

## **Supplementary Information for**

One hundred million years history of bornavirus infections hidden in vertebrate genomes

Junna Kawasaki<sup>a,b</sup>, Shohei Kojima<sup>a,1</sup>, Yahiro Mukai<sup>a,b</sup>, Keizo Tomonaga<sup>a,b,c\*</sup>, Masayuki Horie<sup>a,d\*</sup>

### **\*Corresponding authors**

Masayuki Horie

Email: [horie.masayuki.3m@kyoto-u.ac.jp](mailto:horie.masayuki.3m@kyoto-u.ac.jp)

Keizo Tomonaga

Email: [tomonaga.keizo.5r@kyoto-u.ac.jp](mailto:tomonaga.keizo.5r@kyoto-u.ac.jp)

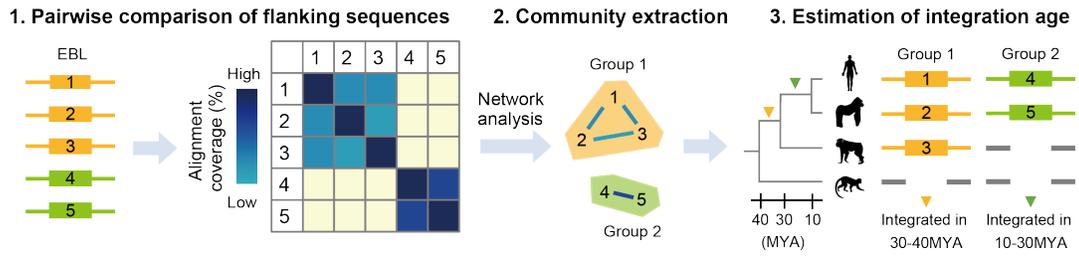
### **This PDF file includes:**

Figures S1 to S4

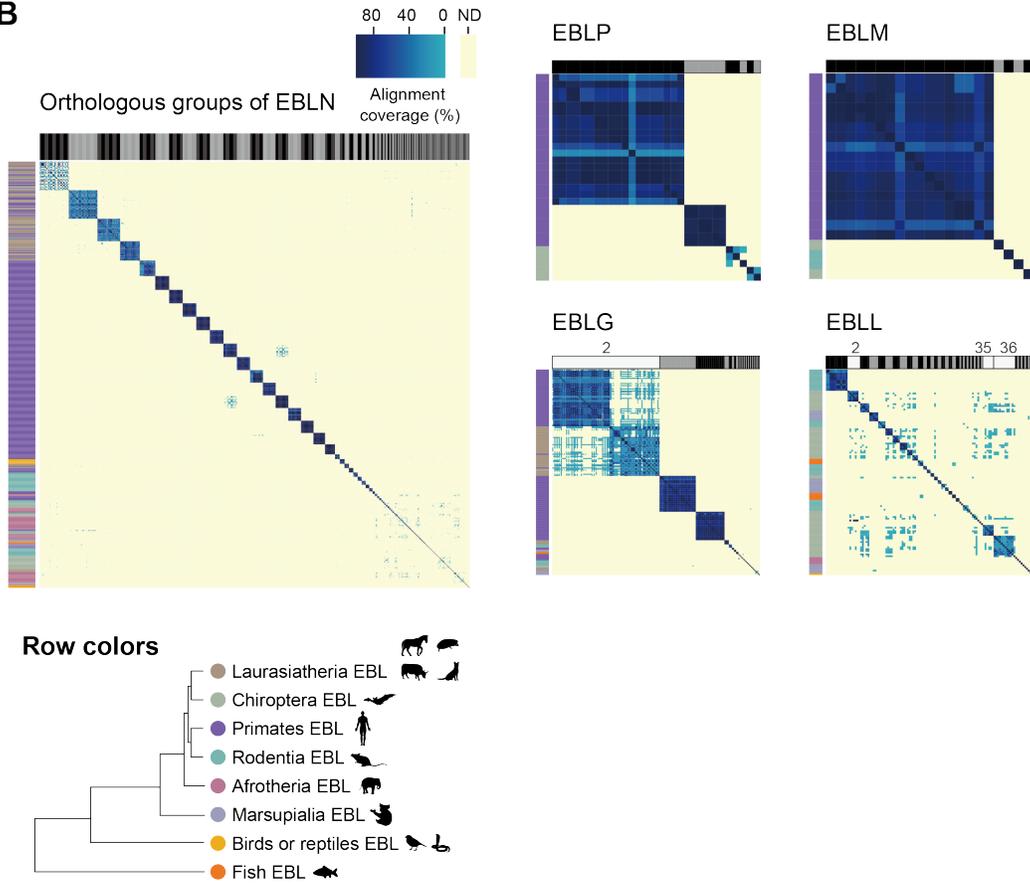
### **Other supplementary materials for this manuscript include the following:**

Datasets S1 to S7

**A**

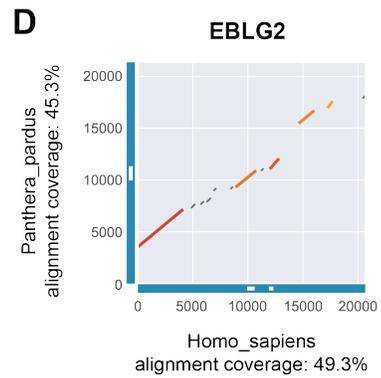
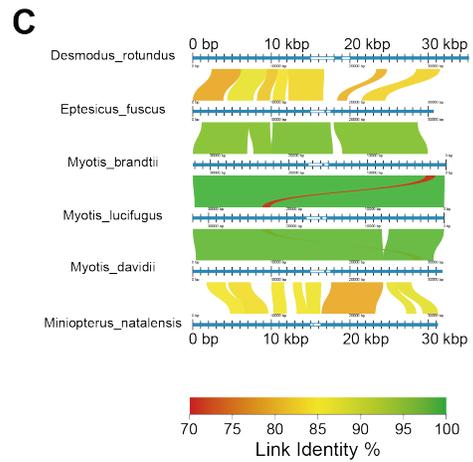
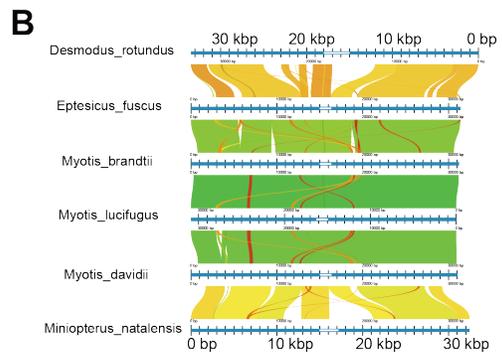
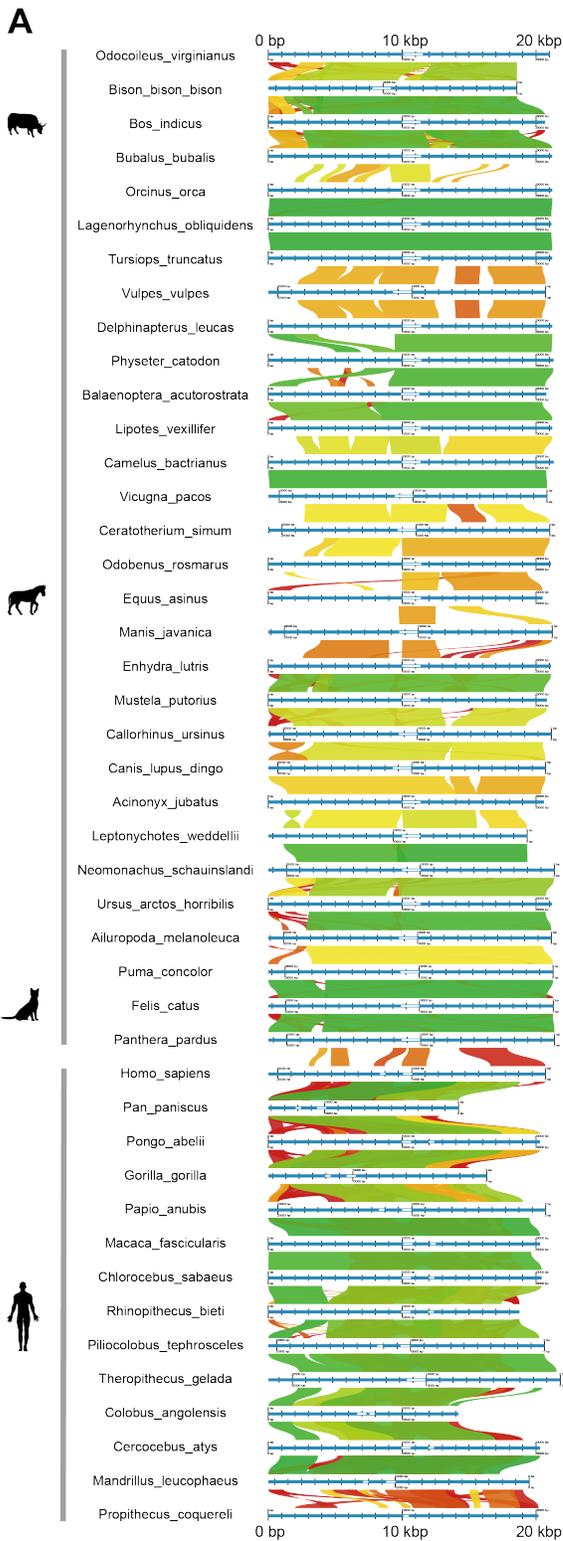


**B**



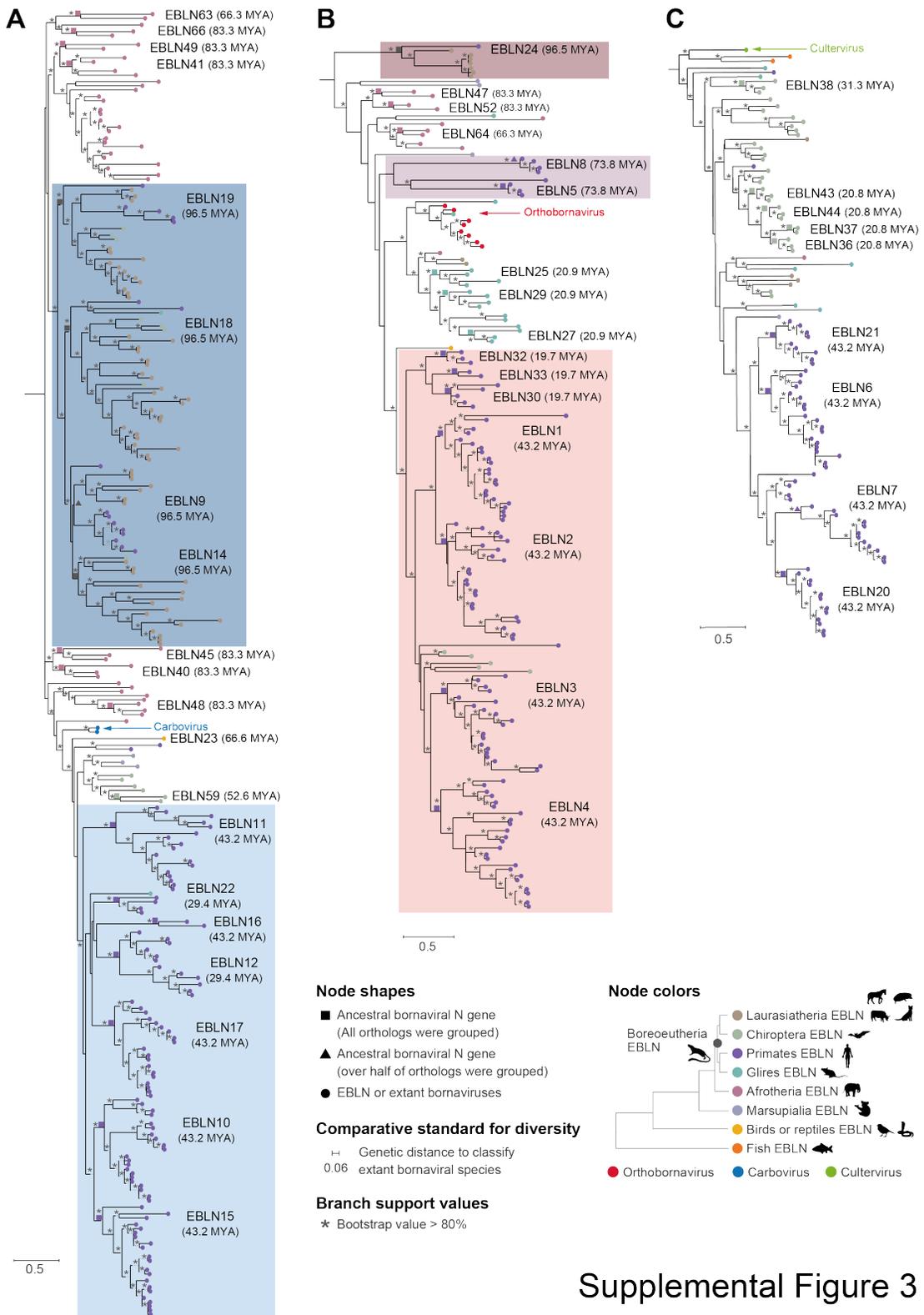
Supplementary Figure 1

1 **Fig. S1. Dating analysis for bornaviral integration events.** (A) Procedure to determine  
2 presence and absence patterns of orthologous EBLs. First, we performed pairwise  
3 sequence comparison among EBL integration sites using BLASTN and made an all-  
4 against-all matrix of their alignment coverages. Second, we constructed a network using  
5 the matrix and grouped EBL loci by extracting community structures from the sequence  
6 network. Finally, the ages of bornavirus integration events were assigned from the  
7 divergence times of the host species with orthologous EBLs. (B) All-against-all matrix of  
8 alignment coverages among EBL integration sites. In the heatmap, the blue color palette  
9 shows the alignment coverage between EBL integration sites (%) and yellow indicates  
10 that sequence similarity was not detected (ND). The column colors indicate EBL groups;  
11 in particular, the white shows manually modified groups (EBLG2, EBLL2, EBLL35, and  
12 EBLL36) (**details in Materials and Methods**). The row colors show host lineages of each  
13 EBL locus.  
14



Supplemental Figure 2

17 **Fig. S2. Alignment quality between EBL integration sites.** (A-C) Schematic images of  
18 alignments between EBL integration sites. The sequence alignments of EBLG2 (A),  
19 EBLL2 (B), and EBLL35 (C) were visualized using AliTV. Blue lines indicate host  
20 chromosomal DNA, and the location of EBLs are shown as white colored portions of the  
21 lines. The black vertical lines are shown for every 1,000 bp. The color palette from red to  
22 green indicates identity scores obtained from lastz. The representative host species are  
23 shown as silhouettes to the left of the alignments. (D) Dot plot between laurasiatherian  
24 and primate EBLG2 integration sites. Line colors except for gray correspond to (A), and  
25 gray lines indicate short fragments aligned by lastz. White portions within the thick blue  
26 lines indicate the positions of EBLG2 in the genomes.  
27



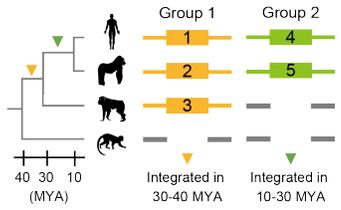
28  
29

Supplemental Figure 3

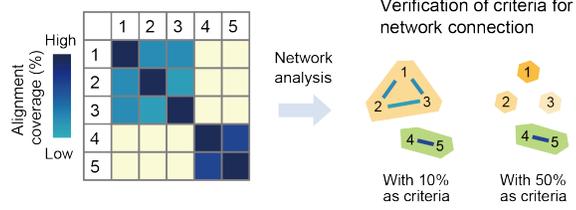
30 **Fig. S3. Phylogenetic tree of EBLNs and modern bornaviral N proteins.** These trees  
31 were constructed by the maximum likelihood method using amino acid sequences of  
32 EBLN and modern bornaviral N genes of the genus *Carbovirus* (A), *Orthobornavirus* (B),  
33 or *Cultervirus* (C). Colored arrows indicate extant bornaviruses. Color of external nodes  
34 indicates the extant bornaviral genus or the host species in which the EBLN was identified,  
35 as shown in the lower right corner. Square or triangle labels on the internal nodes  
36 correspond to the collapsed nodes in **Figs. 3A-C**. Colored boxes highlight the bornaviral  
37 lineages endogenized during primate evolution. Asterisks on the branches indicate that  
38 the bootstrap value based on 1,000 replications is more than 80%. The scale bars show  
39 genetic distances (substitutions per site). The genetic distance to distinguish extant  
40 bornaviral species is shown as the comparative standard for estimating the genetic  
41 diversity of ancient bornaviruses.  
42  
43

**A**

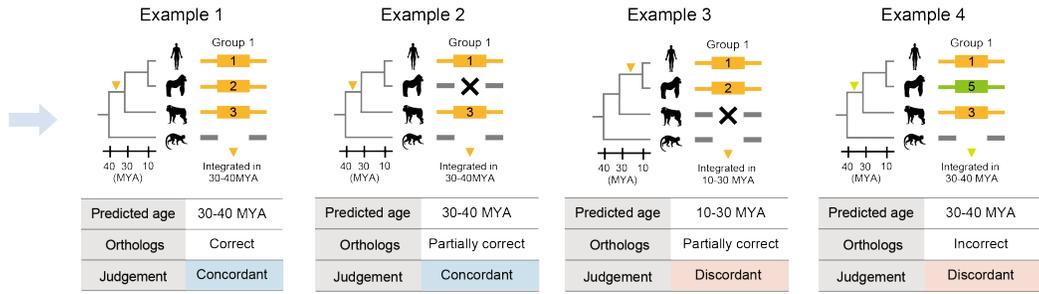
**1. Dating based on genomic alignment**



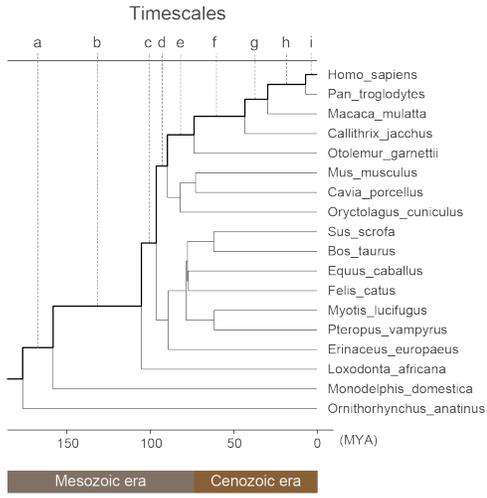
**2. Dating based on network analysis**



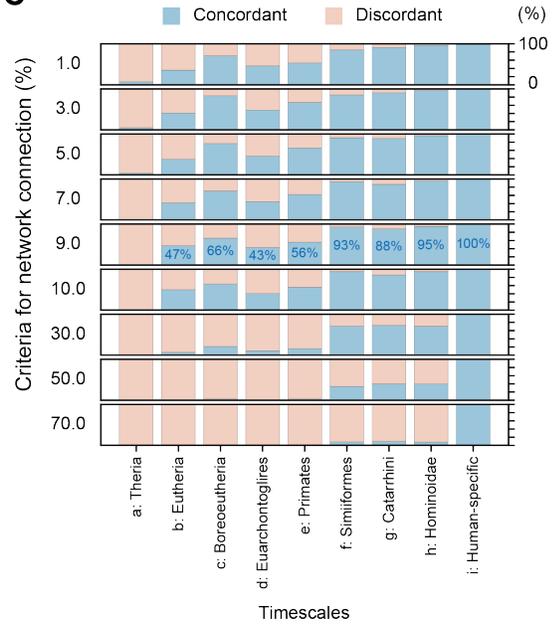
**3. Comparing results of group 1 (genomic alignment-based method vs. network analysis-based method)**



**B**



**C**



Supplemental Figure 4

46 **Fig. S4. Validation of the network-based method for orthologs detection using**  
47 **human transposable elements.** (A) Strategy for evaluating the detection rate of  
48 orthologs using our network-based method. To validate our network-based method, we  
49 compared it with the method of detecting orthologs using genomic alignments. First, we  
50 estimated the integration age of all human transposable elements (TEs) by LiftOver using  
51 the genomic alignment among 18 mammalian species shown in (B). Second, we randomly  
52 sampled 100 loci for each of the nine timescales, shown as a to i in (B), from the dating  
53 results of the genomic alignment-based method. Using these test datasets, we performed  
54 dating analysis by our network-based method. Third, we compared the results between  
55 the two methods by checking the predicted ages and detected orthologs. Example 1:  
56 integration ages coincided between two methods, and our method detected all orthologs  
57 defined by the genomic alignment-based method. Example 2: integration ages coincided  
58 between two methods, but our method detected an incomplete set of orthologs. Example  
59 3: integration ages were mismatched between the two methods. Example 4: estimation  
60 ages were matched between two methods, but there was a contamination of sequence  
61 unrelated to true orthologous relationships. Furthermore, to select the best criteria for our  
62 network-based dating analysis, we evaluated the nine different criteria shown in (C)  
63 **(details in Materials and Methods)**. (B) Phylogenetic tree of mammalian species used  
64 to detect orthologs of human TEs. Genomic alignments among these 18 species were  
65 obtained from the UCSC genome browser. To validate the network-based dating method  
66 for each timescale, we randomly sampled 100 TE loci from nine different timescales (a to  
67 i). (C) Concordant rates between the genomic alignment-based and network-based  
68 methods for estimating integration ages. Each panel shows the result using different  
69 criteria for network construction **(details in Materials and Methods)**. The x-axis indicates  
70 the timescales shown in (B). The y-axis indicates the concordant rate (%) between two  
71 methods. Blue labels indicate concordant rates (%) at the criteria used for the dating  
72 analysis for EBLs.  
73

74 **Dataset S1 (separate file). Genomic position of EBLs**  
75  
76 **Dataset S2 (separate file). Reference list for mammalian biogeography related to**  
77 **ancient bornaviral infections**  
78  
79 **Dataset S3 (separate file). Genetic distances between ancient and extant bornaviral**  
80 **N genes**  
81  
82 **Dataset S4 (separate file). Accession numbers of bornaviral sequences used for the**  
83 **tBLASTn search**  
84  
85 **Dataset S5 (separate file). Chain files and genome assemblies used to validate for**  
86 **our dating method**  
87  
88 **Dataset S6 (separate file). Extant viral sequences used for phylogenetic analyses**  
89  
90 **Dataset S7 (separate file). Bioinformatics tools used in this study**  
91  
92