

1 Supplementary Information for
2
3 Biological Species in the Viral World

4
5 Louis-Marie Bobay and Howard Ochman

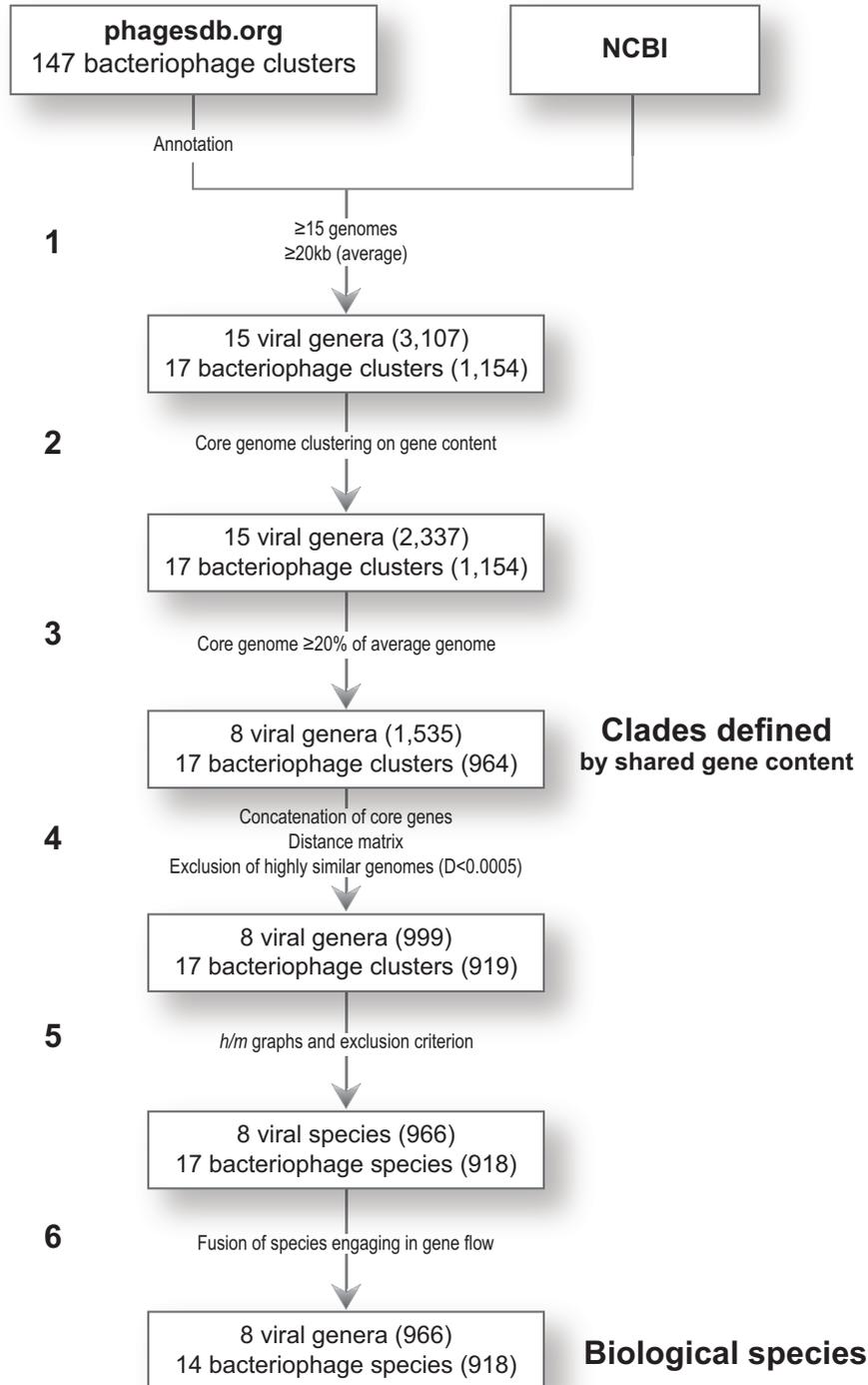
6
7 Louis-Marie Bobay
8 Email: ljbobay@uncg.edu

9
10
11 **This PDF file includes:**

12
13 Figs. S1 to S8
14 SI References

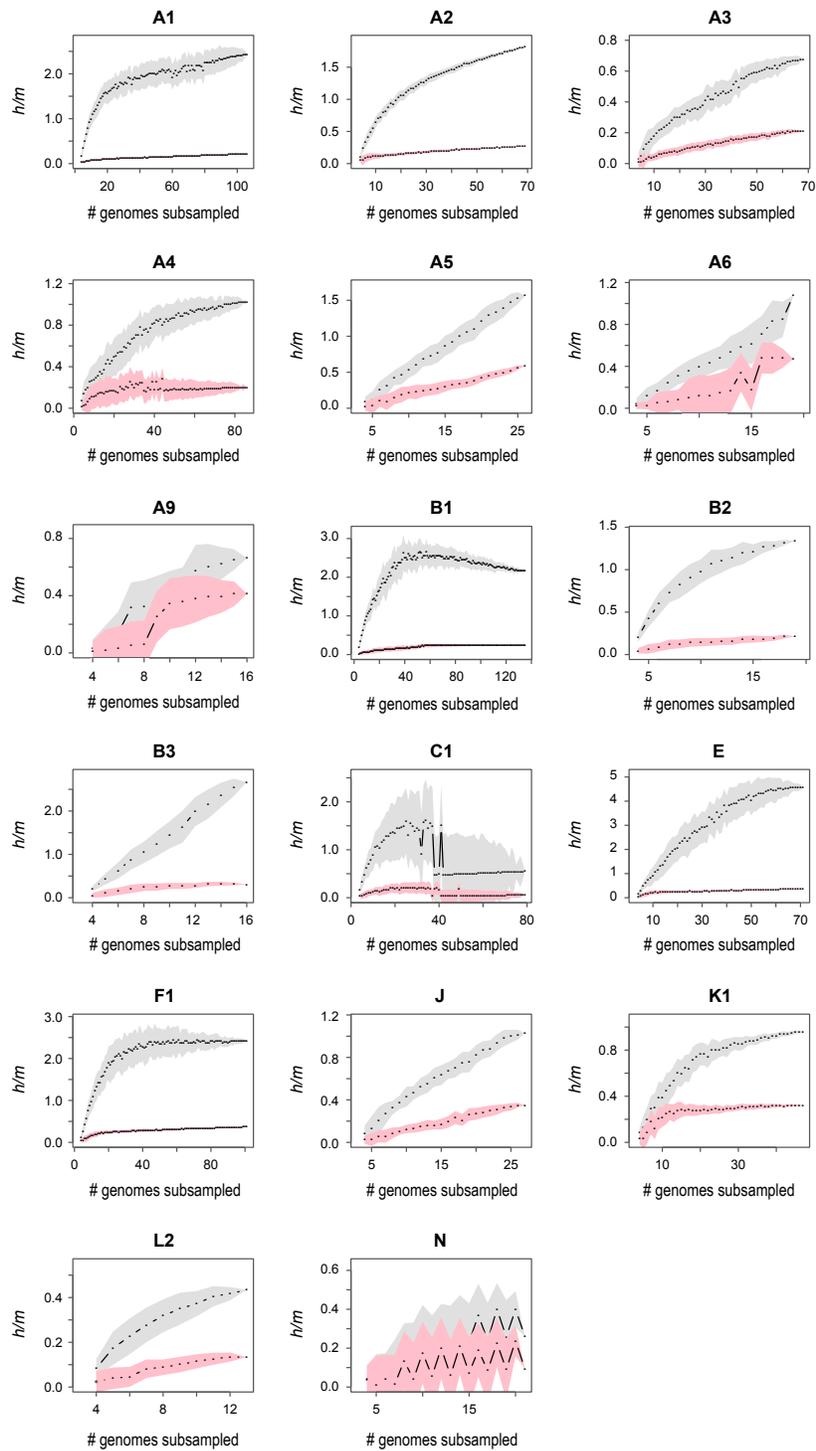
15
16 **Other supplementary materials for this manuscript include the following:**

17
18 Datasets S1 to S2
19



20
21

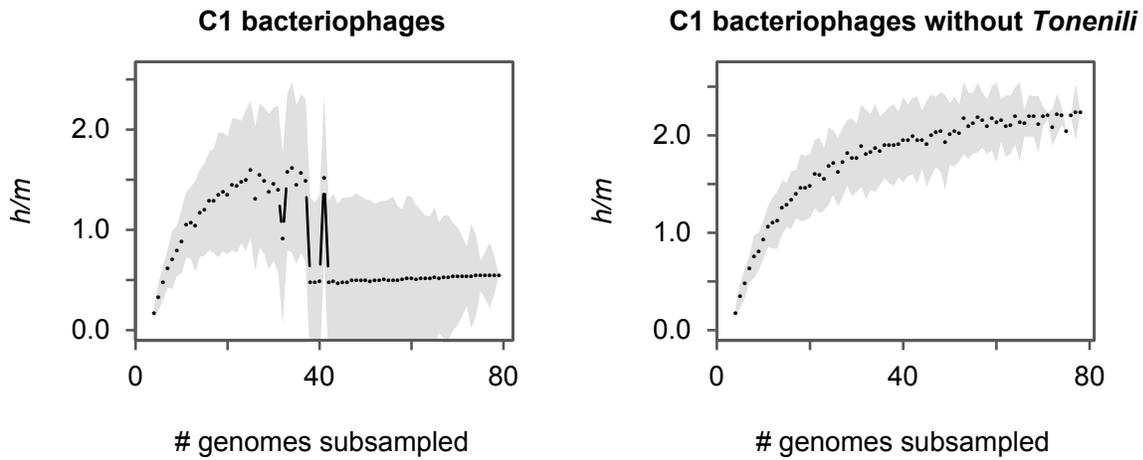
22 **Figure S1. Flow chart depicting each analytical step in defining biological species.** Steps are
23 numbered one to six and correspond to those described in the Methods section. Numbers of viral
24 and bacteriophage clades remaining after each filtering step are indicated within each box, and
25 numbers in parentheses are the total number of viral and bacteriophage genomes at each step.



26

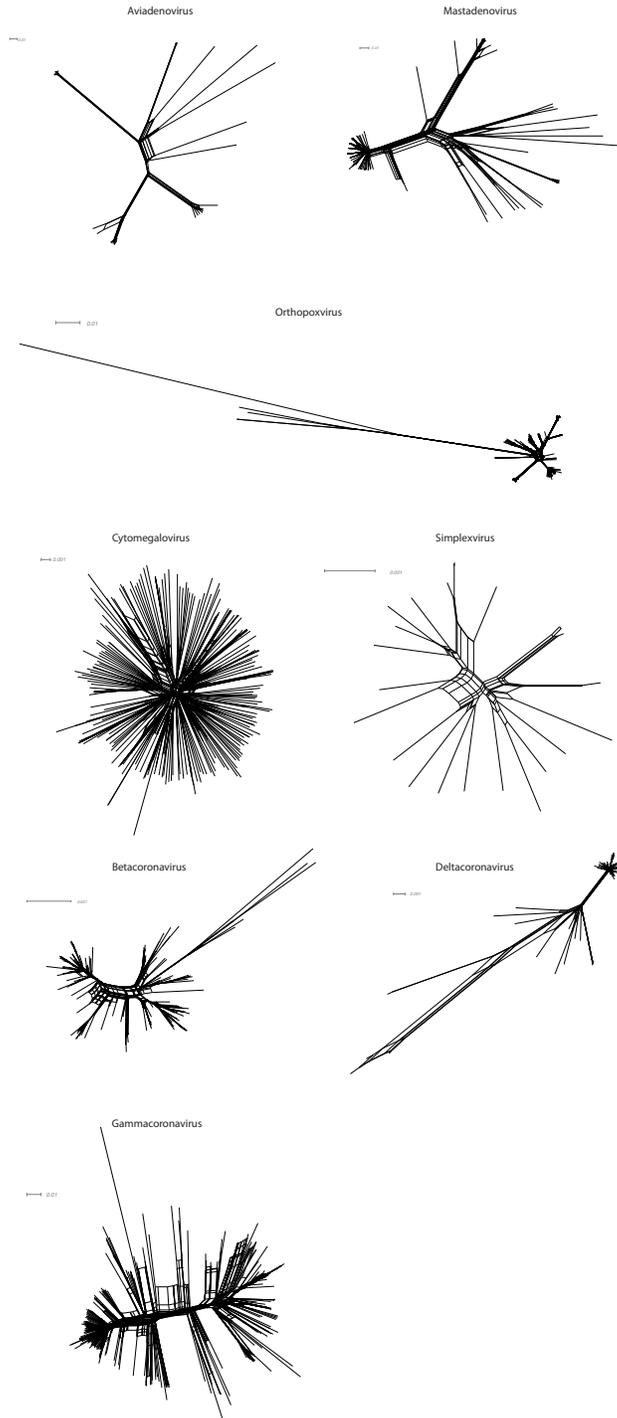
27 **Figure S2. Patterns of gene flow in bacteriophages infecting *Mycobacterium smegmatis*.** For
 28 each cluster, gene flow was estimated by the ratio of homoplastic to non-homoplastic alleles (h/m)
 29 with a re-sampling strategy. For each number i of genomes, 100 combinations of i genomes were
 30 randomly sampled and h/m was computed for each combination. Within the bivariate plots, black
 31 dots are medians and the grey-shaded region is the standard deviation for the indicated number

32 of subsampled combinations of strains, and red dots and pink-shaded regions denote median h/m
33 values and standard deviation for simulations in which all homoplasies are introduced by
34 convergent mutations, as described in the text. Differences between the distributions of observed
35 and simulated h/m values indicate the extent to which homoplasies are introduced by
36 recombination.



37

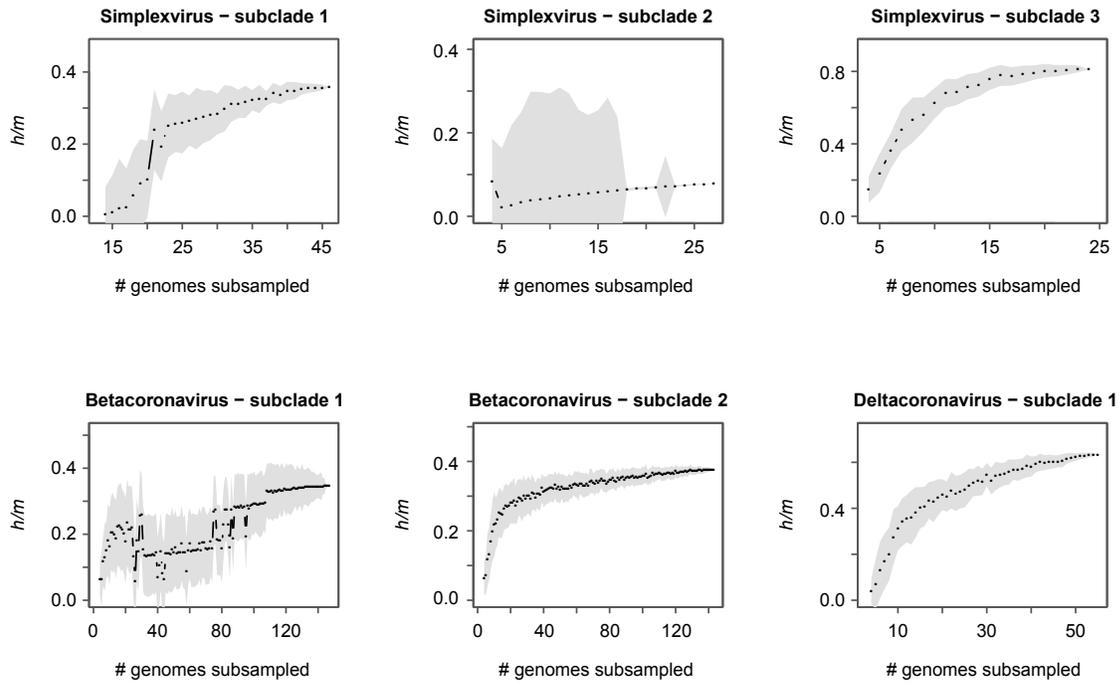
38 **Figure S3. Redefining species membership in bacteriophage cluster C1.** The discontinuity
 39 detected in the graph for the entire set of genomes classified in bacteriophage cluster C1
 40 indicated the presence of multiple species (left). After removal of the sexually isolated genome
 41 (bacteriophage *Tonenili*), the graph was rebuilt (right). Black dots represent the median and the
 42 grey area indicates the standard deviation of h/m for the different combinations of genomes.



43

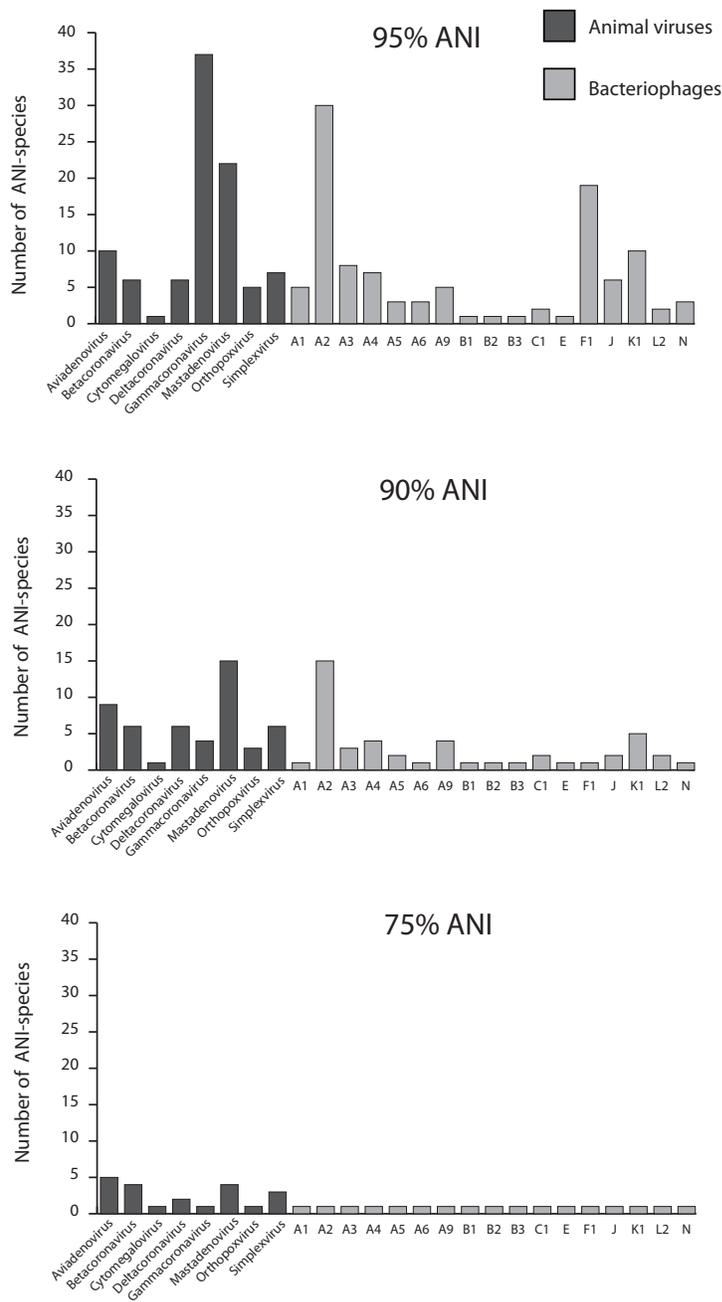
44 **Figure S4. Phylogenetic networks of viral species.** Phylogenetic networks were built using the
 45 core genome of each biological species of viruses with SplitsTree (1). The scale beside each
 46 phylogenetic network indicates substitution rate.

47



52
 53 **Figure S6. Patterns of gene flow in subclades of viral genera.** Subclades within the three viral
 54 genera *Simplexvirus*, *Betacoronavirus* and *Deltacoronavirus* that each presented low signals of
 55 gene flow were tested independently for a signal of gene flow. Subclades were defined from
 56 phylogenetic trees that included all members of a genus, and analyses proceeding progressively
 57 examining smaller subclades, such that subclade 3 is included in subclade 2, which is included in
 58 subclade 1. Black dots represent the median and the grey area indicates the standard deviation of
 59 h/m of the different combinations of genomes.

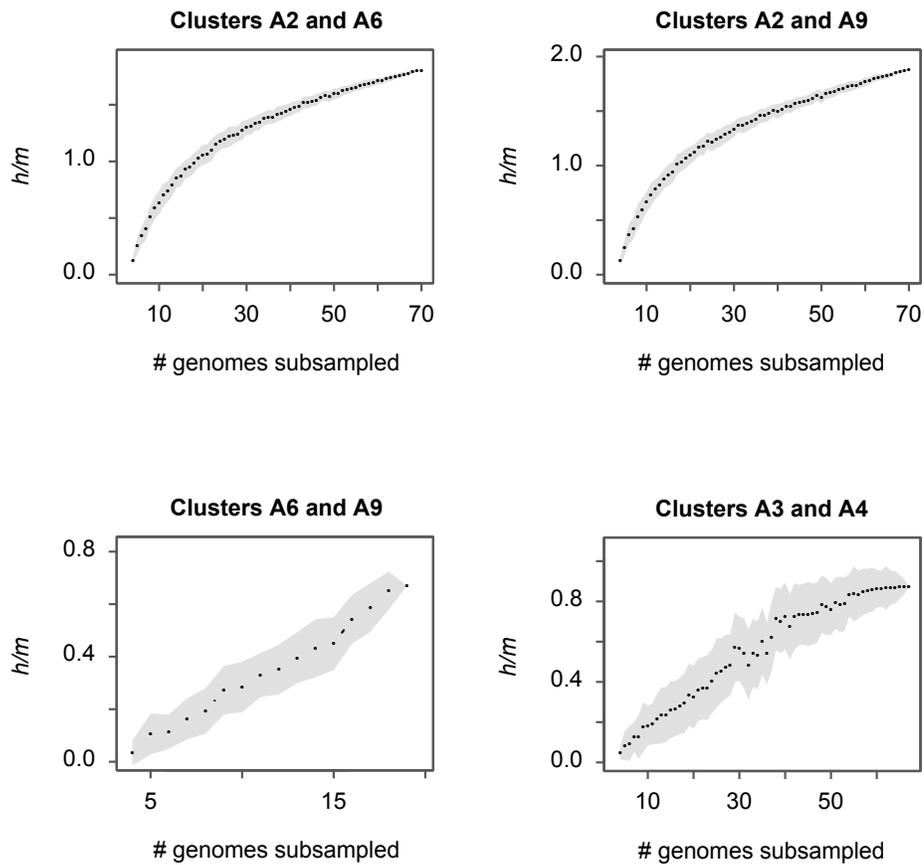
60



61

62 **Figure S7. Number of ANI-species obtained at different sequence-identity thresholds.**

63 Average nucleotide identity (ANI) was computed along the entire core genome for members of
 64 the same biological species. Shown are numbers of groupings (ANI-species) obtained at various
 65 sequence-identity thresholds, selected as follows: 95%, threshold recommended for defining
 66 bacterial species; 90%, threshold recommended for defining bacterial genera; 75%, threshold
 67 representing the maximal divergence observed in biological species of bacteria. Note that the
 68 different sequence-identity thresholds do not partition each biological species in a uniform
 69 manner and that the maximal divergence observed between members of some viral species
 70 exceeds that observed in bacteria.



71

72 **Figure S8. Patterns of gene flow between bacteriophage clusters.** Patterns of gene flow were
 73 analyzed for pairs of clusters with large numbers of shared homologs ($n \geq 16$). Black dots
 74 represent the median and the grey area indicates the standard deviation of h/m of the different
 75 combinations of genomes.

76

77 **SI References**

78

- 79 1. Huson DH, Bryant D (2006) Application of phylogenetic networks in evolutionary
80 studies. *Mol Biol Evol* 23(2):254-267.