1 **Supplementary Text S1**

2 **Dataset of reference viral genomes** Viral genome sequences and taxonomic

3 information of viruses and their hosts are based on the GenomeNet/Virus-Host DB

4 (Mihara, et al., 2016). The Virus-Host DB covers viruses with complete genomes stored

5 in RefSeq/viral and GenBank. The sequence and taxonomic data are downloadable at

6 http://www.genome.jp/virushostdb/. Viral genome sequences and taxonomic

7 information used in ViPTree will be updated every few months following Virus-Host

8 DB updates.

9 **Calculation of genomic similarity ($S_G$) and proteomic tree** The original Phage

10 Proteomic Tree (Rohwer and Edwards, 2002) tested multiple approaches and

11 parameters to calculate genomic distance for evaluation of compatibility between the

12 Phage Proteomic Tree and taxonomical system of the International Committee on

13 Taxonomy of Viruses. In contrast, ViPTree performs proteomic tree calculation based

14 on tBLASTx as reported previously by other several studies (Bellas, et al., 2015;

15 Bhunchoth, et al., 2016; Mizuno, et al., 2013). Specifically, normalized tBLASTx scores

16 ($S_G$; $0 \leq S_G \leq 1$) between viral genomes are calculated as described (Bhunchoth, et al.,

17 2016). For simple calculation of $S_G$ of segmented viruses, sequences of each segmented

18 viral genome were concatenated into one sequence by inserting a sequence of 100

19 ambiguous nucleotides (N) at each concatenation site. In addition, genome sequences

20 longer than 100 kb are split into each 100 kb fragment just for faster tBLASTx

21 calculation. Genomic distance is calculated as $1-S_G$. Proteomic tree generation is

22 performed by BIONJ (Gascuel, 1997) based on the genomic distance, using the R

23 package 'Ape'. The root of a tree is defined by mid-point rooting.

24 **Gene finding and gene annotation** Gene finding of uploaded sequences is performed

25    by GeneMarkS (Besemer, et al., 2001) without using its self-training method and with

26    using an option "-p 0" to prohibit gene overlaps. Automatic functional annotations of

27    predicted genes are performed by protein similarity searches against NCBI/nr using

28    GHOSTX (Suzuki, et al., 2014).

29

30    **References**

31    Bellas, C.M., Anesio, A.M. and Barker, G. Analysis of virus genomes from glacial

32    environments reveals novel virus groups with unusual host interactions. *Front Microbiol*

33    2015;6:656.

34    Besemer, J., Lomsadze, A. and Borodovsky, M. GeneMarkS: a self-training method for

35    prediction of gene starts in microbial genomes. Implications for finding sequence motifs in

36    regulatory regions. *Nucleic Acids Res* 2001;29(12):2607-2618.

37    Bhunchoth, A*., et al.* Two asian jumbo phages, ΦRSL2 and ΦRSF1, infect Ralstonia

38    solanacearum and show common features of ΦKZ-related phages. *Virology* 2016;494:56-66.

39    Gascuel, O. BIONJ: an improved version of the NJ algorithm based on a simple model of

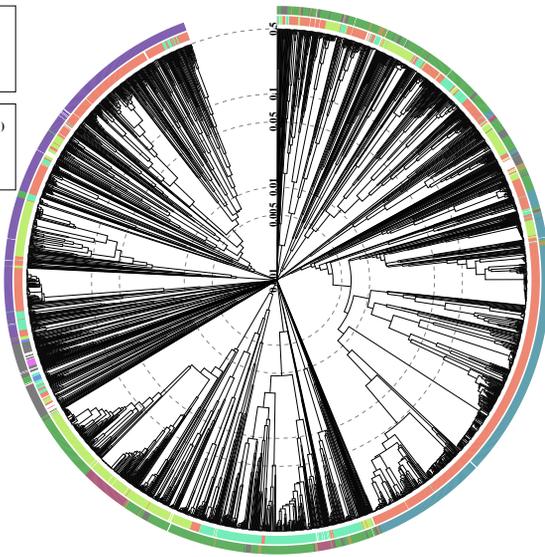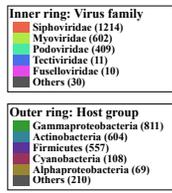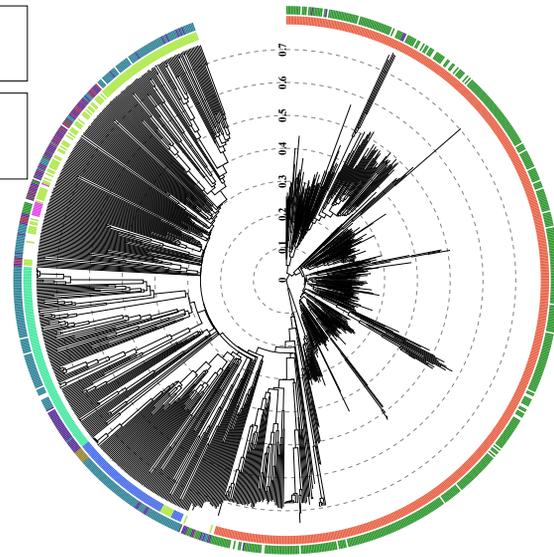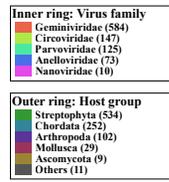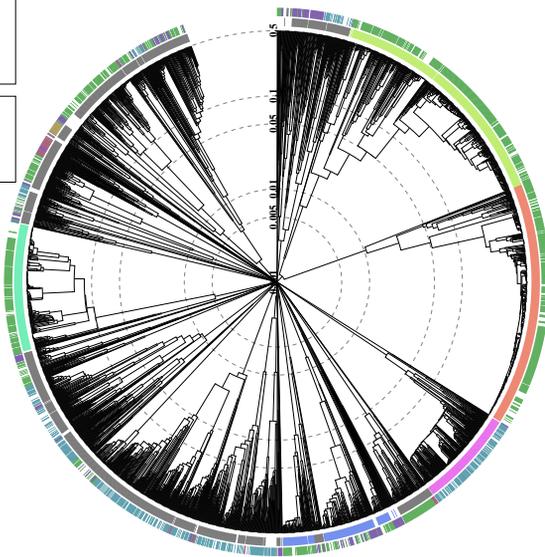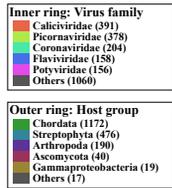40    sequence data. *Mol Biol Evol* 1997;14(7):685-695.

41    Mihara, T*., et al.* Linking Virus Genomes with Host Taxonomy. *Viruses* 2016;8(3):66.

42    Mizuno, C.M*., et al.* Expanding the marine virosphere using metagenomics. *PLoS Genet*

43    2013;9(12):e1003987.

44    Rohwer, F. and Edwards, R. The Phage Proteomic Tree: a genome-based taxonomy for phage.

45    *J Bacteriol* 2002;184(16):4529-4535.

46    Suzuki, S*., et al.* GHOSTX: an improved sequence homology search algorithm using a query

47    suffix array and a database suffix array. *PLoS One* 2014;9(8):e103833.
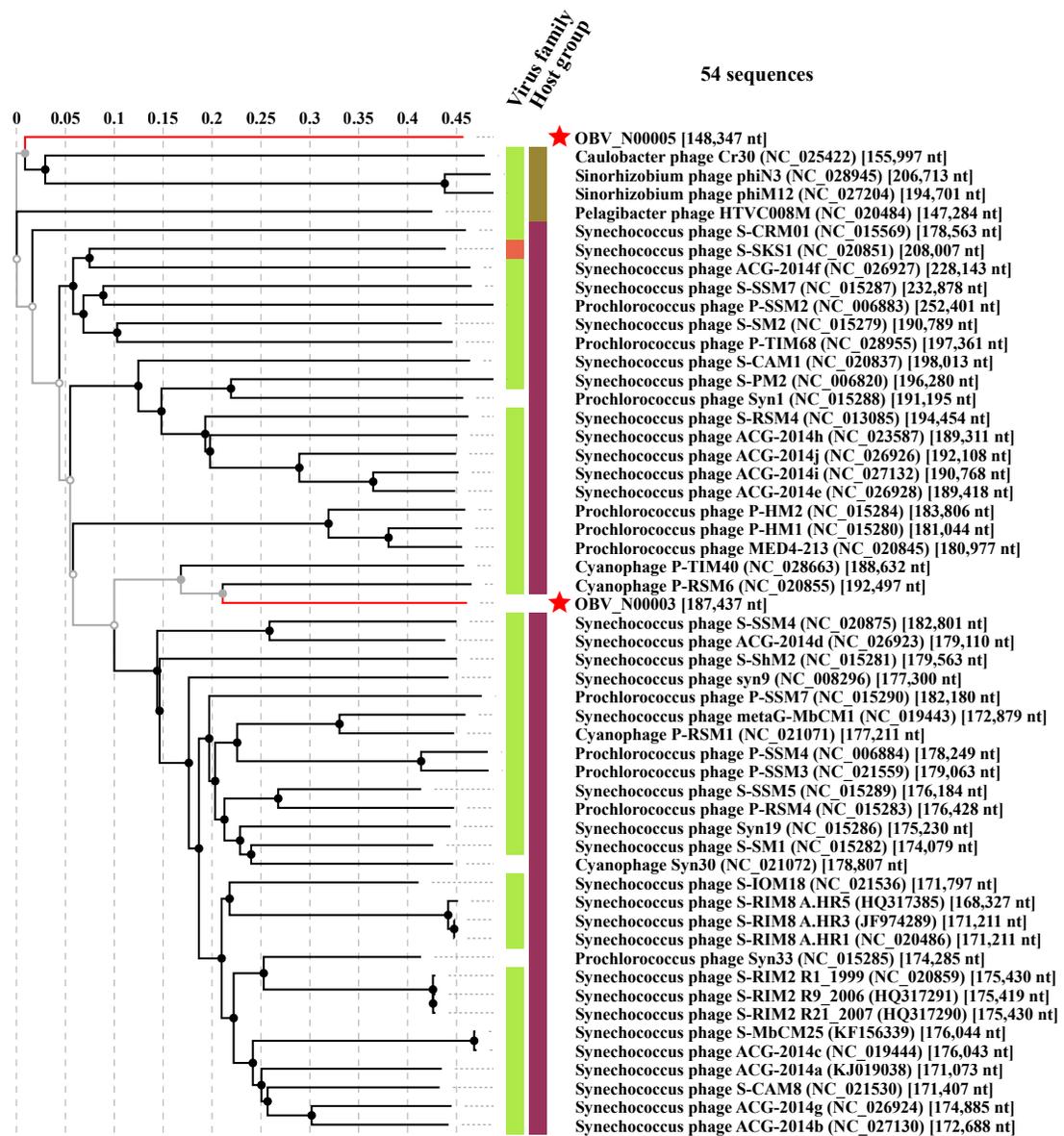
48

**A**

Inner ring: Virus family
- Siphoviridae (1214)
- Myoviridae (602)
- Podoviridae (409)
- Tectiviridae (11)
- Fuselloviridae (10)
- Others (30)

Outer ring: Host group
- Gammaproteobacteria (811)
- Actinobacteria (604)
- Firmicutes (557)
- Cyanobacteria (108)
- Alphaproteobacteria (69)
- Others (210)

**B**

Inner ring: Virus family
- Geminiviridae (584)
- Circoviridae (147)
- Parvoviridae (125)
- Anelloviridae (73)
- Nanoviridae (10)

Outer ring: Host group
- Streptophyta (534)
- Chordata (252)
- Arthropoda (102)
- Mollusca (29)
- Ascomycota (9)
- Others (11)

**C**

Inner ring: Virus family
- Caliciviridae (391)
- Picornaviridae (378)
- Coronaviridae (204)
- Flaviviridae (158)
- Potyviridae (156)
- Others (1060)

Outer ring: Host group
- Chordata (1172)
- Streptophyta (476)
- Arthropoda (190)
- Ascomycota (40)
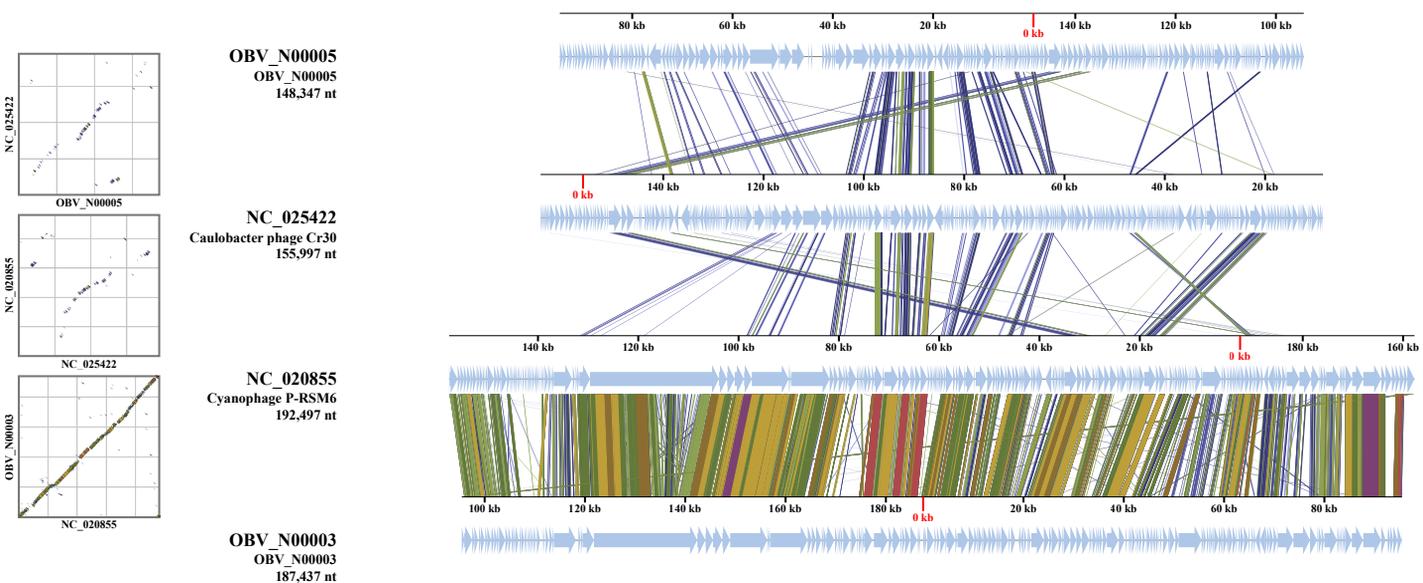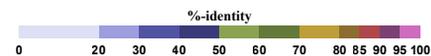- Gammaproteobacteria (19)
- Others (17)

**Supplementary Figure S1.** Viral proteomic trees of reference genomes represented in the circular view. Color rings indicate virus families (inner rings) and host groups (at a level of phylum except for Proteobacteria; outer rings). These trees are calculated by BIONJ based on genomic distance matrixes, and mid-point rooted. Branch lengths are log-scaled (A, C) and linearly scaled (B). The sequence and taxonomic data is based on Virus-Host DB (release 78). (A) A tree of prokaryotic dsDNA viruses (2,384 genomes). (B) A tree of eukaryotic ssDNA viruses (1,013 genomes). (C) A tree of all ssRNA viruses (2,485 genomes).

**Supplementary Figure S2.** A Proteomic tree of dsDNA viruses represented in the rectangular view. This tree includes 52 reference viral genomes (48 *Myoviridae*, one *Siphoviridae*, and three unclassified viruses) and two uploaded viral genomes that are highlighted by red branches and stars (OBV_N00003 and OBV_N00005; BioProject: PRJDB4437). The tree is constructed by BIONJ based on genomic distance matrixes, and mid-point rooted. Branch lengths are linearly scaled. In the ViPTree server, where inner nodes of a tree are shown as filled circles, each of them links to a genomic alignment of sequences included in its subtree.

**Supplementary Figure S3.** An example of the genomic alignment view of four viral genomes (two reference and two environmental viral genomes) that are included in Supplementary Figure S2. Pairwise dot plots of these genomes are also shown. Colored lines in the alignment and the dot plots indicate tBLASTx results (E-value < 1e-2). Grid lines in the dot plots indicate 40 kb intervals. Positions of each sequences are automatically adjusted (i.e., circular permuted and reverse stranded) for clear representation of colinearity between genomes.