock-based analysis of contemporary filovirus genomes [75]. Here, we found that TAPV, an  
11 exogenous virus of snakes, robustly separates two clades of mammalian filoviruses in  
12 phylogenetic reconstructions. Since transmission of filoviruses between reptiles and  
13 mammals is likely quite rare, and both lineages contain ancient EVEs (**Fig. 5, Table S2**),  
14 these findings support the long-term existence of two highly distinct filovirus lineages  
15 ('mammal 1' and 'mammal 2') in mammals. Notably, basal taxa within the 'mammal 1'  
16 lineage – which also includes all known contemporary filoviruses of mammals – disclose  
17 associations with Southern Hemisphere continents (Australia, South America) that were  
18 largely isolated throughout extensive periods of the Cenozoic Era. These data suggest that  
19 filoviruses were present in ancestral mammals inhabiting Gondwanaland (an ancient  
20 supercontinent comprised of South America, Africa, India, and Australia) and diversified into  
21 at least two major lineages as mammalian populations became compartmentalised in  
22 distinct continental regions during the early to mid-Cenozoic. An interesting question is  
23 whether the 'mammal 2' group represents filoviruses that evolved in Northern hemisphere-  
24 associated, boreoeutherian mammals (magnorder Boreoeutheria), while 'mammal 1'  
25 represents filoviruses that initially evolved in Southern hemisphere-associated marsupials  
26 (infraclass Marsupialia) and xenarthrans (magnorder Xenarthra) before disseminating  
27 throughout the globe (possibly in association with volant mammals – i.e., bats).

28

1 While several previous studies have described EVE diversity in vertebrates [33, 36, 76], our  
2 investigation is significantly larger in scale and breadth. Furthermore, for non-retroviral  
3 viruses, we introduced a higher level of order to EVE data, making use of the DIGS  
4 framework to discriminate orthologous versus paralogous EVE loci, and to identify intra-  
5 genomically amplified EVE lineages. This allowed us to establish a panoramic view of  
6 germline incorporation by non-retroviral viruses during vertebrate evolution (**Fig. 7**).  
7 Furthermore, discriminating orthologous and paralogous EVEs enabled us to infer the rates  
8 of germline infiltration by non-retroviral virus families with greater accuracy than in previous  
9 studies (**Fig. 6, Fig. 7**). Notably, we did not find strong evidence for a reduced rate of  
10 germline infiltration in avian genomes, as suggested by a previous study [77]. Incorporation  
11 of DNART viruses is higher in birds than in any other vertebrate class (**Fig. 6**), while  
12 acquisition of EVEs derived from ssDNA-ve viruses does appear to be limited in this group,  
13 they closely resemble reptiles in this respect. Furthermore, avian hosts appear similar to  
14 reptiles with regard to ERV RT copy number (**Fig. 3**).

15

16 The absence, or near absence, of many virus groups from our catalogue of vertebrate EVEs  
17 is noteworthy. For example, vertebrates are infected with many ssRNA+ve viruses, but  
18 EVEs derived from these viruses are extremely rare, while EVEs derived from viruses with  
19 circular RNA genomes, or double-stranded RNA genomes, were not detected at all (**Table**  
20 **3**). All other virus genome types were represented by EVEs in the vertebrate germline, but  
21 their distribution is patchy and limited to a relatively small number of virus families (**Fig. 4,**  
22 **Fig 7**). For example, among ssRNA-ve viruses, only mononegaviruses were found to be  
23 present – we found no evidence for germline integration of segmented ssRNA-ve viruses  
24 such as orthomyxoviruses and bunyaviruses. The scarcity of EVEs derived from these virus  
25 groups suggests that aspects of their biology strictly limit their capacity to for germline  
26 incorporation. These likely include cell tropism (whether germline cells are typically infected)  
27 and the site of cellular replication (with viruses replicating in the nucleus more likely to be



1 incorporated). Additionally, vertebrate germline cells may present strong intrinsic barriers to  
2 the replication of certain virus groups.

3

4 The catalogue of EVE loci generated here provides a foundation for further investigations in  
5 virology, genomics, and human health. From the virology perspective, EVEs provide  
6 information about the long-term evolutionary history viruses, which greatly influences how  
7 we understand their biology. As well as enabling future studies of vertebrate 'paleoviruses',  
8 the EVE catalogue can inform efforts to identify and characterise new viruses (both by  
9 providing ecological and evolutionary insights [49], and by helping identify 'false positive' hits  
10 arising from genomic DNA). From the genomics side, EVEs are of interest due to their  
11 important roles in physiology and genome evolution [78]. These include roles antiviral  
12 immunity [11, 79, 80], and a diverse range of other physiological processes [18, 58, 59, 81-  
13 83]. Notably, we identified numerous non-retroviral EVEs encoding ORFs longer than 300aa  
14 (**Fig. S6**), indicating that their coding capacity has been conserved during vertebrate  
15 evolution. One of these - a chuvirus-derived L-protein identified in livebearers – adds to  
16 previous evidence that viral RdRp sequences have been co-opted by vertebrate genomes  
17 [43]. Mapping of EVE loci can also inform efforts to develop new medical treatments - in a  
18 recent study, EVE loci identified using DIGS were used to identify potential genomic safe  
19 harbours for human transgene therapy applications [84].

20

21 The EVE screen performed here has several important limitations. Firstly, it relied on  
22 published WGS data generated for extant species. Secondly, our results have likely been  
23 influenced by aspects of our screening configuration, such as the composition of the probe  
24 set with respect to viral taxa and polypeptide probe length [85, 86]. This might mean that we  
25 failed to detect some of the potentially recognisable EVE loci present in our TDb. For  
26 example, counts of RT-encoding ERV loci were found to be generally lower in ray-finned fish  
27 and jawless fish (**Fig. 3**), but previous studies have shown that RT loci related to other  
28 families of reverse-transcribing virus, such as metaviruses (family *Metaviridae*) [87] and

1 'lokiretroviruses' [88] are relatively common in these hosts. These would likely have been  
2 missed in our search because they were not included in our RT RSL. Finally, previous  
3 studies have indicated that vertebrate genomes contain EVEs that lack any clear homology  
4 to extant viruses [89], and these would not be detected using a sequence similarity-based  
5 approach.

6  
7 As vertebrate genome sequencing progresses, further opportunities to identify novel EVEs  
8 will arise, since: (i) any novel genome could in theory contain a lineage-specific EVE, and;  
9 (ii) ongoing characterization of exogenous virus diversity may allow for detection of  
10 previously undetectable EVEs. The DIGS project created here, which is openly available  
11 online, can be reused to accommodate newly sequenced vertebrate genomes (TDb  
12 expansion) and newly discovered vertebrate virus diversity (RSL/probe set expansion). In  
13 addition, similar projects can readily be created to screen for EVEs in other host groups.

14  
15 The use of DIGS is not limited to investigations of EVEs. DIGS can be used to investigate  
16 any sufficiently conserved genome feature lurking within the dark genome, including both  
17 coding and non-coding elements (**Table 1**). Many of the most interesting genes have  
18 evolved relatively rapidly and are difficult to annotate reliably using automated approaches  
19 [90]. Furthermore, even relatively conserved genes may be incompletely annotated by  
20 automated pipelines. DIGS has previously been used to broadly survey the distribution of  
21 interferon stimulated genes in mammals [ref], and for in-depth investigation of specific genes  
22 and gene families, such as OAS1 [91] and APOBEC3 [92]. While DIGS is best suited to  
23 investigations of genome features that comprise a single contiguous unit and contain  
24 relatively long, easily recognised regions, it can also be used to investigate genome features  
25 that are shorter or are comprised of several short sub-components, providing that a careful  
26 approach is used. For example, when investigating interferon lambda (IFNL) genes, which  
27 are expressed from multiple, short exons, we included conserved flanking features in our  
28 RSL and probe set [93]. This enabled more confident matching of IFNL exons based on their

1 positional relationships relative to conserved markers. We have also used DIGS in functional  
2 genomics studies to investigate the locations of short nucleotide motifs identified in binding  
3 assays (e.g. CHIP-seq) relative to other genomic features such as ERVs [94, 95].

4

5 The framework described here or implementing DIGS could be further developed and  
6 improved. For example, by including the option to use of other sequence similarity search  
7 tools, such as Diamond [96] and ElasticBLAST [97], and RNA structure based search tools  
8 such as INFERNAL [98]. Integrating with functional genomics resources could provide  
9 further dimensionality to the kinds of investigations that may be performed using DIGS [99].

10

## 11 **CONCLUSIONS**

12 We demonstrate how a relational database management system can be linked to a similarity  
13 search-based screening pipeline to investigate the dark genome *in silico*. Using this  
14 approach, we catalogue and analyse EVEs throughout vertebrate genomes, providing a  
15 broad range of novel insights into the evolution of ancient viruses and their interactions with  
16 host species.

## 1 **MATERIALS & METHODS**

### 2 Whole genome sequence and taxonomic data

3 Whole genome shotgun (WGS) sequence assemblies of 874 vertebrate species were  
4 obtained from the NCBI genomes resource [100]. Taxonomic data for the vertebrate species  
5 included in our screen and the viruses in our reference sequence library were obtained from  
6 the NCBI taxonomy database [101], using PERL scripts included with the DIGS tool  
7 package.

8

### 9 Database-integrated screening for RT-encoding ERVs

10 An RT RSL was collated to represent diversity within the *Retroviridae*. We included  
11 representatives of previously identified ERV lineages and exogenous retrovirus species. A  
12 subset of these sequences was used as probes in similarity search-based screens [42]. For  
13 initial screening we used a bitscore cut-off of 60. For comparisons of ERV RT copy number  
14 across species we filtered initial results using a more conservative bitscore cut-off of 90. Our  
15 previous, DIGS-based studies of ERVs have shown that spurious matches (i.e. to  
16 sequences other than retroviral RTs) do not arise when this cut-off is applied, although some  
17 genuine ERV RT hits may be excluded [67].

18

### 19 Database-integrated screening for non-retroviral EVEs

20 We obtained an RSL representing the proteome of eukaryotic viruses from the NCBI virus  
21 genomes database [39]. We supplemented this with sequence likely to cross-match to virus  
22 probes during screening. These included the teratorn transposon found in fish, which is  
23 known to contain multiple alloherpesvirus-derived genes [102]. We included the polypeptide  
24 sequences of these genes, obtained from the subtype 1 Teratorn reference (Accession #:  
25 LC199500) in our RSL. We also included representatives of the maverick/polinton lineage of  
26 transposons, derived from sequences defined in a previous study [103]. Since these element  
27 derive from a group of midsize eukaryotic linear dsDNA (MELD) viruses provisionally named  
28 'adintoviruses' [71]. Probes constituted a subset of 685 sequences contained within our

1 RSL, and incorporated polypeptide sequences representing all major protein-coding genes  
2 of representative species of all recognised vertebrate virus families. We also included  
3 representative sequences of maverick/polinton elements in our probe set. We used a bit  
4 score cut-off of 60 as a threshold for counting non-retroviral EVE loci. This threshold was  
5 established through previous experience searching for non-retroviral EVEs using DIGS [48,  
6 50, 51, 61]. Experience from previous studies had shown that nearly 100% of matches with  
7 bit scores  $\geq 60$  were either virus-derived or represented genuine similarity between virus  
8 genes and their cellular orthologs. By contrast, investigation of a subset of 100 hits with bit  
9 scores of b 40-59 showed that ~50% could not be confidently confirmed as having a viral  
10 origin (data not shown).

11

12 Artefactual hits to host DNA can occur since some virus genomes contain genes that have  
13 cellular homologs [104], some virus genomes contain captured host DNA [105]. To  
14 distinguish host from virus-derived DNA in these cases, we exported such hits from the  
15 screening database and virtually translated them to obtain a polypeptide sequence. We then  
16 used the translated sequences as query input to online BLAST searches of GenBank's non-  
17 redundant (nr) database. If searches revealed closer matching to host genes than to known  
18 viral genes, the input sequences were assumed to be host derived. Wherever this occurred  
19 we incorporated representatives of the matching host sequences into the RSL, so that they  
20 would be recognised as host hits on reclassification. By updating hit classifications in this  
21 way, we could progressively filter out host-derived hits from our final screening output.

22

### 23 *Filtering sequences-derived from exogenous viruses*

24 Sequences derived from exogenous viruses are occasionally incorporated into WGS  
25 assemblies. We used SQL queries to identify and exclude these sequences based on hit  
26 characteristics. Where hits derived from virus species or species groups that have been  
27 sequenced previously, they could be discriminated on the basis of sequence identity (i.e.,  
28 98-99% nucleotide-level identity known viruses. The 'extract start' field could be used to

1 identify sequences that lacked flanking genomic sequences, indicating a potential  
2 exogenous origin. We also examined the virtually translated sequences to look for evidence  
3 of long-term presence in the host germline (e.g., stop codons, frameshifting mutations).

4

#### 5 Filtering of cross-matching retrovirus-derived sequences

6 Hits that match more closely to virus genomes than to host DNA, and are clearly inserted  
7 into host DNA, are most likely *bona fide* EVE sequences. However, they may not necessarily  
8 be non-retroviral EVEs, because some filoviruses and arenaviruses (family *Arenaviridae*)  
9 contain glycoprotein genes that are distantly related to those found in certain retroviruses  
10 [106, 107]. When such hits were investigated and found to correspond to ERVs (established  
11 through the presence of proviral genome features adjacent to the hit) we included the  
12 putative sequences of glycoproteins encoded by these ERVs into our RSL and reclassified  
13 hits, so that spurious matches could be recognized as ERV-derived.

14

#### 15 Genomic analysis

16 Previous studies of presence/absence patterns have shown that non-retroviral EVEs are  
17 present in many genomes due to orthology (ancient insertions) rather than paralogy (recent  
18 independent insertion) [48, 50-52]. To differentiate orthologs of previously described EVEs  
19 from newly identified paralogs, we substituted expanded our RSL to include  
20 consensus/reference sequences representing unique EVE loci. This set of EVE loci was  
21 comprised of insertions identified in previous studies [48, 50, 51, 53, 108], as well as a set of  
22 clearly novel EVEs identified in the present screen. For high-copy number, amplified  
23 lineages within this set, we only included a single reference sequence, rather than  
24 attempting to represent each individual ortholog, since it was clear that these elements  
25 derive from a single germline incorporation event (see **Fig. S8**). EVEs were considered  
26 novel if: (i) they derived from a virus group not previously reported in the host group in which  
27 they were identified, or (ii) occurred in species only distantly related to species in which  
28 similar EVEs had been identified previously (e.g. an entirely distinct host class). Whenever

1 novel EVEs were defined, results were reclassified using the updated RSL (see **Fig. 2**).  
2 Orthologs of previously identified EVEs could be inferred by using SQL queries to  
3 summarise screening results, as they disclosed high similarity to these EVE sequences and  
4 occurred in host species relatively closely related to the species in which the putatively  
5 orthologous EVEs had previously been identified. By contrast, novel paralogs either  
6 disclosed only limited similarity to previously identified EVE sequences or occurred in  
7 distantly related host species. This approach to discriminating between paralogs and  
8 orthologs has limitations, but can guide further investigations that use more reliable  
9 approaches (e.g. via investigation of flanking sequences, or phylogeny) to infer orthology  
10 [51]. Se-AL (version 2.0a11) was used to inspect multiple sequence alignments of EVEs and  
11 genomic flanking sequences. Minimum age estimates were obtained for orthologous EVEs  
12 by using host species divergence time estimates collated in TimeTree [109]. We identified  
13 open reading frames and open coding regions within EVEs using PERL scripts available on  
14 request.

15

#### 16 *Phylogenetic analysis*

17 Phylogenies were reconstructed using the maximum likelihood approach implemented in  
18 RAxML (version 8.2.12) [110] and model parameters selected using IQ-TREE model  
19 selection function [111]. Support for phylogenies was assessed via 1000 non-parametric  
20 bootstrap replicates. A time-calibrated vertebrate phylogeny was obtained via TimeTree, an  
21 open database of species divergence time estimates [109]. To determine germline infiltration  
22 rate, we divided the total number of distinct EVE orthologs identified in each vertebrate class  
23 by the total amount of branch length sampled for that class (obtained from the time-  
24 calibrated phylogeny).

25

#### 26 *Application of standardised nomenclature to EVE loci*

27 We assigned all non-retroviral EVEs identified in our study unique identifiers (IDs), following  
28 a convention developed for ERVs [112], Each was assigned a unique identifier (ID)

1 constructed from three components. The first component is a classifier denoting the type of  
2 EVE. The second component comprises: (i) the name of the taxonomic group of viruses the  
3 element derived from and; (ii) a numeric ID that uniquely identifies a specific integration  
4 locus, or for multicopy lineages, a unique founding event. The final component denotes the  
5 taxonomic distribution of the element. This approach has been applied in several previous  
6 studies of vertebrate EVEs [50, 51, 53, 61] and we maintained consistency with these  
7 studies with respect to the numeric ID. Where our study revealed new information about the  
8 taxonomic relationship of an EVE to contemporary viruses, or its distribution across taxa, the  
9 ID was updated accordingly.



1 **FIGURE LEGENDS**

2

3 **Figure 1. Exploring the dark genome using in silico screening.**

4 **(a) Overview of sequence similarity search-based screening.** Screening aims to identify  
5 and classify sequences similar to a set of query sequences within a target database (TDb)  
6 comprising whole genome sequence assemblies of multiple species. The schematic shows  
7 the steps that comprise a single round of screening, as follows: (i) A BLAST search is  
8 performed using a probe sequence selected from a curated ‘reference sequence library’  
9 (RSL) and a ‘target’ file is selected from the TDb; (ii) Matching sequences (referred to as  
10 ‘hits’) identified in this screen are classified via similarity search-based comparison to the  
11 RSL; (iii) A non-redundant set of classified hits is compiled, incorporating hits from previous  
12 rounds of screening.

13 **(b) Comparative analysis of screen output.** Sequences recovered via screening can be  
14 investigated using a wide range of *in silico*, comparative approaches, as follows: (i) analysis  
15 of feature distribution – e.g. annotating host phylogeny to show frequency of occurrence  
16 (coloured circles); (ii) phylogenetic screening, in which sequences obtained via similarity  
17 search-based screening are investigated in phylogenetic reconstructions (e.g. to identify  
18 novel lineages not present in the RSL, as shown here); (iii) Pairwise sequence comparisons  
19 – these can be used to identify specific differences in sequence obtained via screening,  
20 relative to reference sequences; (iv) Comparative phylogenetic analysis - the genetic  
21 properties of novel homologs can be inferred via comparative analysis (e.g. pairwise  
22 comparisons), while their phenotypic properties can potentially be investigated  
23 experimentally (e.g., via transcriptome sequencing).

1 **Figure 2. The database-integrated genome screening (DIGS) process as implemented**  
2 **in the DIGS tool. (i) Screening. (a)** On initiation of screening a list of searches, composed  
3 of each query sequence versus each target database (TDb) file is composed based on the  
4 probe and TDb paths supplied to the DIGS program. Subsequently, screening proceeds  
5 systematically as follows: **(b)** the status table of the project-associated screening database is  
6 queried to determine which searches have yet to be performed. if there are no outstanding  
7 searches then screening ends, otherwise it proceeds to step **(b)** wherein the next  
8 outstanding search of the TDb is performed using the selected probe and the appropriate  
9 BLAST+ program. Results are recorded in the data processing table ('active set'); **(c)**  
10 Results in the processing table are compared to those (if any) obtained previously to derive  
11 a non-redundant set of non-overlapping loci, and an updated set of non-redundant hits is  
12 created, with each hit being represented by a single results table row. To create this non-  
13 redundant set, hits that overlap, or occur within a given range of one another, are merged to  
14 create a single entry; **(d)** Nucleotide sequences associated with results table rows are  
15 extracted from TDb files and stored in the results table; **(e)** extracted sequences are  
16 classified via BLAST-based comparison to the RSL using the appropriate BLAST program;  
17 **(f)** The header-encoded details of the best-matching sequence (species name, gene name)  
18 are recorded in the results table. **(f)** The status table is updated to create a record of the  
19 search having been performed, and the next round of screening is initiated. **(ii)**  
20 Reclassification: Hits in the results table can be reclassified following an update to the  
21 reference sequence library.

22

23 **Figure 3. Counts of ERV RT loci identified by identified via database integrated**  
24 **genome screening of 874 vertebrate species.**

25 Box plots showing the distribution of endogenous retrovirus (ERV) reverse transcriptase  
26 (RT) counts in distinct vertebrate classes. Median and range of values are indicated. Circles  
27 indicate counts for individual species. Counts are shown against a log scale. Figure created  
28 in R using ggplot2 and geom\_boxplot.

1

2 **Figure 4. Distribution of virus families across vertebrate hosts.** Circles indicate the  
3 presence of exogenous viruses. Shaded boxes indicate the presence of confirmed  
4 endogenous viral elements. Abbreviations: DNArt = reverse transcribing DNA viruses;  
5 DNAss = single-stranded DNA viruses; DNAds double-stranded DNA viruses; RNAds =  
6 double-stranded RNA viruses; RNAss-ve = single-stranded negative sense RNA viruses;  
7 RNAss+ve = RNAss-ve = single-stranded positive sense RNA viruses. Information on the  
8 distribution of exogenous viruses was obtained from the NCBI virus genomes resource [39],  
9 supplemented with information obtained from recently published papers [56, 113-118].

10

11 **Figure 5. Evolutionary relationships of filoviruses and filovirus-derived EVEs**

12 Bootstrapped maximum likelihood phylogenies showing the evolutionary relationships  
13 between filoviruses and filovirus EVEs in the nucleoprotein (NP) and viral protein 35 (VP35)  
14 genes. Phylogenies were constructed using maximum likelihood as implemented in RAXML,  
15 and codon-aligned nucleotides for each gene. Numbers adjacent internal nodes indicate  
16 bootstrap support (1000 bootstrap replicates). The scale bar indicates evolutionary distance  
17 in substitutions per site. Virus taxon names are shown in regular font, EVE names are  
18 shown bold. EVE names follow standardised nomenclature (see Methods). Brackets to the  
19 right of each tree indicate virus genera (*italics*) and major lineages (**bold**). Silhouettes  
20 indicate host groups following the key. For Ebola virus, Bundibugyo virus, and Tai Forest  
21 virus, the main reservoir hosts are unknown. The inset box adjacent these taxa shows host  
22 species in which one or more of these viruses has been isolated [73, 119], following the key.

1 **Figure 6. Comparison of germline infiltration rates in five vertebrate classes.**

2 Infiltration rates represent the rate of incorporation and fixation per million years (MY) of  
3 species branch length sampled. Rates are shown for each non-retroviral family represented  
4 by vertebrate EVEs. Colours indicate reverse transcribing DNA (DNArt) viruses, single  
5 stranded DNA (DNAss) viruses, single stranded negative sense RNA (RNAss-ve) viruses  
6 and single stranded positive sense RNA (RNAss-ve) viruses, following the key.

7

8 **Figure 7. Overview of germline incorporation in vertebrates.**

9 A time-calibrated phylogeny of vertebrate species examined in this study, obtained via  
10 TimeTree [109]. Minimum ages of endogenization events are indicated by diamonds on  
11 internal nodes for EVE loci present as orthologs in multiple species. The presence of EVE  
12 sequences in each species genome is indicated by circles at phylogeny tips. Circles and  
13 diamonds nodes are scaled by the number of sequences detected and color-coded by virus  
14 family as indicated in legend. For circles, scaling indicates the total number of EVE  
15 sequences detected within each species genome, including both unique and shared  
16 endogenization events. In panel **(a)** the distribution of 10 families of viruses is shown across  
17 vertebrates separately. In panel **(b)** all virus families are shown on the same tree.

**Table 1. Examples of published studies utilising database-integrated screening.**

Genome feature	Target database	Reference sequence library & probes*	Reference	Year
<b>Non-coding DNA</b>				
ZP3AR (and SFP819)	Rodents	ZP3AR*. ZFP819* & related genes	[95] <sup>b</sup>	2022
SHIN (and IAP elements)	Rodents	SHIN*, IAP subgroups*, <i>Retroviridae</i>	[94] <sup>b</sup>	2023
<b>Genes</b>				
OAS1 gene	Mammals	OAS1* & related genes	[91] <sup>b</sup>	2021
APOBEC3 (A3) genes	Mammals	APOBEC3* & related genes	[92] <sup>a</sup>	2020
Interferon stimulated genes (ISGs)	Vertebrates	ISGs* & related genes	[120]	2017
Interferon lambda (IFNL) genes	Vertebrates	IFNLs* & locus marker genes*	[93] <sup>a,b</sup>	2023
<b>Endogenous viral elements</b>				
Family <i>Flaviviridae</i>	Metazoa	AVP, <i>Flaviviridae</i> * & EFVs	[61] <sup>a</sup>	2022
Family <i>Parvoviridae</i>	Vertebrates	AVP, <i>Parvoviridae</i> * & EPVs	[50] <sup>a</sup>	2022
Family <i>Parvoviridae</i>	Vertebrates	AVP, <i>Parvoviridae</i> * & EPVs	[84] <sup>b</sup>	2023
Genus <i>Protoparvovirus</i>	Mammals	AVP, protoparvoviruses* & EPVs	[121] <sup>a,b</sup>	2019
Family <i>Hepadnaviridae</i>	Metazoa	AVP, <i>Hepadnaviridae</i> * & eHBVs	[51] <sup>a</sup>	2021
Family <i>Circoviridae</i>	Metazoa	AVP, <i>Circoviridae</i> * & ECVs	[48] <sup>a</sup>	2019
<b>Endogenous retroviruses</b>				
Genus <i>Lentivirus</i>	Rodents	Lentiviruses*, other XRVs, & ERVs	[122] <sup>a</sup>	2022
Family <i>Retroviridae</i>	Perissodactyls	<i>Retroviridae</i> *, Retroelements	[67] <sup>a</sup>	2018
HERV-T	Hominids	Class I HERVs*, <i>Retroviridae</i>	[123] <sup>a,b</sup>	2017
MuERV-L	Mice	Class III ERVs*, <i>Retroviridae</i>	[124] <sup>b</sup>	2018

**Footnote:** <sup>a</sup>) DIGS was used as part of 'phylogenetic screening' pipeline (see Fig 1c) <sup>b</sup>) DIGS-based investigations were allied to experimental or functional genomics investigations. \* Indicates subset of the RSL from which probes were derived (note that *Retroviridae* here denotes both endogenous and exogenous retroviruses). **Abbreviations:** zinc finger protein (ZFP); 2'-5'-oligoadenylate synthetase 1 (OAS1); intracisternal A-type particle (IAP); endogenous flaviviral element (EFV); endogenous parvoviral element (EPV); endogenous hepadnavirus (eHBV); endogenous circoviral element (ECV). Human endogenous retrovirus (HERV). Murine endogenous retrovirus (muERV); NCBI all virus proteins set (AVP).

**Table 2. ERV RT loci identified via *in silico* screening**

		Retrovirus clade					
Vertebrate class	# WGS	Clade I		Clade II		Clade III	
		Total #	Average #	Total #	Average #	Total #	Average #
Agnatha	3	32	10.67	1*	0.33	300	100.00
Chondrichthyes	6	2018	336.33	0	0.00	2843	473.83
Actinopterii	173	8514	49.21	64*	0.37	2177	12.58
Actinistia	1	0	0.00	0	0.00	97	97.00
Amphibia	34	17319	509.38	973	28.62	8019	235.85
Reptiles	92	13676	148.65	12120	131.74	20197	219.53
Aves	143	17951	125.53	20797	145.43	42014	293.80
Mammalia	452	13676	30.26	174549	386.17	143364	317.18

**Legend.** WGS=whole genome sequence assemblies screened. \* Hits likely due to contamination.

**Table 3.** Number of non-retroviral EVE sequence identified and estimated number of germline incorporation events in distinct vertebrate classes.

Virus Group	# EVEs identified (Estimated # germline incorporation events)															
	Total		Mammalia		Aves		Reptiles		Amphibia		Actinopterygii		Chondrichthyes		Agnatha	
<b>ssRNA-ve</b>																
<i>Bornaviridae</i>	2566	(383)	2434	(292)	30	(11)	27	(14)	52	(44)	22	(21)	-	-	1	(1)
<i>Chuviridae</i>	182	(164)	24	(24)	-	-	23	(23)	9	(9)	119	(108)	-	-	7	-
<i>Filoviridae</i>	390	(69)	389	(68)	-	-	-	-	1	(1)	-	-	-	-	-	-
<i>Paramyxoviridae</i>	19	(17)	-	-	-	-	-	-	4	(3)	14	(3)	1	(1)	-	-
<b>ssRNA+ve</b>																
<i>Flaviviridae</i>	8	(11)	1	(1)	-	-	-	-	-	-	7	(10)	-	-	-	-
<b>DNArt</b>																
<i>Hepadnaviridae</i>	993	(108)	-	-	897	(89)	93	(17)	2	(1)	-	-	1	(1)	-	-
<b>DNAss</b>																
<i>Circoviridae</i>	1198	(131)	918	(29)	32	(15)	91	(19)	82	(24)	68	(38)	-	-	7	(6)
<i>Parvoviridae</i>	689	(238)	534	(199)	34	(10)	34	(13)	12	(7)	12	(6)	3	(3)	-	-
<b>DNAds</b>																
<i>Herpesviridae</i>	13	(8)	11	(6)	1	(1)	1	(1)	-	-	-	-	-	-	-	-
<i>Alloherpesviridae</i>	28	(8)	-	-	-	-	-	-	15	(3)	13	(5)	-	-	-	-
<b>Total</b>	6087	(1137)	4311	(619)	994	(126)	269	(87)	177	(92)	255	(191)	5	(5)	15	(7)

**Legend:** Germline incorporation here implies both integration into the germline and fixation.

## **DECLARATIONS**

### **Ethics approval and consent to participate**

Not applicable

### **Consent for publication**

Not applicable

### **Availability of data and materials**

All data are openly available in the DIGS-for-EVEs repository hosted on GitHub:

<https://github.com/giffordlabcvr/DIGS-for-EVEs>

<https://twitter.com/DigsTool>

### **Competing interests**

None declared.

### **Funding**

DB-M is supported by the V Foundation for Cancer Research and the Searle Scholars Program. TD, SL, JH, and RJG, were funded by the Medical Research Council of the United Kingdom (MC\_UU\_12014/12).

### **Acknowledgments**

We thank Connor Bamford, Paul Bieniasz, and Jamie Henzy for helpful discussions. Additional thanks to Ade Tukur (Aaron Diamond AIDS Research Centre) and Scott Arkinson (MRC-University of Glasgow Centre for Virus Research) for bioinformatics support.



## REFERENCES

1. Margulies, E.H. and E. Birney, *Approaches to comparative sequence analysis: towards a functional view of vertebrate genomes*. Nat Rev Genet, 2008. **9**(4): p. 303-13.
2. Cheng, J.F., J.R. Priest, and L.A. Pennacchio, *Comparative genomics: a tool to functionally annotate human DNA*. Methods Mol Biol, 2007. **366**: p. 229-51.
3. Nobrega, M.A. and L.A. Pennacchio, *Comparative genomic analysis as a tool for biological discovery*. J Physiol, 2004. **554**(Pt 1): p. 31-9.
4. Guan, D. and M.A. Lazar, *Shining light on dark matter in the genome*. Proc Natl Acad Sci U S A, 2019. **116**(50): p. 24919-24921.
5. Wright, B.W., et al., *The dark proteome: translation from noncanonical open reading frames*. Trends Cell Biol, 2022. **32**(3): p. 243-258.
6. Eisenstein, M., *Drug hunters uncloak the non-coding 'hidden' genome*. Nat Biotechnol, 2021. **39**(10): p. 1169-1171.
7. Katzourakis, A. and R.J. Gifford, *Endogenous viral elements in animal genomes*. PLoS Genet, 2010. **6**(11): p. e1001191.
8. Chiba, S., et al., *Widespread endogenization of genome sequences of non-retroviral RNA viruses into plant genomes*. PLoS Pathog, 2011. **7**(7): p. e1002146.
9. Diop, S.I., et al., *Tracheophyte genomes keep track of the deep evolution of the Caulimoviridae*. Sci Rep, 2018. **8**(1): p. 572.
10. Soucy, S.M., J. Huang, and J.P. Gogarten, *Horizontal gene transfer: building the web of life*. Nat Rev Genet, 2015. **16**(8): p. 472-82.
11. Parrish, N.F. and K. Tomonaga, *Endogenized viral sequences in mammals*. Curr Opin Microbiol, 2016. **31**: p. 176-183.
12. de Tomás, C. and C.M. Vicient, *Genome-wide identification of Reverse Transcriptase domains of recently inserted endogenous plant pararetrovirus (Caulimoviridae)*. Front Plant Sci, 2022. **13**: p. 1011565.
13. Gong, Z., Y. Zhang, and G.Z. Han, *Molecular fossils reveal ancient associations of dsDNA viruses with several phyla of fungi*. Virus Evol, 2020. **6**(1): p. veaa008.
14. Bellas, C., et al., *Large-scale invasion of unicellular eukaryotic genomes by integrating DNA viruses*. Proc Natl Acad Sci U S A, 2023. **120**(16): p. e2300465120.
15. Dewannieux, M. and T. Heidmann, *Endogenous retroviruses: acquisition, amplification and taming of genome invaders*. Curr Opin Virol, 2013. **3**(6): p. 646-56.
16. Geis, F.K. and S.P. Goff, *Silencing and Transcriptional Regulation of Endogenous Retroviruses: An Overview*. Viruses, 2020. **12**(8).
17. Srinivasachar Badarinarayan, S. and D. Sauter, *Switching Sides: How Endogenous Retroviruses Protect Us from Viral Infections*. J Virol, 2021. **95**(12).
18. Fujino, K., et al., *A Human Endogenous Bornavirus-Like Nucleoprotein Encodes a Mitochondrial Protein Associated with Cell Viability*. J Virol, 2021. **95**(14): p. e0203020.
19. Ophinni, Y., et al., *piRNA-Guided CRISPR-like Immunity in Eukaryotes*. Trends Immunol, 2019. **40**(11): p. 998-1010.
20. Patel, M.R., M. Emerman, and H.S. Malik, *Paleovirology - ghosts and gifts of viruses past*. Curr Opin Virol, 2011. **1**(4): p. 304-9.
21. Holmes, E.C., *The evolution of endogenous viral elements*. Cell Host Microbe, 2011. **10**(4): p. 368-77.
22. 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.
23. 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.
24. Camacho, C., et al., *BLAST+: architecture and applications*. BMC Bioinformatics, 2009. **10**: p. 421.

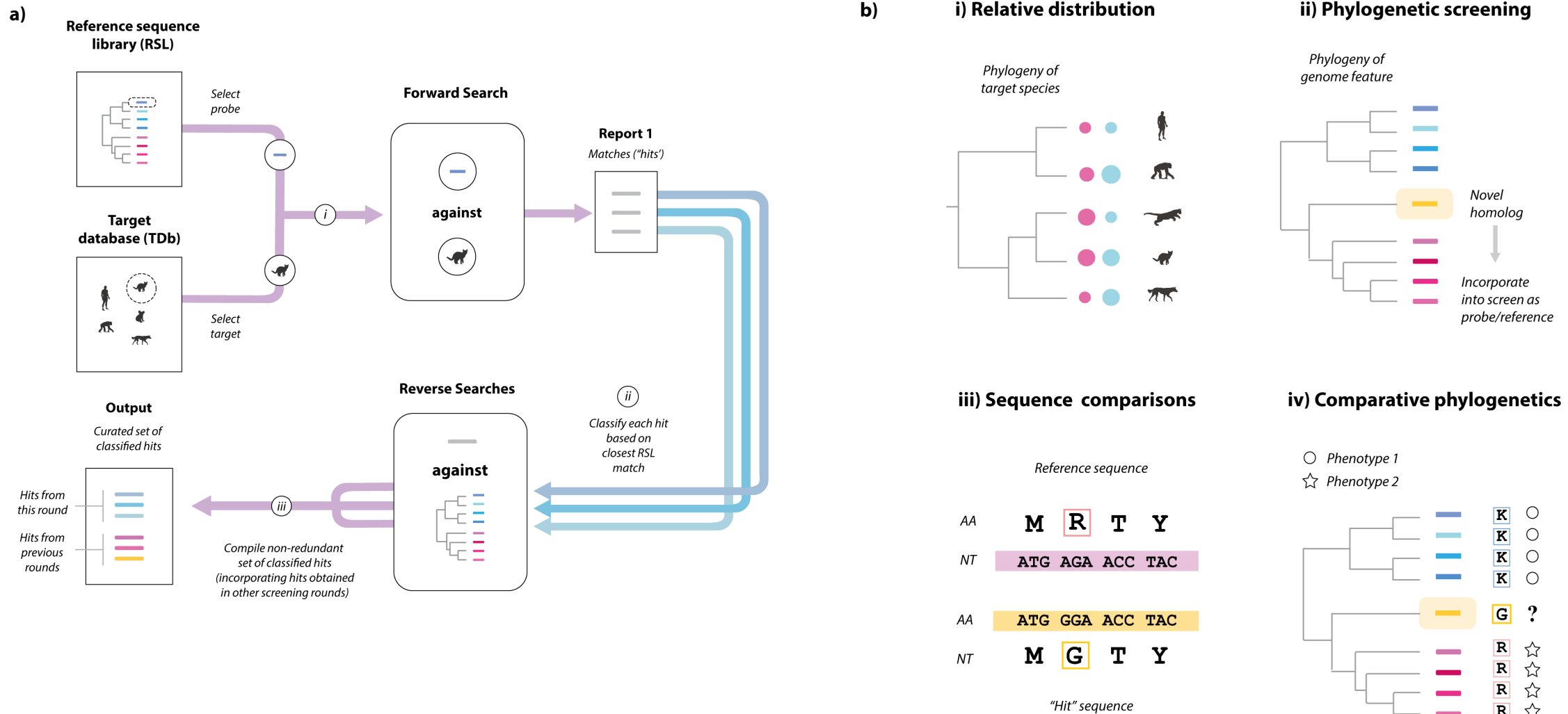
25. Pearson, W.R. and A.J. Mackey, *Using SQL Databases for Sequence Similarity Searching and Analysis*. Curr Protoc Bioinformatics, 2017. **59**: p. 9.4.1-9.4.22.
26. Belyi, V.A., A.J. Levine, and A.M. Skalka, *Sequences from ancestral single-stranded DNA viruses in vertebrate genomes: the parvoviridae and circoviridae are more than 40 to 50 million years old*. J Virol, 2010. **84**(23): p. 12458-62.
27. Heusinger, E., et al., *Early Vertebrate Evolution of the Host Restriction Factor Tetherin*. J Virol, 2015. **89**(23): p. 12154-65.
28. Blanco-Melo, D., S. Venkatesh, and P.D. Bieniasz, *Origins and Evolution of tetherin, an Orphan Antiviral Gene*. Cell Host Microbe, 2016. **20**(2): p. 189-201.
29. Waterhouse, R.M., et al., *OrthoDB: a hierarchical catalog of animal, fungal and bacterial orthologs*. Nucleic Acids Res, 2013. **41**(Database issue): p. D358-65.
30. Cunningham, F., et al., *Ensembl 2015*. Nucleic Acids Res, 2015. **43**(Database issue): p. D662-9.
31. Gifford, R.J. *Database-integrated genome screening (DIGS) tool*. 2022; Available from: <https://giffordlabcvr.github.io/DIGS-tool/>.
32. Johnson, W.E., *Origins and evolutionary consequences of ancient endogenous retroviruses*. Nat Rev Microbiol, 2019. **17**(6): p. 355-370.
33. Hayward, A., M. Grabherr, and P. Jern, *Broad-scale phylogenomics provides insights into retrovirus-host evolution*. Proc Natl Acad Sci U S A, 2013. **110**(50): p. 20146-51.
34. Xiong, Y. and T.H. Eickbush, *Origin and evolution of retroelements based upon their reverse transcriptase sequences*. Embo j, 1990. **9**(10): p. 3353-62.
35. Tristem, M., *Identification and characterization of novel human endogenous retrovirus families by phylogenetic screening of the human genome mapping project database*. J Virol, 2000. **74**(8): p. 3715-30.
36. Hayward, A., C.K. Cornwallis, and P. Jern, *Pan-vertebrate comparative genomics unmasks retrovirus macroevolution*. Proc Natl Acad Sci U S A, 2015. **112**(2): p. 464-9.
37. Han, G.Z., *Extensive retroviral diversity in shark*. Retrovirology, 2015. **12**: p. 34.
38. Xu, X., et al., *Endogenous retroviruses of non-avian/mammalian vertebrates illuminate diversity and deep history of retroviruses*. PLoS Pathog, 2018. **14**(6): p. e1007072.
39. Brister, J.R., et al., *NCBI viral genomes resource*. Nucleic Acids Res, 2015. **43**(Database issue): p. D571-7.
40. Koonin, E.V., M. Krupovic, and N. Yutin, *Evolution of double-stranded DNA viruses of eukaryotes: from bacteriophages to transposons to giant viruses*. Ann N Y Acad Sci, 2015. **1341**(1): p. 10-24.
41. Inoue, Y. and H. Takeda, *Teratorn and its relatives - a cross-point of distinct mobile elements, transposons and viruses*. Front Vet Sci, 2023. **10**: p. 1158023.
42. Gifford, R.J. *DIGS-for-EVEs*. 2023; Available from: <https://github.com/giffordlabcvr/DIGS-for-EVEs>.
43. Horie, M., et al., *An RNA-dependent RNA polymerase gene in bat genomes derived from an ancient negative-strand RNA virus*. Scientific Reports, 2016. **6**(1): p. 25873.
44. Ho, A.L.F.C., C.L. Pruett, and J. Lin, *Phylogeny and biogeography of Poecilia (Cyprinodontiformes: Poeciliinae) across Central and South America based on mitochondrial and nuclear DNA markers*. Molecular Phylogenetics and Evolution, 2016. **101**: p. 32-45.
45. Aswad, A. and A. Katzourakis, *The first endogenous herpesvirus, identified in the tarsier genome, and novel sequences from primate rhadinoviruses and lymphocryptoviruses*. PLoS Genet, 2014. **10**(6): p. e1004332.
46. Aswad, A., et al., *Evolutionary History of Endogenous Human Herpesvirus 6 Reflects Human Migration out of Africa*. Mol Biol Evol, 2021. **38**(1): p. 96-107.
47. Liu, X., et al., *Endogenization and excision of human herpesvirus 6 in human genomes*. PLoS Genet, 2020. **16**(8): p. e1008915.
48. Dennis, T.P.W., et al., *The evolution, distribution and diversity of endogenous circoviral elements in vertebrate genomes*. Virus Res, 2019. **262**: p. 15-23.

49. Dennis, T.P.W., et al., *Insights into Circovirus Host Range from the Genomic Fossil Record*. J Virol, 2018. **92**(16).
50. Campbell, M.A., et al., *Comparative analysis reveals the long-term co-evolutionary history of parvoviruses and vertebrates*. bioRxiv, 2022: p. 2021.10.25.465781.
51. 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.
52. Suh, A., et al., *Early mesozoic coexistence of amniotes and hepadnaviridae*. PLoS Genet, 2014. **10**(12): p. e1004559.
53. Kawasaki, J., et al., *100-My history of bornavirus infections hidden in vertebrate genomes*. Proc Natl Acad Sci U S A, 2021. **118**(20).
54. Horie, M., et al., *Endogenous non-retroviral RNA virus elements in mammalian genomes*. Nature, 2010. **463**(7277): p. 84-7.
55. Hyndman, T.H., et al., *Isolation and molecular identification of Sunshine virus, a novel paramyxovirus found in Australian snakes*. Infect Genet Evol, 2012. **12**(7): p. 1436-46.
56. Harvey, E., et al., *Divergent hepaciviruses, delta-like viruses and a chu-like virus in Australian marsupial carnivores (dasyurids)*. Virus Evolution, 2023.
57. Gorbalenya, A.E. and C. Lauber, *Phylogeny of Viruses*. Reference Module in Biomedical Sciences, 2017.
58. Edwards, M.R., et al., *Conservation of Structure and Immune Antagonist Functions of Filoviral VP35 Homologs Present in Microbat Genomes*. Cell Rep, 2018. **24**(4): p. 861-872.e6.
59. Kondoh, T., et al., *Putative endogenous filovirus VP35-like protein potentially functions as an IFN antagonist but not a polymerase cofactor*. PLoS One, 2017. **12**(10): p. e0186450.
60. Li, Y., et al., *Endogenous Viral Elements in Shrew Genomes Provide Insights into Pestivirus Ancient History*. Mol Biol Evol, 2022. **39**(10).
61. Bamford, C.G.G., et al., *Comparative analysis of genome-encoded viral sequences reveals the evolutionary history of flavivirids (family Flaviviridae)*. Virus Evolution, 2022. **8**(2).
62. Qin, X.C., et al., *A tick-borne segmented RNA virus contains genome segments derived from unsegmented viral ancestors*. Proc Natl Acad Sci U S A, 2014. **111**(18): p. 6744-9.
63. Birney, E., et al., *Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project*. Nature, 2007. **447**(7146): p. 799-816.
64. Schattner, P., *Automated querying of genome databases*. PLoS Comput Biol, 2007. **3**(1): p. e1.
65. Obbard, D.J., *Expansion of the metazoan virosphere: progress, pitfalls, and prospects*. Curr Opin Virol, 2018. **31**: p. 17-23.
66. Zhang, Y.Z., M. Shi, and E.C. Holmes, *Using Metagenomics to Characterize an Expanding Virosphere*. Cell, 2018. **172**(6): p. 1168-1172.
67. Zhu, H., R.J. Gifford, and P.R. Murcia, *Distribution, Diversity, and Evolution of Endogenous Retroviruses in Perissodactyl Genomes*. J Virol, 2018. **92**(23).
68. Reus, K., et al., *HERV-K(OLD): ancestor sequences of the human endogenous retrovirus family HERV-K(HML-2)*. J Virol, 2001. **75**(19): p. 8917-26.
69. Pavlíček, A., et al., *Processed pseudogenes of human endogenous retroviruses generated by LINEs: their integration, stability, and distribution*. Genome Res, 2002. **12**(3): p. 391-9.
70. Belshaw, R., et al., *High copy number in human endogenous retrovirus families is associated with copying mechanisms in addition to reinfection*. Mol Biol Evol, 2005. **22**(4): p. 814-7.
71. Starrett, G.J., et al., *Adintoviruses: a proposed animal-tropic family of midsize eukaryotic linear dsDNA (MELD) viruses*. Virus Evolution, 2021. **7**(1).
72. Mahanty, S. and M. Bray, *Pathogenesis of filoviral haemorrhagic fevers*. Lancet Infect Dis, 2004. **4**(8): p. 487-98.

73. Mari Saez, A., et al., *Investigating the zoonotic origin of the West African Ebola epidemic*. EMBO Mol Med, 2015. **7**(1): p. 17-23.
74. Taylor, D.J., et al., *Evidence that ebolaviruses and cuevaviruses have been diverging from marburgviruses since the Miocene*. PeerJ, 2014. **2**: p. e556.
75. Carroll, S.A., et al., *Molecular evolution of viruses of the family Filoviridae based on 97 whole-genome sequences*. J Virol, 2013. **87**(5): p. 2608-16.
76. Kryukov, K., et al., *Systematic survey of non-retroviral virus-like elements in eukaryotic genomes*. Virus Res, 2019. **262**: p. 30-36.
77. Cui, J., et al., *Low frequency of paleoviral infiltration across the avian phylogeny*. Genome Biol, 2014. **15**(12): p. 539.
78. Frank, J.A. and C. Feschotte, *Co-option of endogenous viral sequences for host cell function*. Curr Opin Virol, 2017. **25**: p. 81-89.
79. Aswad, A. and A. Katzourakis, *Paleovirology and virally derived immunity*. Trends Ecol Evol, 2012. **27**(11): p. 627-36.
80. Bravo, A., et al., *Antiviral Activity of an Endogenous Parvoviral Element*. Viruses, 2023. **15**(7).
81. Lavielle, C., et al., *Paleovirology of 'syncytins', retroviral env genes exapted for a role in placentation*. Philos Trans R Soc Lond B Biol Sci, 2013. **368**(1626): p. 20120507.
82. 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).
83. Pastuzyn, E.D., et al., *The Neuronal Gene Arc Encodes a Repurposed Retrotransposon Gag Protein that Mediates Intercellular RNA Transfer*. Cell, 2018. **172**(1-2): p. 275-288.e18.
84. Quezada-Ramírez, M.A., et al., *Identification of genome safe harbor loci for human gene therapy based on evolutionary biology and comparative genomics*. bioRxiv, 2023: p. 2023.09.08.556857.
85. Hu, G. and L. Kurgan, *Sequence Similarity Searching*. Curr Protoc Protein Sci, 2019. **95**(1): p. e71.
86. Pearson, W.R., *An introduction to sequence similarity ("homology") searching*. Curr Protoc Bioinformatics, 2013. **Chapter 3**: p. Unit3.1.
87. Miller, K., et al., *Identification of multiple Gypsy LTR-retrotransposon lineages in vertebrate genomes*. J Mol Evol, 1999. **49**(3): p. 358-66.
88. Wang, J. and G.Z. Han, *A Sister Lineage of Sampled Retroviruses Corroborates the Complex Evolution of Retroviruses*. Mol Biol Evol, 2021. **38**(3): p. 1031-1039.
89. Kojima, S., et al., *Virus-like insertions with sequence signatures similar to those of endogenous nonretroviral RNA viruses in the human genome*. Proc Natl Acad Sci U S A, 2021. **118**(5).
90. Bruno, M., M. Mahgoub, and T.S. Macfarlan, *The Arms Race Between KRAB-Zinc Finger Proteins and Endogenous Retroelements and Its Impact on Mammals*. Annu Rev Genet, 2019. **53**: p. 393-416.
91. Wickenhagen, A., et al., *A prenylated dsRNA sensor protects against severe COVID-19*. Science, 2021. **374**(6567): p. eabj3624.
92. Ito, J., R.J. Gifford, and K. Sato, *Retroviruses drive the rapid evolution of mammalian APOBEC3 genes*. Proc Natl Acad Sci U S A, 2020. **117**(1): p. 610-618.
93. Bamford, C.G.G., et al., *Partial Gene Conversion Shapes the Emergence of Functional Novelty in the Placental Mammal Interferon Lambda System.*, in *Infectious Diseases Through an Evolutionary Lens*. 2023: British Medical Association House, Tavistock Square, London.
94. Enriquez-Gasca, R., et al., *Co-option of endogenous retroviruses through genetic escape from TRIM28 repression*. Cell Rep, 2023. **42**(6): p. 112625.
95. Fernandes, L.P., et al., *A satellite DNA array barcodes chromosome 7 and regulates totipotency via ZFP819*. Science Advances, 2022. **8**(43): p. eabp8085.
96. Buchfink, B., K. Reuter, and H.G. Drost, *Sensitive protein alignments at tree-of-life scale using DIAMOND*. Nat Methods, 2021. **18**(4): p. 366-368.

97. Camacho, C., et al., *ElasticBLAST: accelerating sequence search via cloud computing*. BMC Bioinformatics, 2023. **24**(1): p. 117.
98. Nawrocki, E.P., D.L. Kolbe, and S.R. Eddy, *Infernal 1.0: inference of RNA alignments*. Bioinformatics, 2009. **25**(10): p. 1335-7.
99. Grabowski, P. and J. Rappsilber, *A Primer on Data Analytics in Functional Genomics: How to Move from Data to Insight?* Trends Biochem Sci, 2019. **44**(1): p. 21-32.
100. Kitts, P.A., et al., *Assembly: a resource for assembled genomes at NCBI*. Nucleic Acids Res, 2016. **44**(D1): p. D73-80.
101. Schoch, C.L., et al., *NCBI Taxonomy: a comprehensive update on curation, resources and tools*. Database (Oxford), 2020. **2020**.
102. Inoue, Y., et al., *Fusion of piggyBac-like transposons and herpesviruses occurs frequently in teleosts*. Zoological Lett, 2018. **4**: p. 6.
103. Barreat, J.G.N. and A. Katzourakis, *Phylogenomics of the Maverick Virus-Like Mobile Genetic Elements of Vertebrates*. Mol Biol Evol, 2021. **38**(5): p. 1731-1743.
104. Koonin, E.V., *On the origin of cells and viruses: primordial virus world scenario*. Ann N Y Acad Sci, 2009. **1178**(1): p. 47-64.
105. Becher, P. and N. Tautz, *RNA recombination in pestiviruses: cellular RNA sequences in viral genomes highlight the role of host factors for viral persistence and lethal disease*. RNA Biol, 2011. **8**(2): p. 216-24.
106. Benit, L., P. Dessen, and T. Heidmann, *Identification, phylogeny, and evolution of retroviral elements based on their envelope genes*. Journal of Virology, 2001. **75**(23): p. 11709-11719.
107. Gallaher, W.R., C. DiSimone, and M.J. Buchmeier, *The viral transmembrane superfamily: possible divergence of Arenavirus and Filovirus glycoproteins from a common RNA virus ancestor*. BMC Microbiol, 2001. **1**: p. 1.
108. Hildebrandt, E., et al., *Evolution of dependoparvoviruses across geological timescales – implications for design of AAV-based gene therapy vectors*. Virus Evolution, 2020.
109. Kumar, S., et al., *TimeTree 5: An Expanded Resource for Species Divergence Times*. Mol Biol Evol, 2022. **39**(8).
110. Stamatakis, A., *RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with thousands of taxa and mixed models*. Bioinformatics, 2006. **22**(21): p. 2688-90.
111. Minh, B.Q., et al., *IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic Era*. Molecular Biology and Evolution, 2020. **37**(5): p. 1530-1534.
112. Gifford, R.J., et al., *Nomenclature for endogenous retrovirus (ERV) loci*. Retrovirology, 2018. **15**(1): p. 59.
113. Ariel, E., *Viruses in reptiles*. Veterinary Research, 2011. **42**(1): p. 100.
114. Waller, S.J., et al., *Cloacal virome of an ancient host lineage - The tuatara (*Sphenodon punctatus*) - Reveals abundant and diverse diet-related viruses*. Virology, 2022. **575**: p. 43-53.
115. Soto, E., et al., *First Isolation of a Novel Aquatic Flavivirus from Chinook Salmon (*Oncorhynchus tshawytscha*) and Its In Vivo Replication in a Piscine Animal Model*. J Virol, 2020. **94**(15).
116. Koda, S.A., et al., *Complete genome sequences of infectious spleen and kidney necrosis virus isolated from farmed albino rainbow sharks *Epalzeorhynchus frenatum* in the United States*. Virus Genes, 2021. **57**(5): p. 448-452.
117. Harding, E.F., et al., *Revealing the uncharacterised diversity of amphibian and reptile viruses*. ISME Communications, 2022. **2**(1): p. 95.
118. Harvey, E. and E.C. Holmes, *Diversity and evolution of the animal virome*. Nat Rev Microbiol, 2022. **20**(6): p. 321-334.
119. Leroy, E.M., et al., *Multiple Ebola virus transmission events and rapid decline of central African wildlife*. Science, 2004. **303**(5656): p. 387-90.

120. Shaw, A.E., et al., *Fundamental properties of the mammalian innate immune system revealed by multispecies comparison of type I interferon responses*. PLoS Biol, 2017. **15**(12): p. e2004086.
121. 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).
122. Kambol, R., A. Gatseva, and R.J. Gifford, *An endogenous lentivirus in the germline of a rodent*. Retrovirology, 2022. **19**(1): p. 30.
123. Blanco-Melo, D., R.J. Gifford, and P.D. Bieniasz, *Co-option of an endogenous retrovirus envelope for host defense in hominid ancestors*. Elife, 2017. **6**.
124. Blanco-Melo, D., R.J. Gifford, and P.D. Bieniasz, *Reconstruction of a replication-competent ancestral murine endogenous retrovirus-L*. Retrovirology, 2018. **15**(1): p. 34.



**Figure 1.**

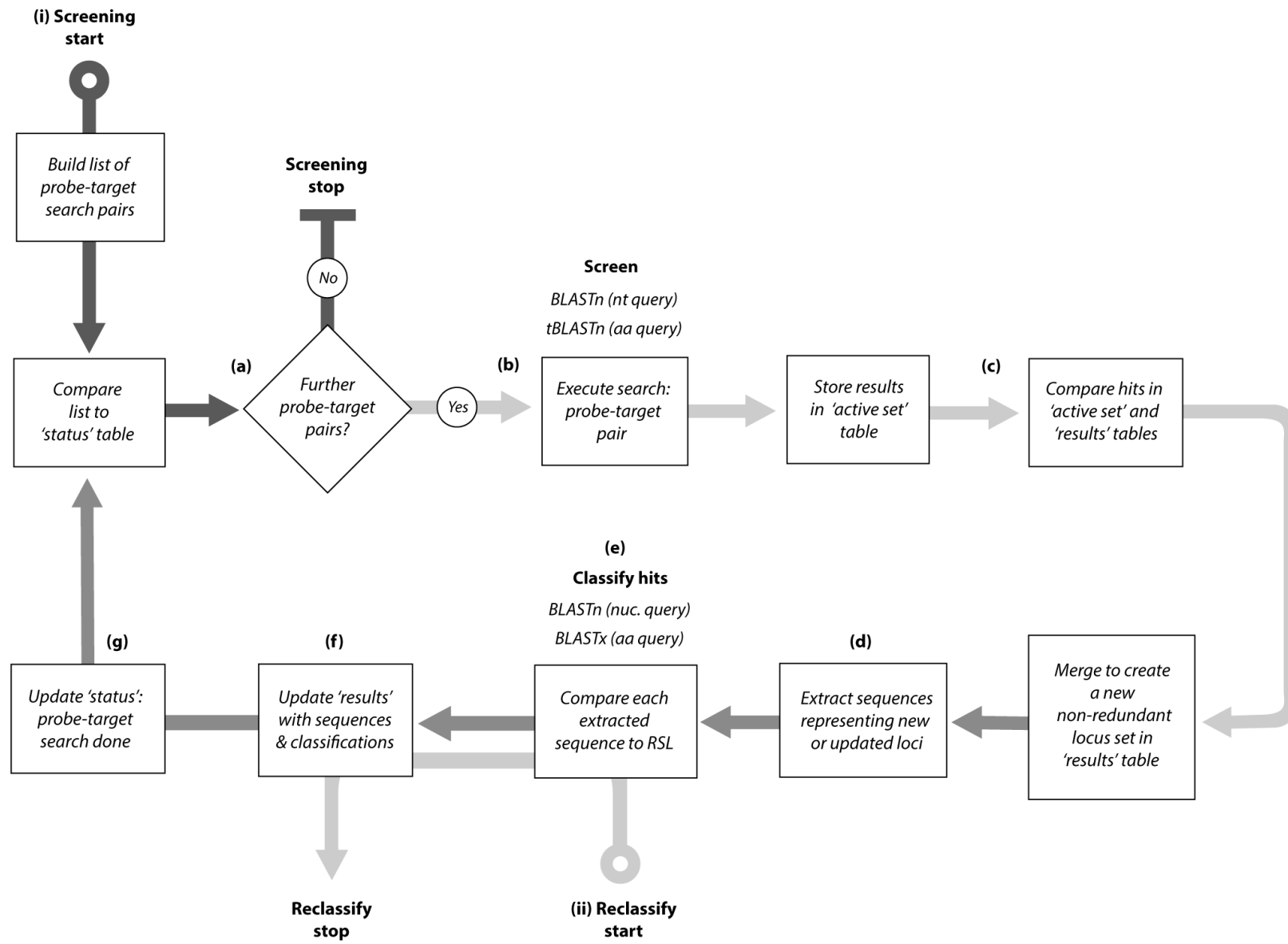
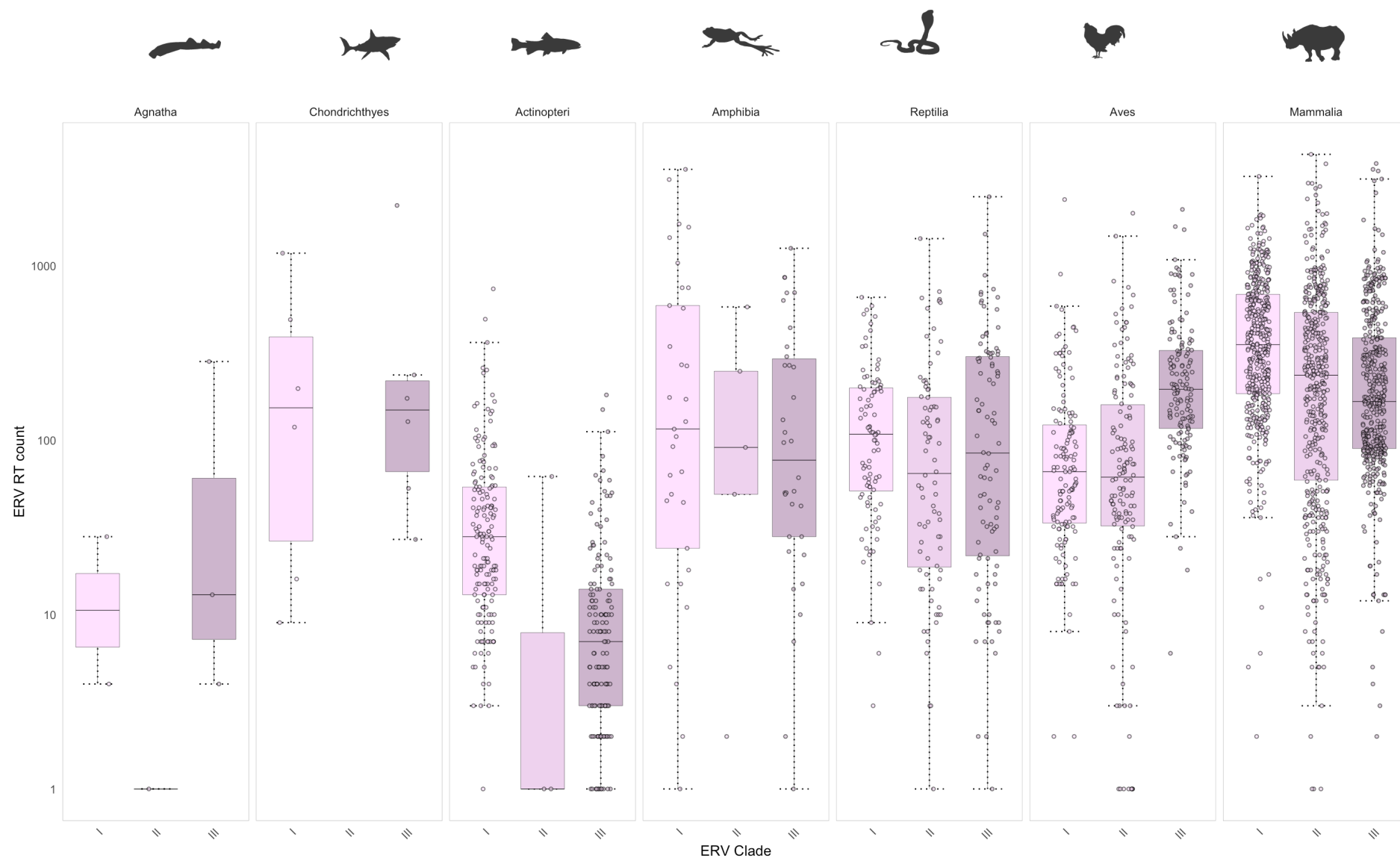


Figure 2.





**Figure 3.**

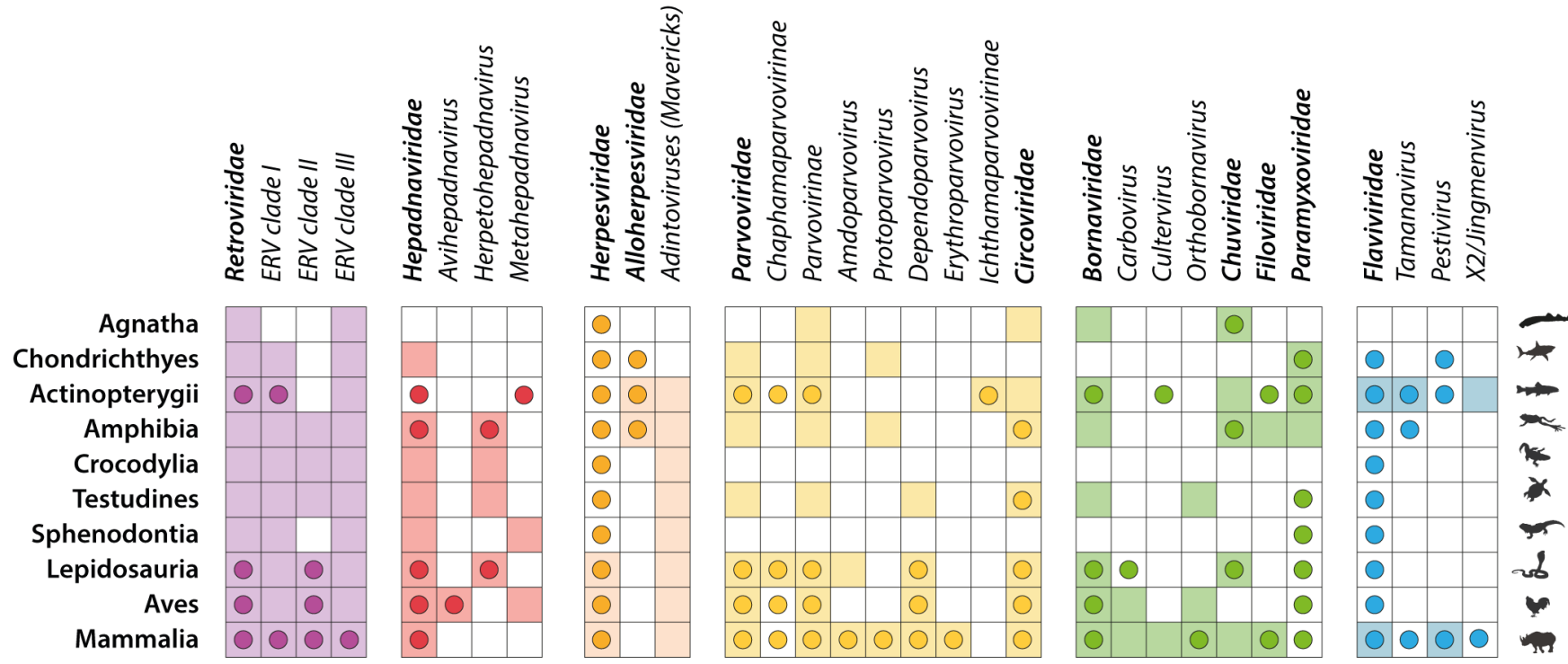
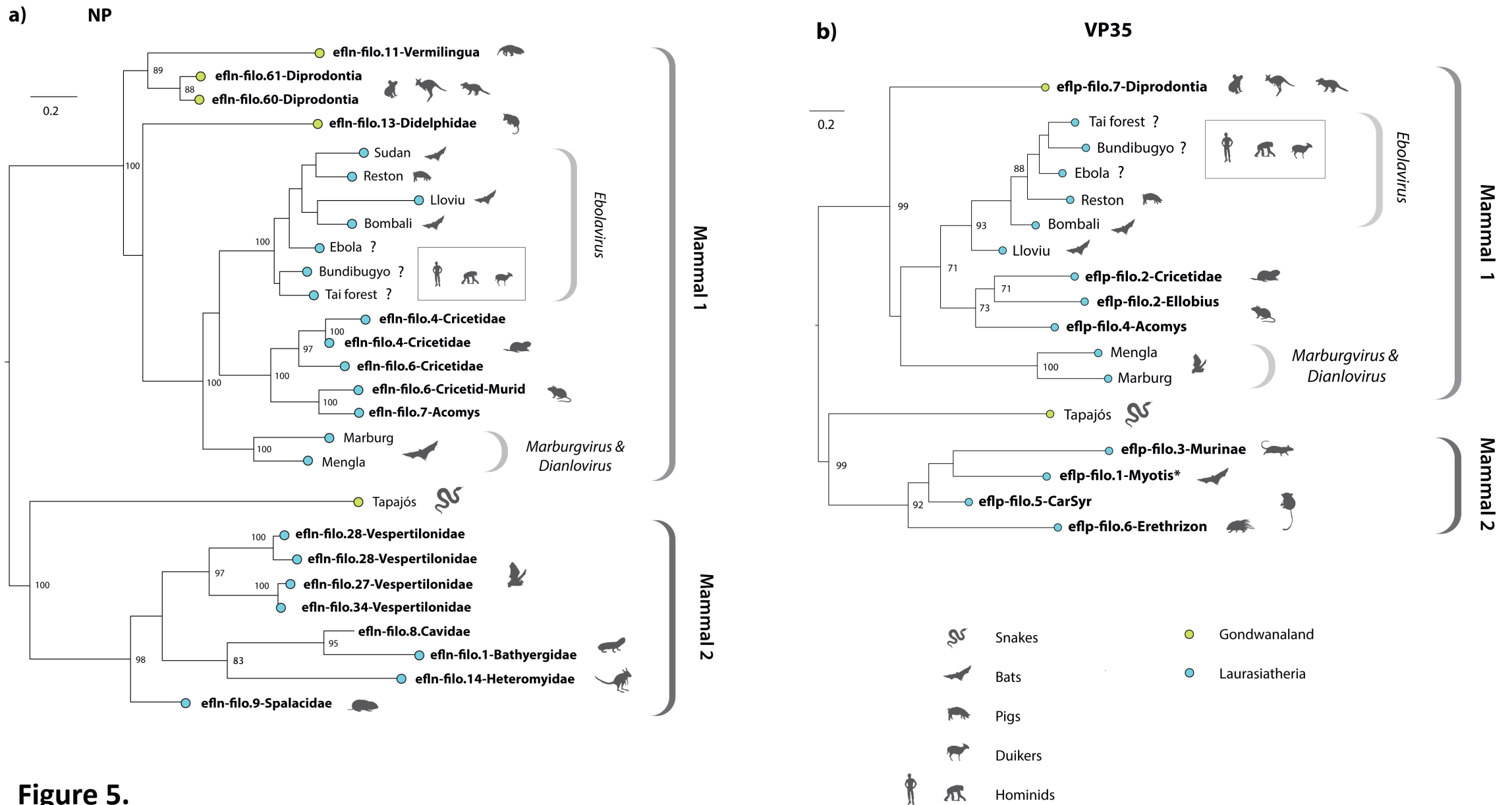
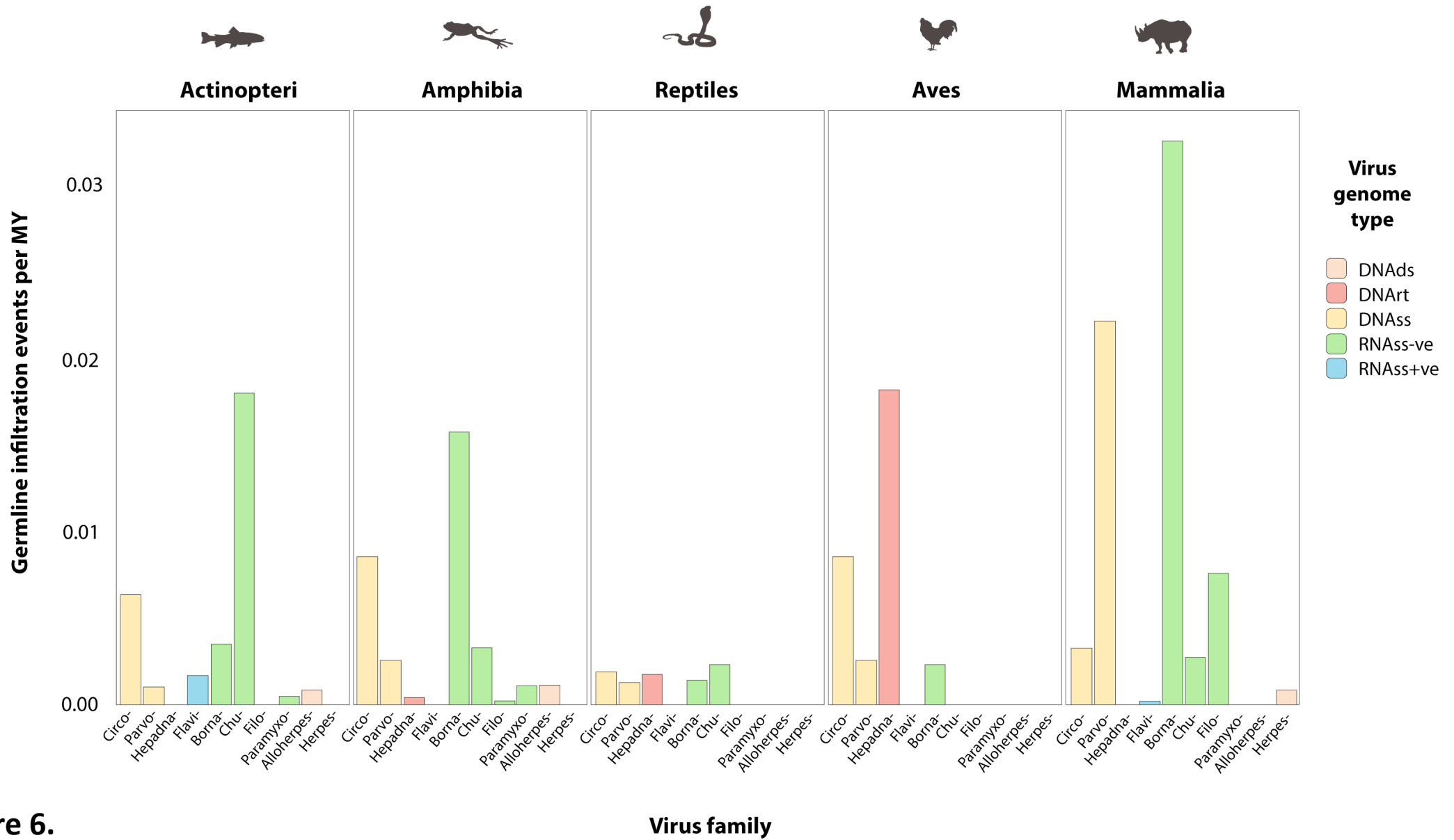


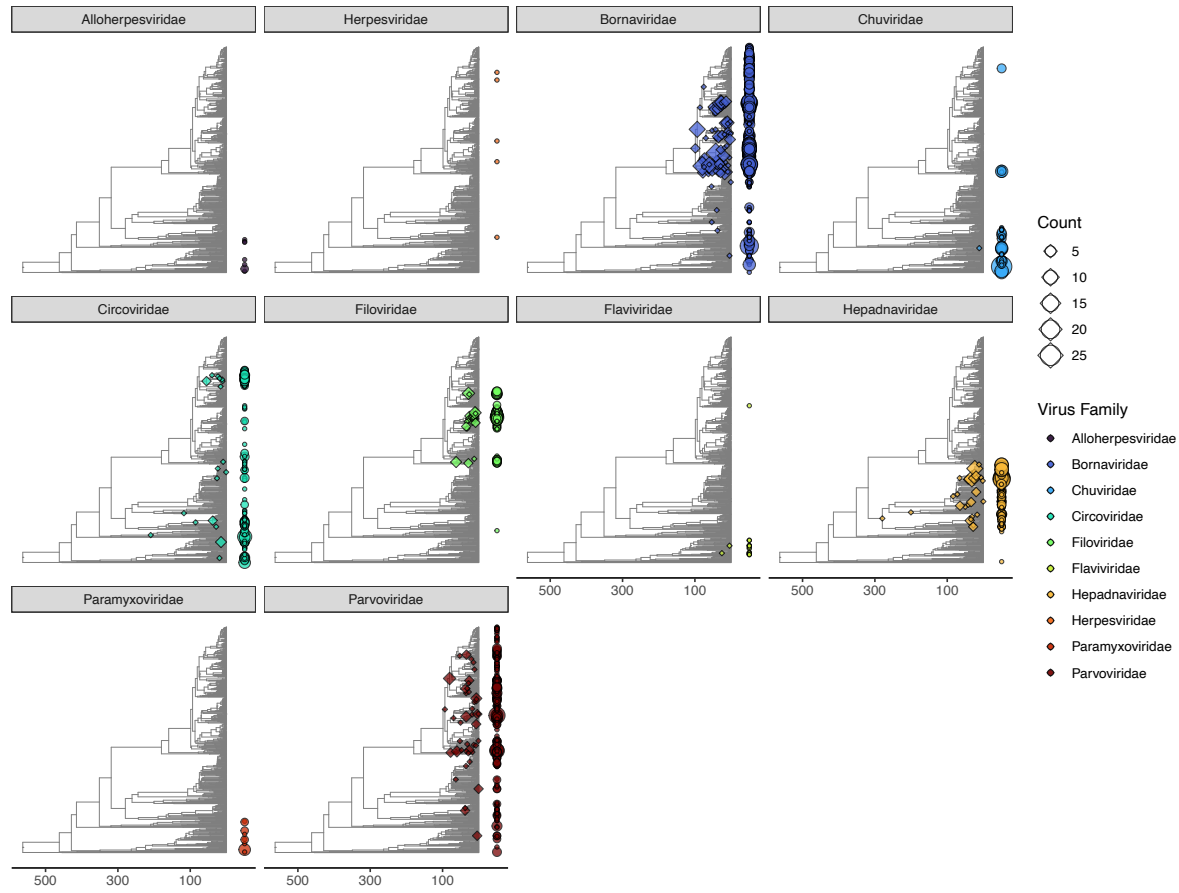
Figure 4.



**Figure 5.**



a)



b)

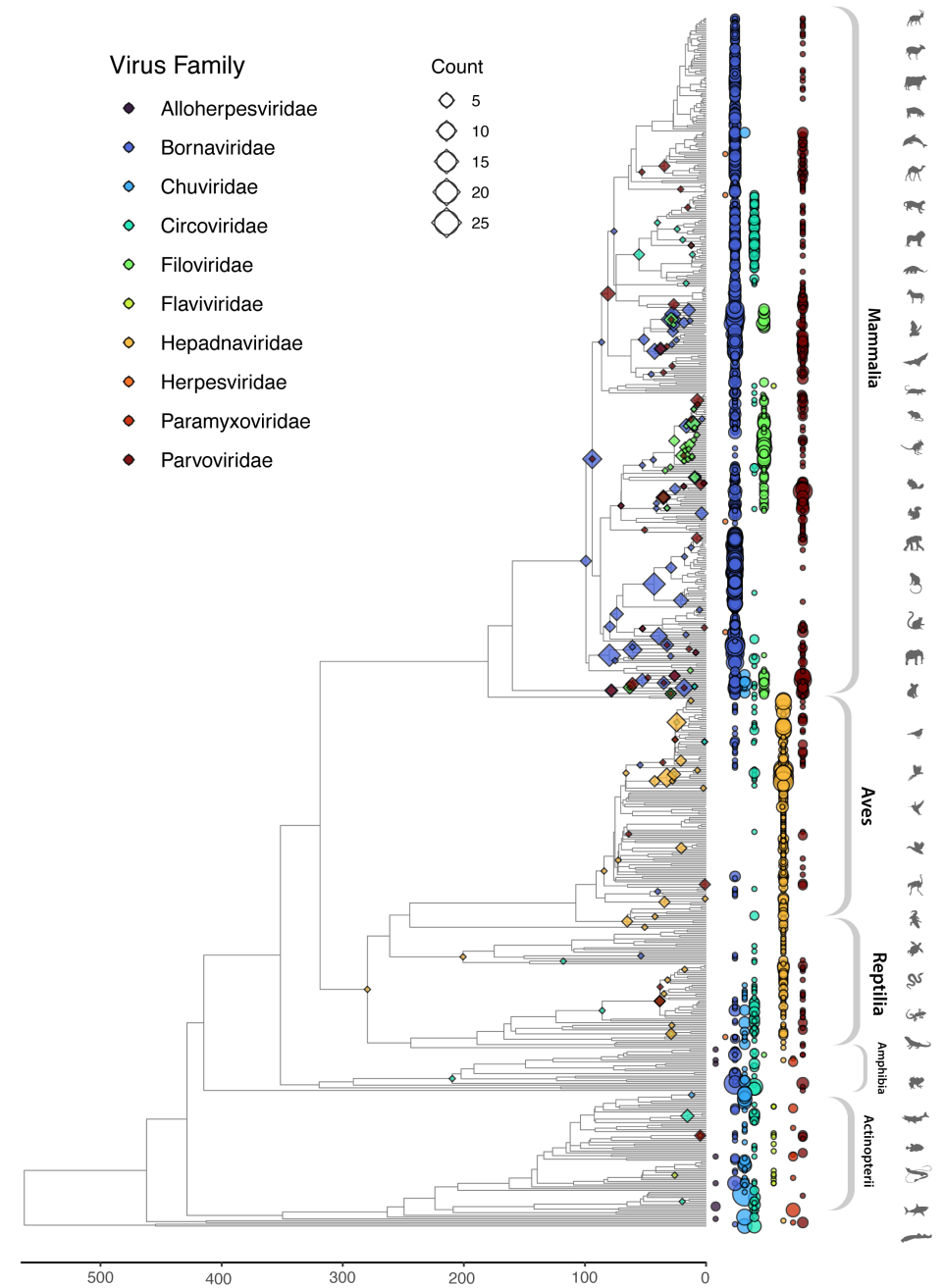


Figure 7.