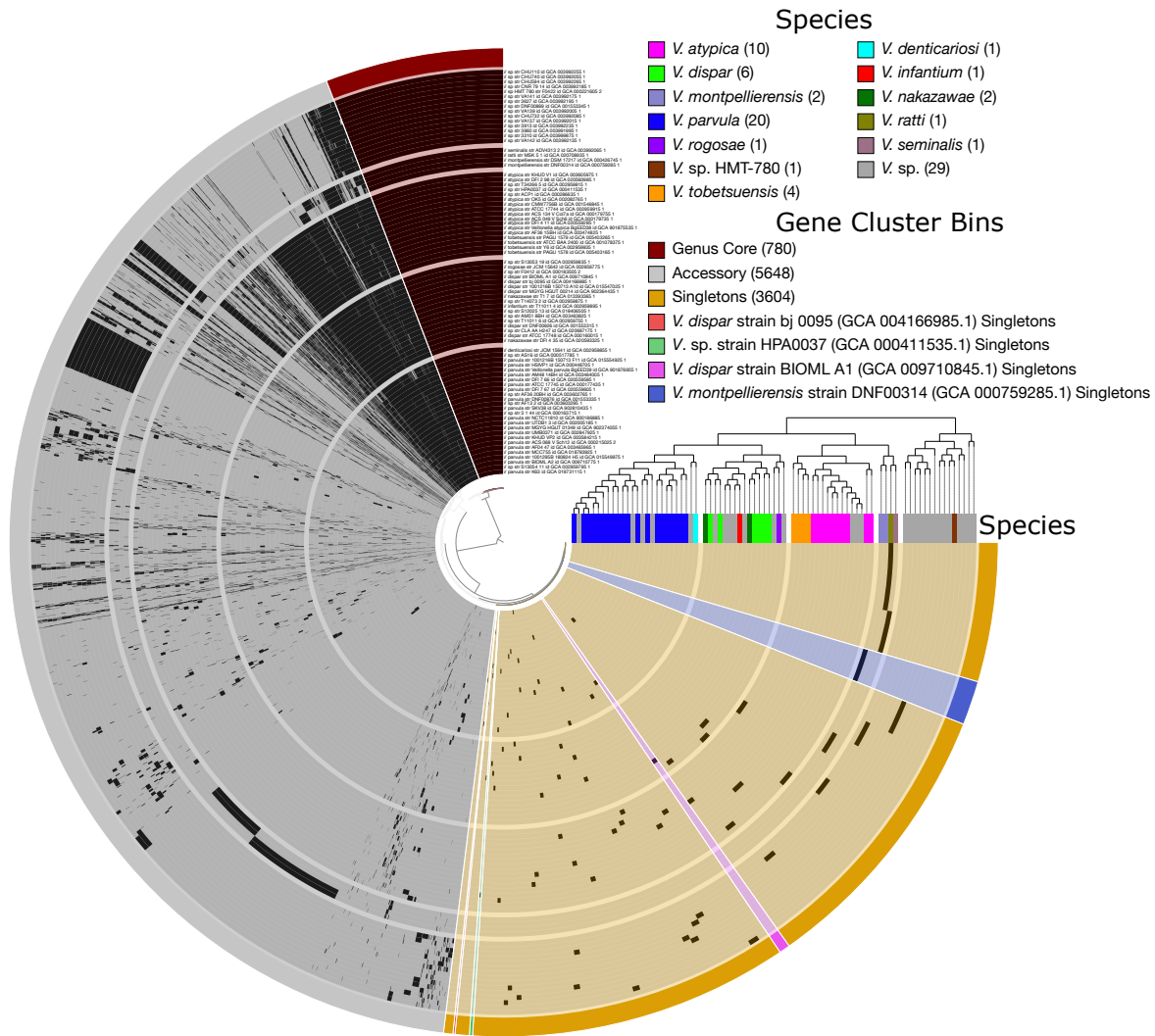2

3   **Figure S1**. A pangenome of human oral *Veillonella* reference genomes (n = 79) was constructed

4   using a set genomes dereplicated based on a 99% ANI threshold to test the effect of a 98% ANI

5   threshold on a comparative pangenome analysis. Genomes are hierarchically clustered based on

6   gene cluster frequency (i.e., the number of representatives of each gene cluster present in each

7   genome). Gene clusters are arranged based on their presence or absence across the genomes and

8   colored according to their presence in all genomes (core; red; n = 780), a subset of genomes

9   (accessory; gray; n = 5648), or unique to a single genome (singletons; gold; n = 3604). Singleton

10    gene clusters unique to a genome that was added because of increasing the ANI threshold are
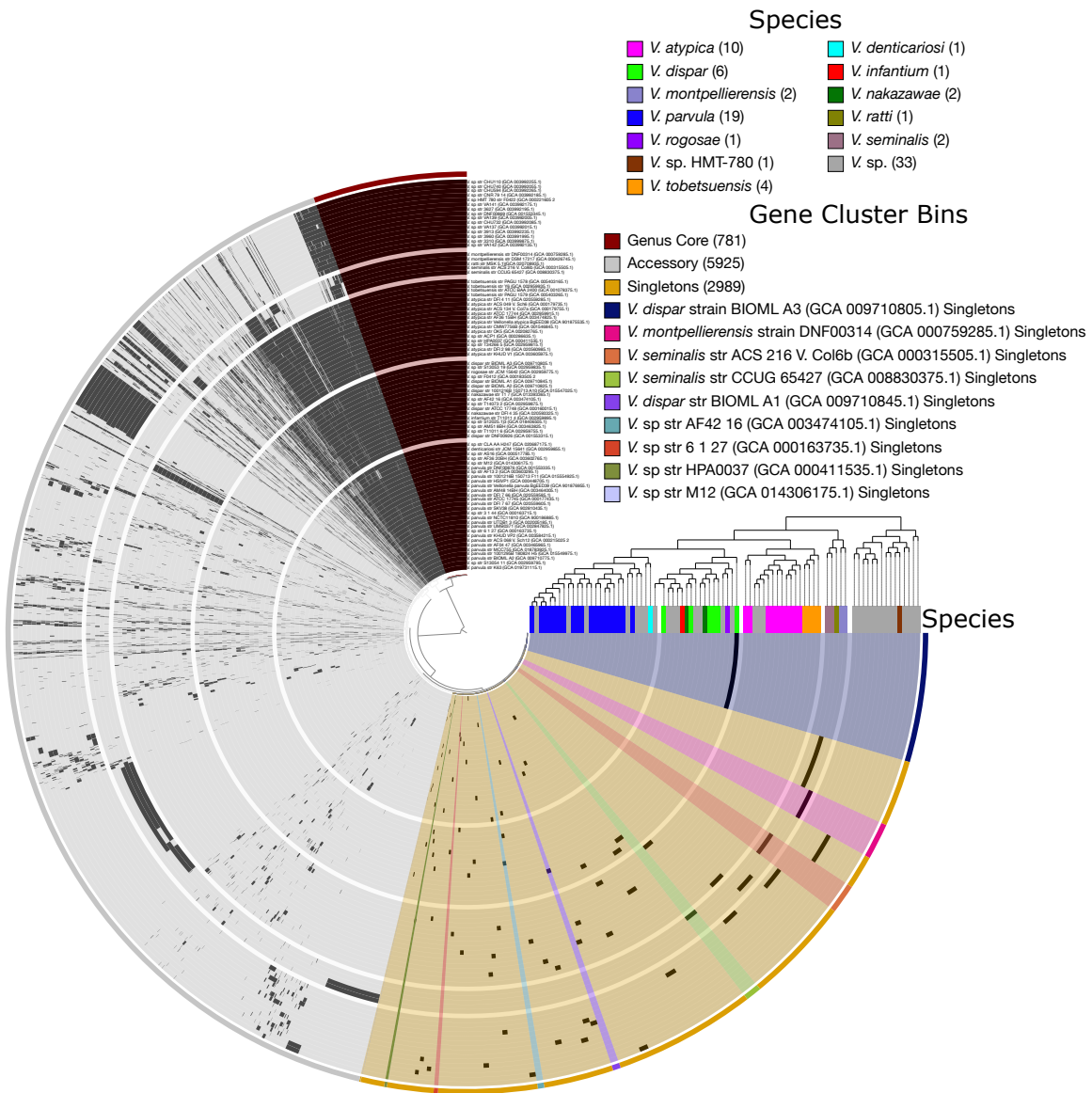
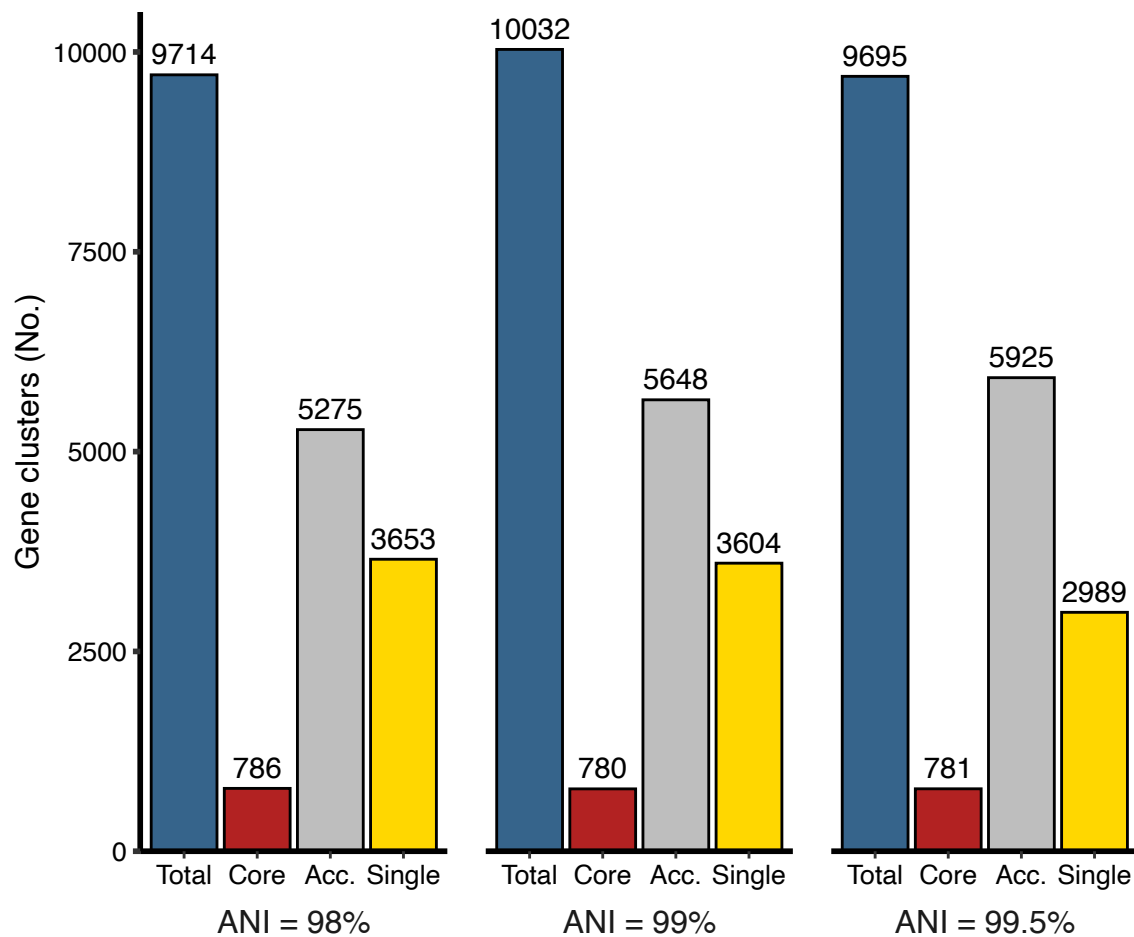11    colored uniquely.

12

13

14

15

16

17

**Figure S2**. A pangenome of human oral *Veillonella* reference genomes (n = 82) was constructed

using a set genomes dereplicated based on a 99.5% ANI threshold to test the effect of a 98%

ANI threshold on a comparative pangenome analysis. Genomes and gene clusters are arranged as

in Figure S1.

24

**Figure S3**. Comparison of the number of gene cluster bins for three pangenomes of human

*Veillonella* species that were constructed using a set of genomes dereplicated based on a 98%,

99% or a 99.5% ANI threshold. The distinct gene clusters of each pangenome include core genes

that occur in every genome (red), singleton genes that occur in only a single genome (gold) and

accessory genes that occur in more than one, but not all genomes (gray).

30

31

32

33

34    **SUPPLEMENTAL TABLE LEGENDS**

35

36    **Table S1.** Metadata for *Veillonella* reference genomes.

37

38    **Table S2.** Whole-genome average nucleotide percent identity for *Veillonella* reference genomes.

39

40    **Table S3.** Metabolic pathway enrichment test results.

41

42    **Table S4.** Results for enrichment of COG20, Pfams and KEGG associations.

43

44    **Table S5.** ANI dereplication test blast results. We analyzed how additional genomes altered

45    accessory gene content by inspecting the singleton gene clusters unique to each of the additional

46    genomes. For each pangenome, we manually binned singleton gene clusters for each of the

47    additional genomes, extracted their amino acid sequences using anvi-get-sequences-for-gene-

48    clusters and blasted each sequence against the NCBI non-redundant protein database (1), which

49    includes protein sequences from GenPept, Swissprot, PIR, PDF, PDB, and NCBI RefSeq.

50

51    **REFERENCES**

52    1.  O'Leary NA, Wright MW, Brister JR, Ciufo S, Haddad D, McVeigh R, Rajput B, Robbertse

53        B, Smith-White B, Ako-Adjei D. 2016. Reference sequence (RefSeq) database at NCBI:

54        current status, taxonomic expansion, and functional annotation. Nucleic Acids Res 44:D733–

55        D745.

56