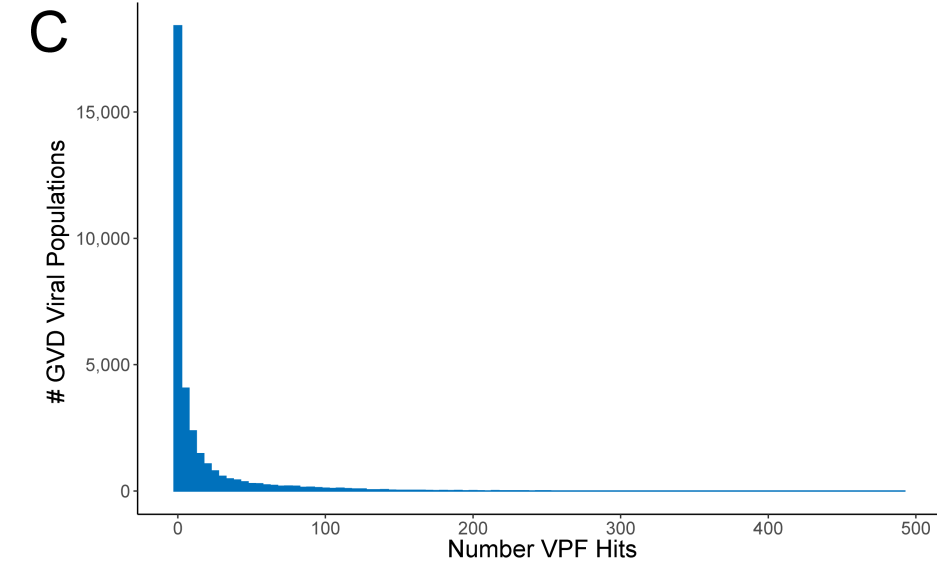
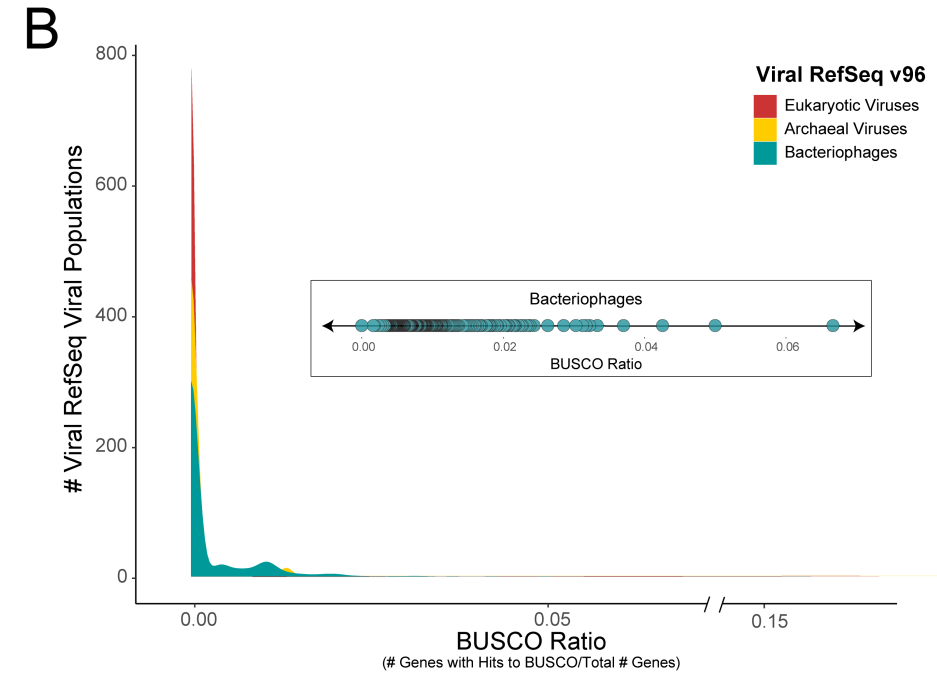
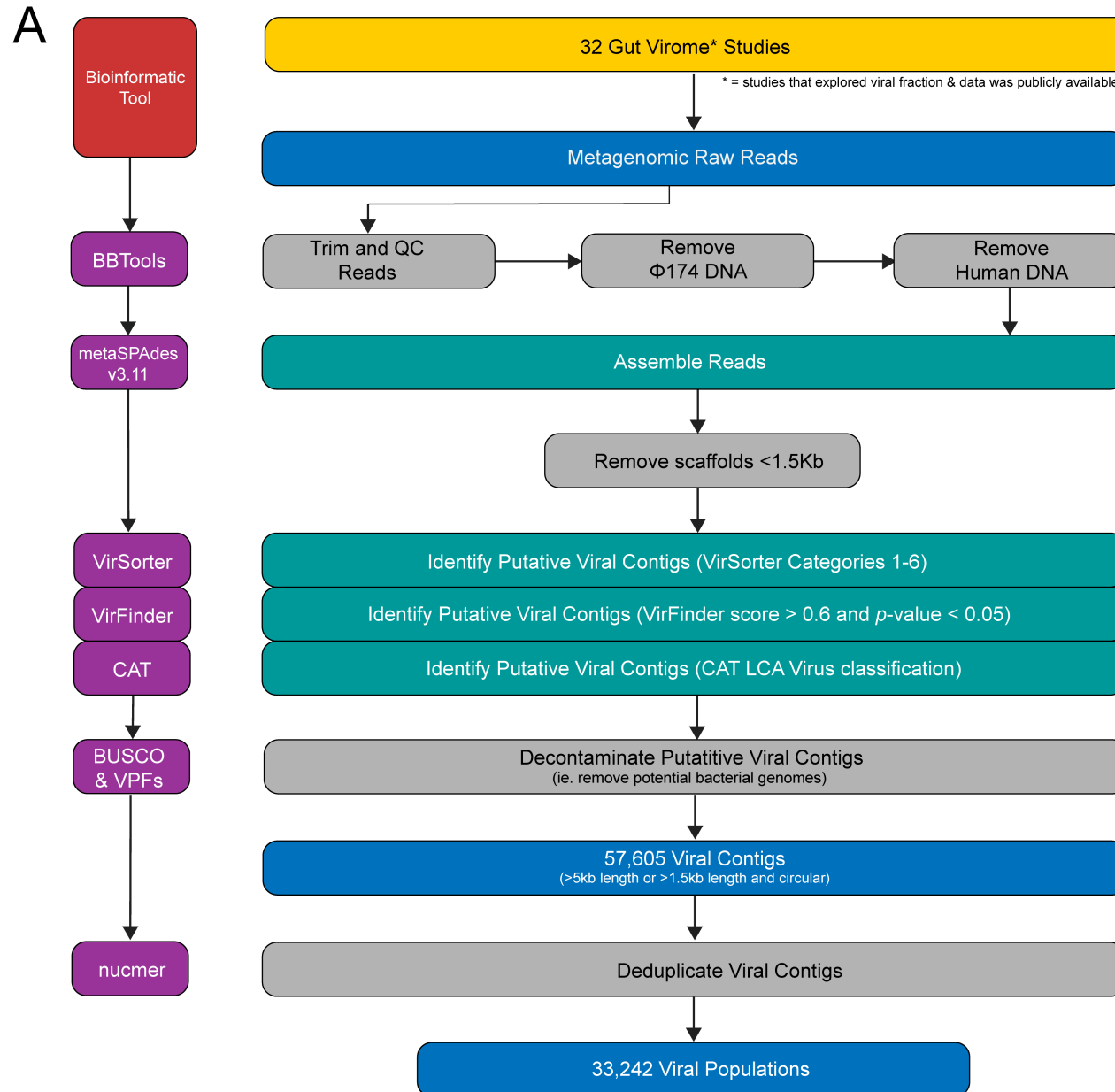


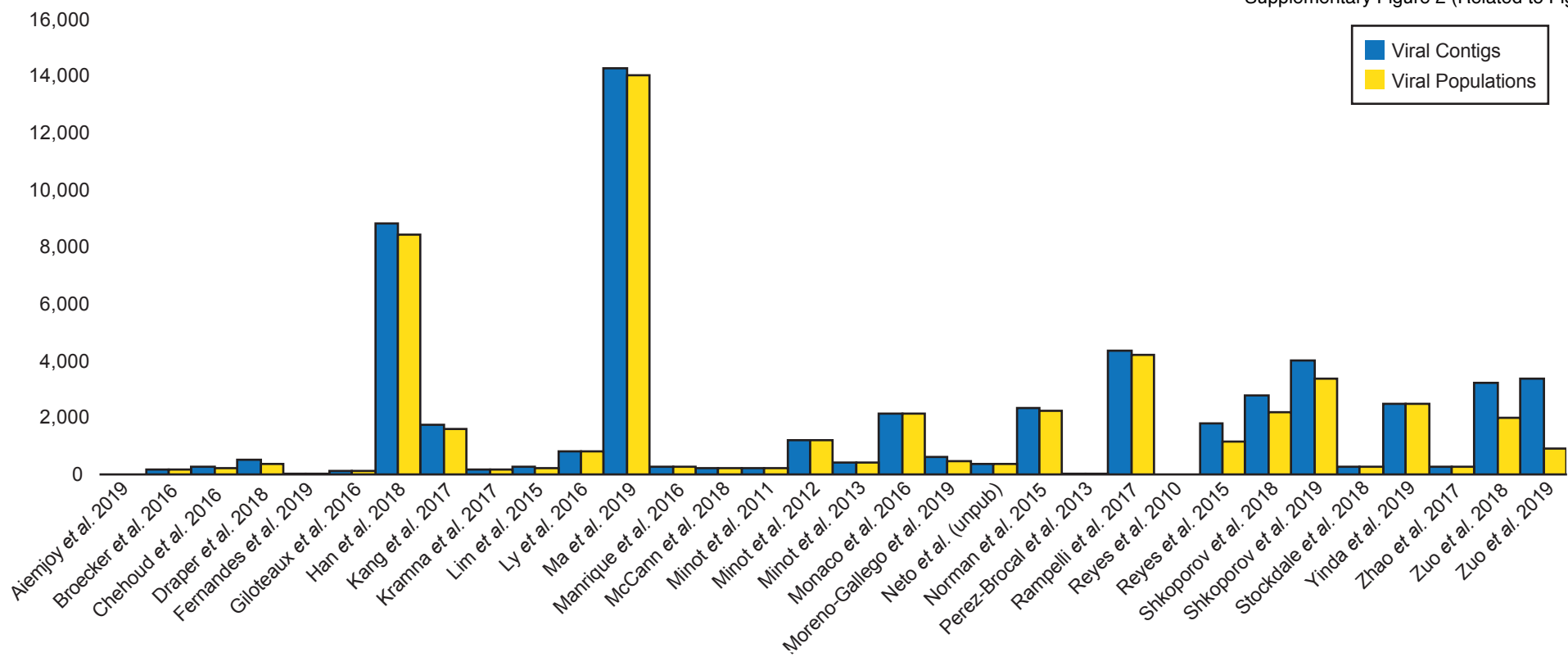
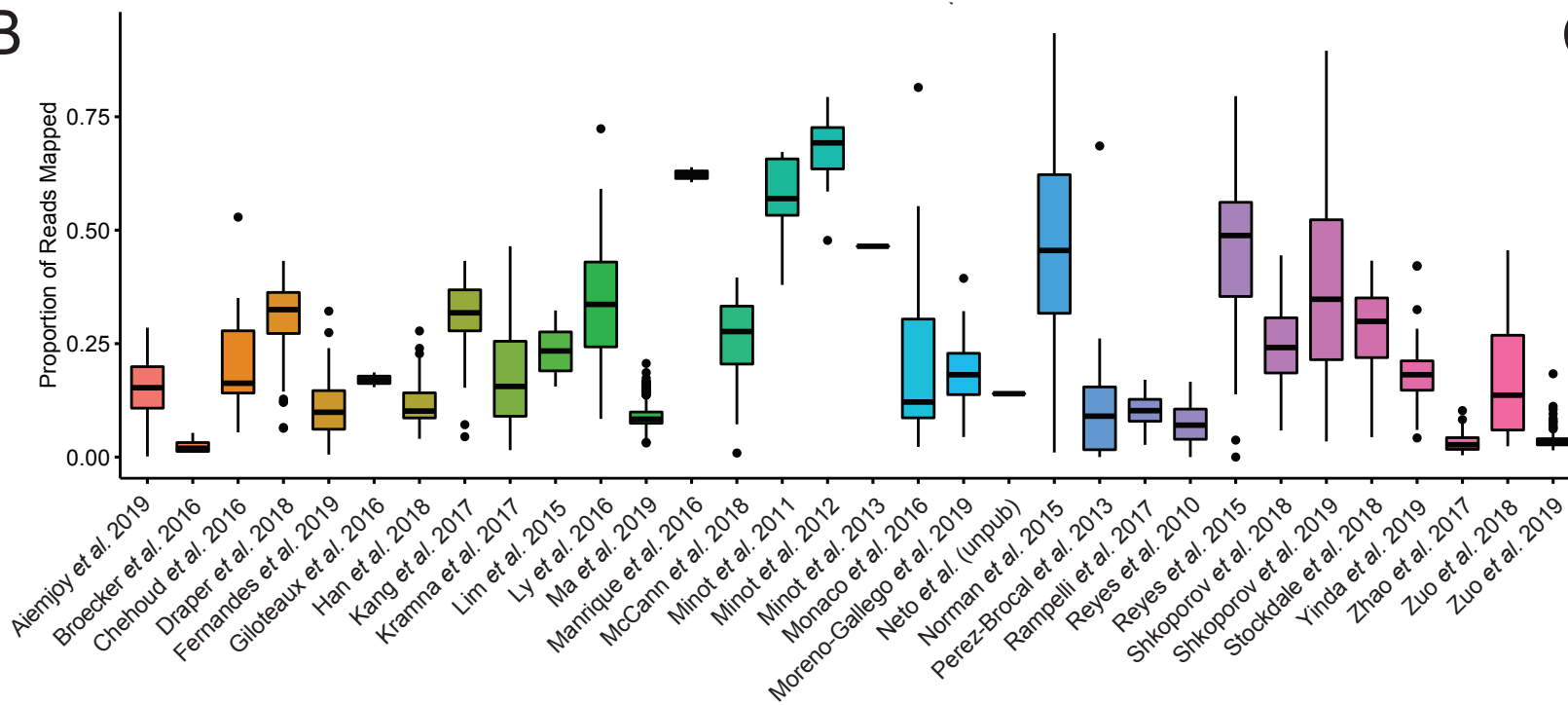
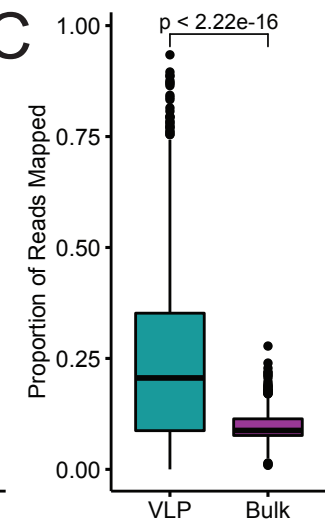
Cell Host & Microbe, Volume 28

Supplemental Information

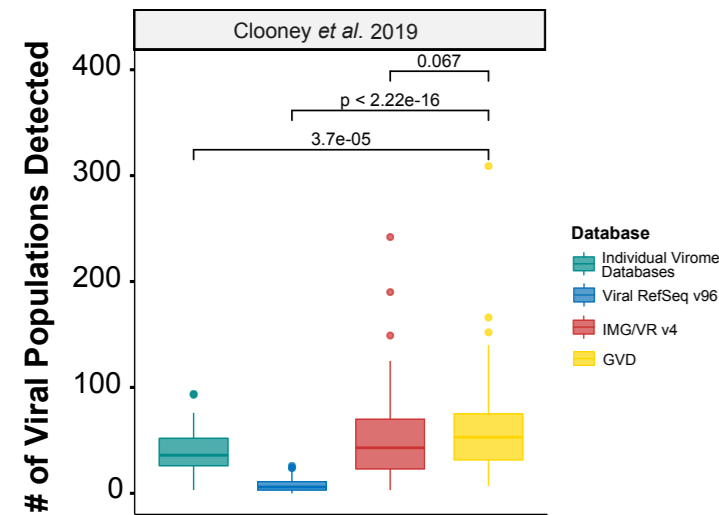
**The Gut Virome Database Reveals Age-Dependent
Patterns of Virome Diversity in the Human Gut**

Ann C. Gregory, Olivier Zablocki, Ahmed A. Zayed, Allison Howell, Benjamin Bolduc, and Matthew B. Sullivan

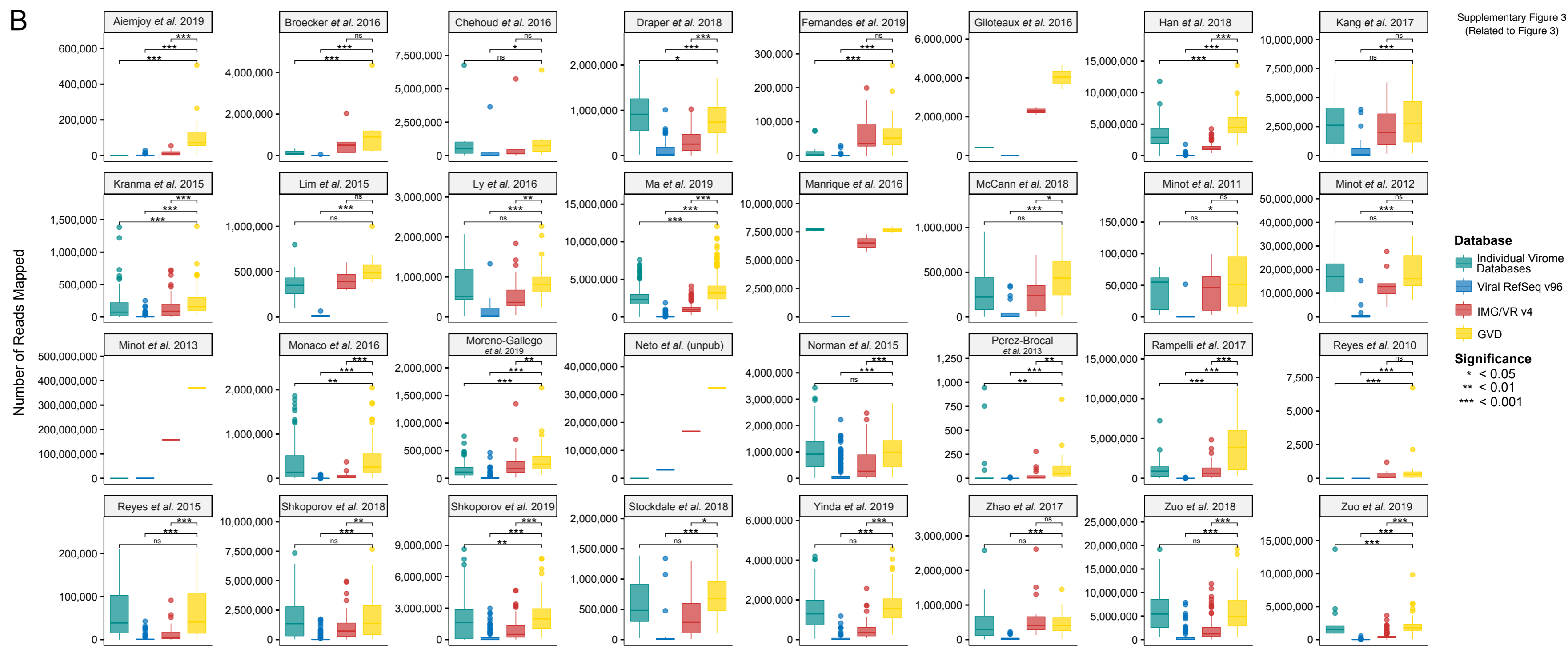


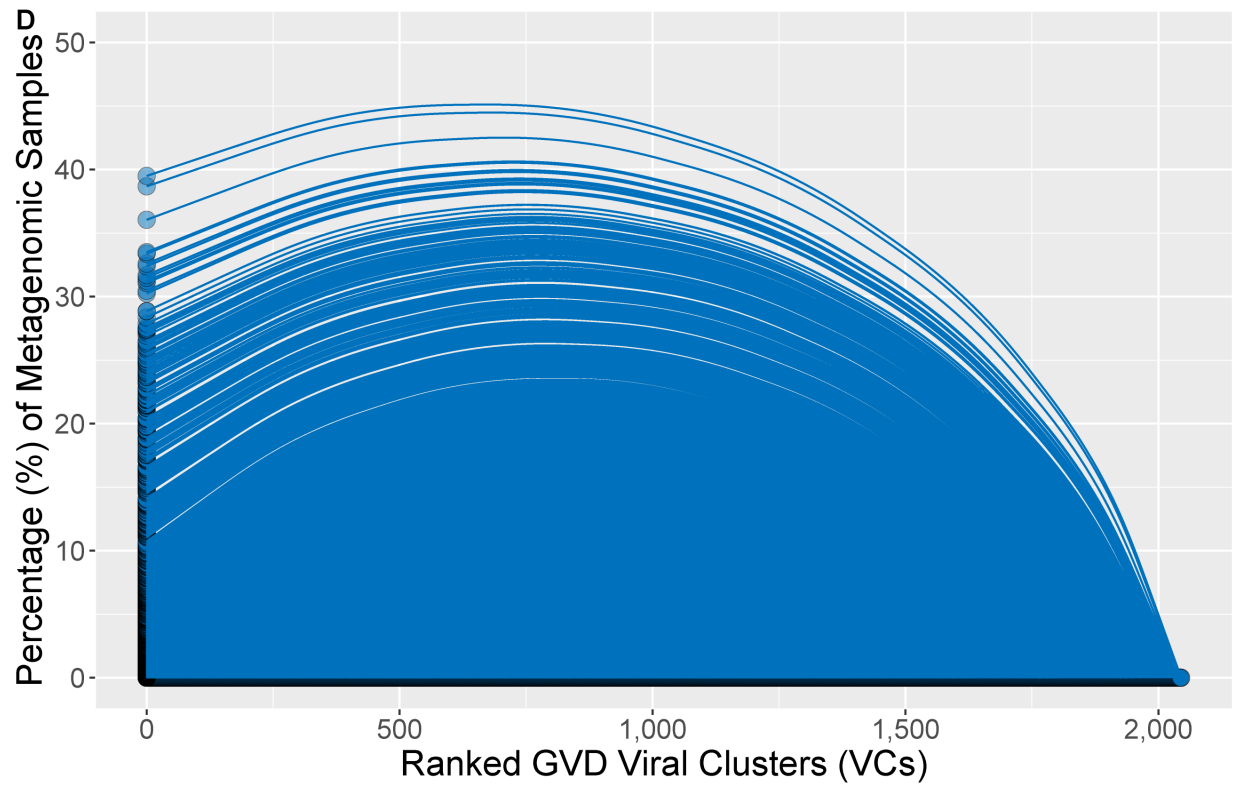
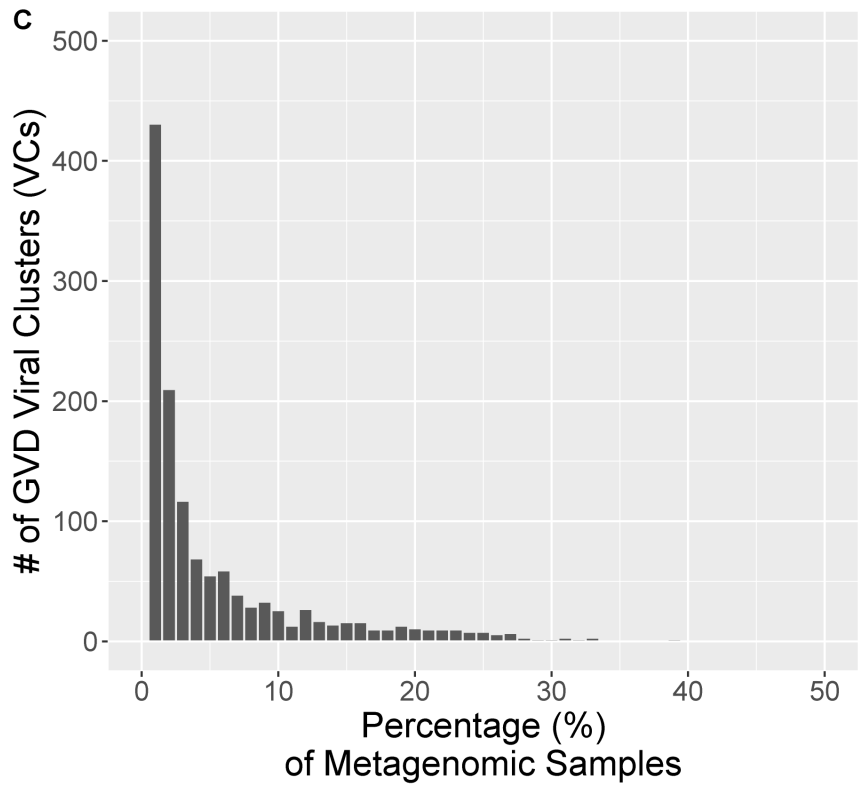
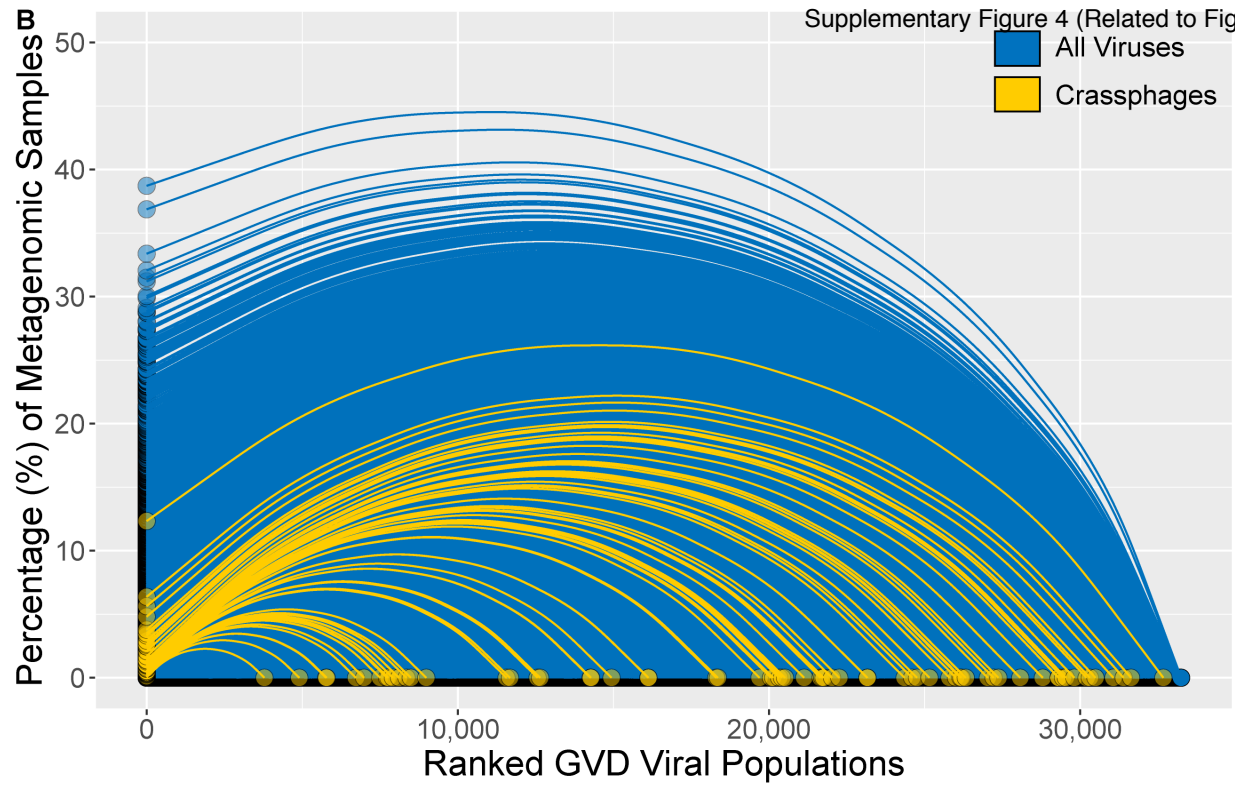
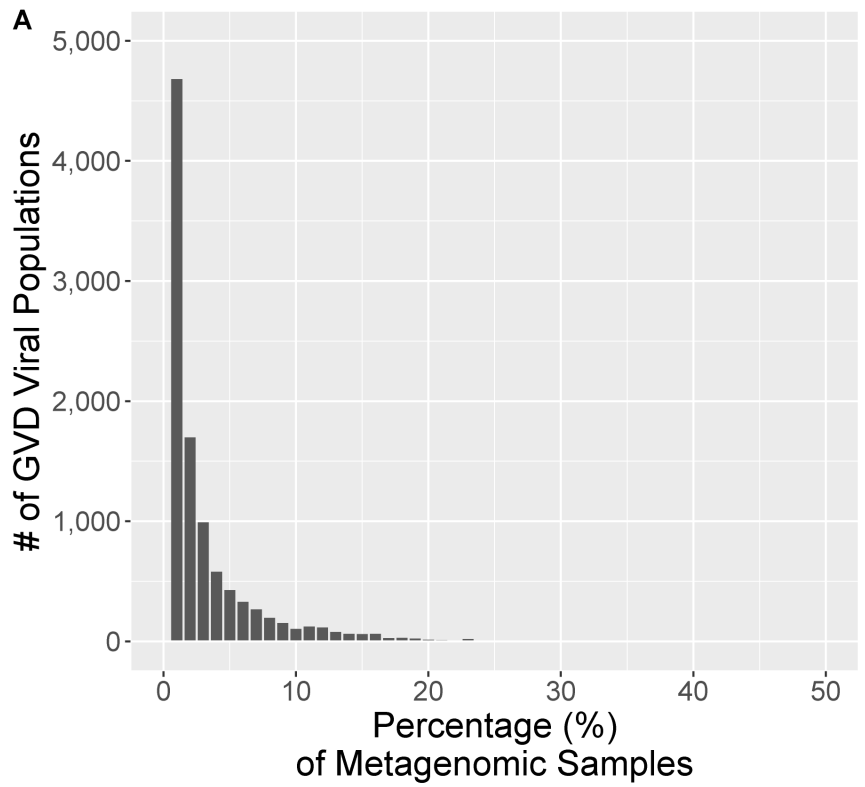
A**B****C**

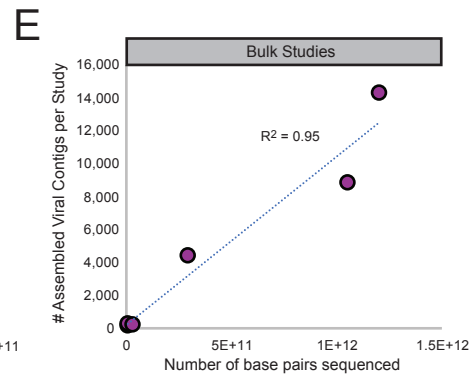
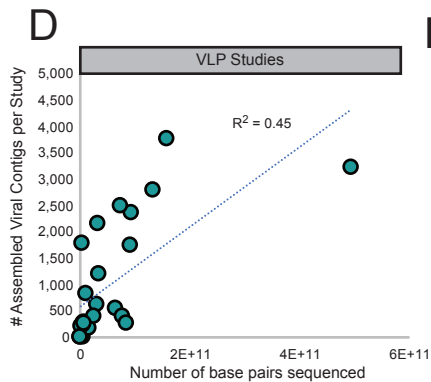
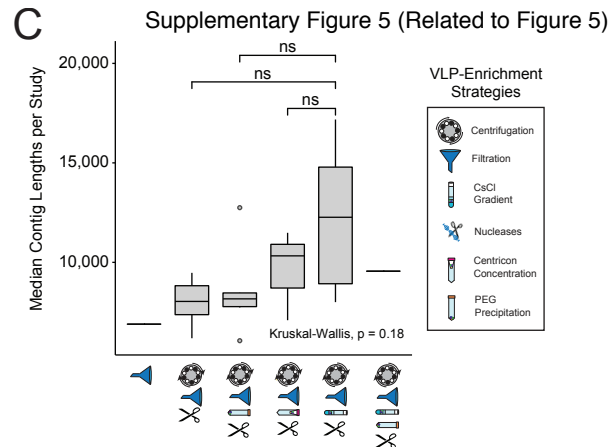
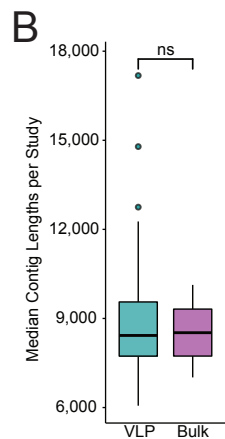
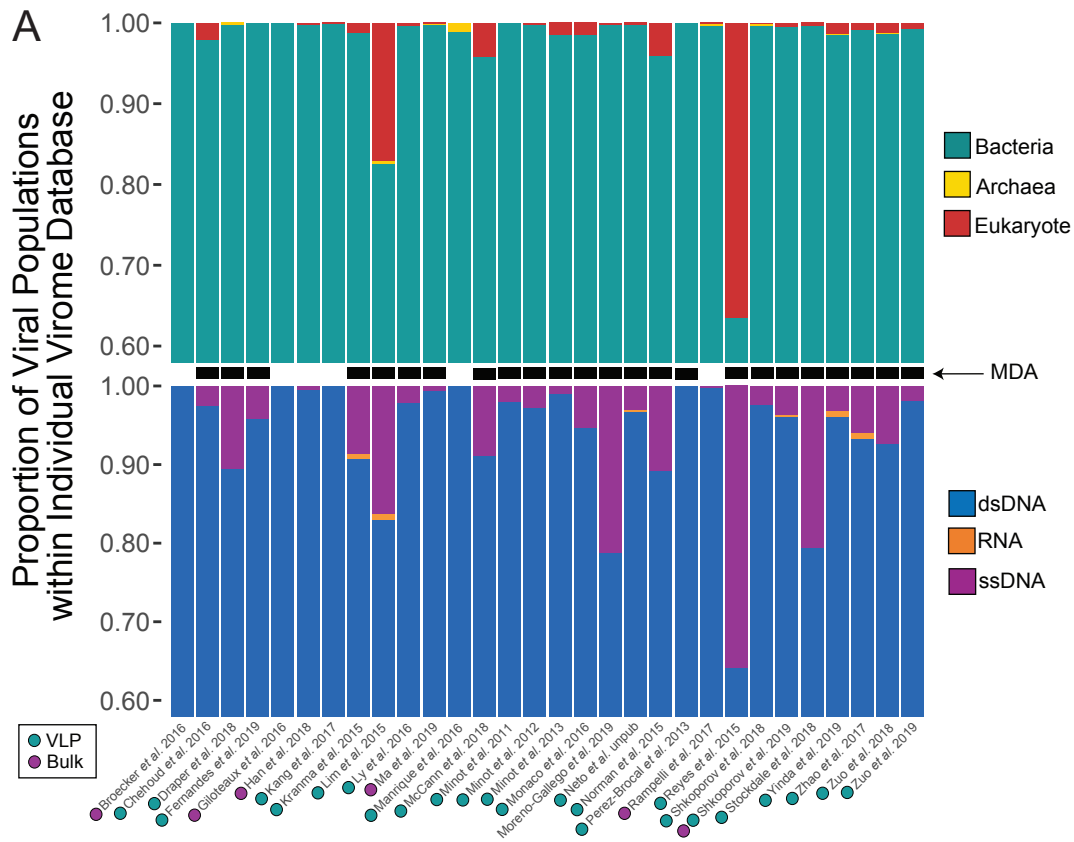
A

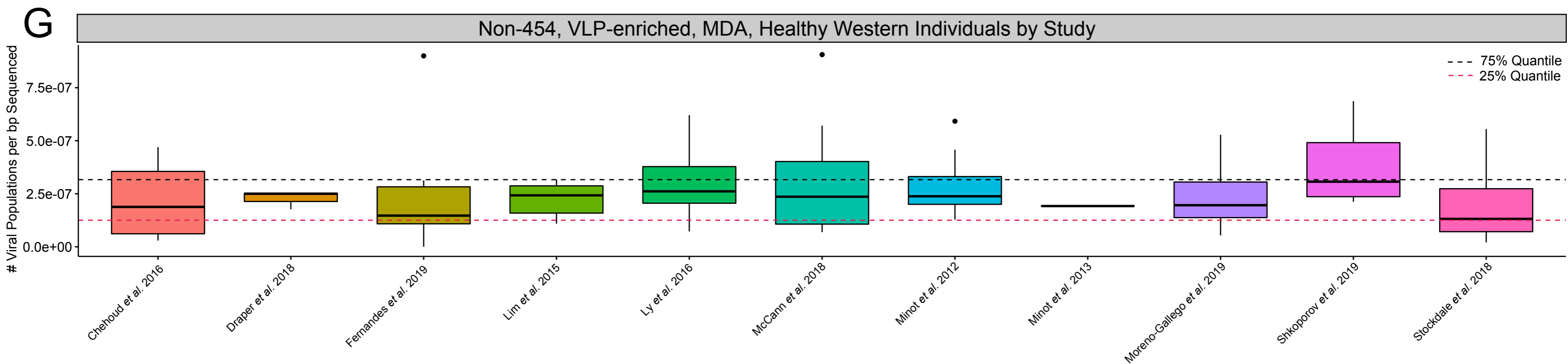
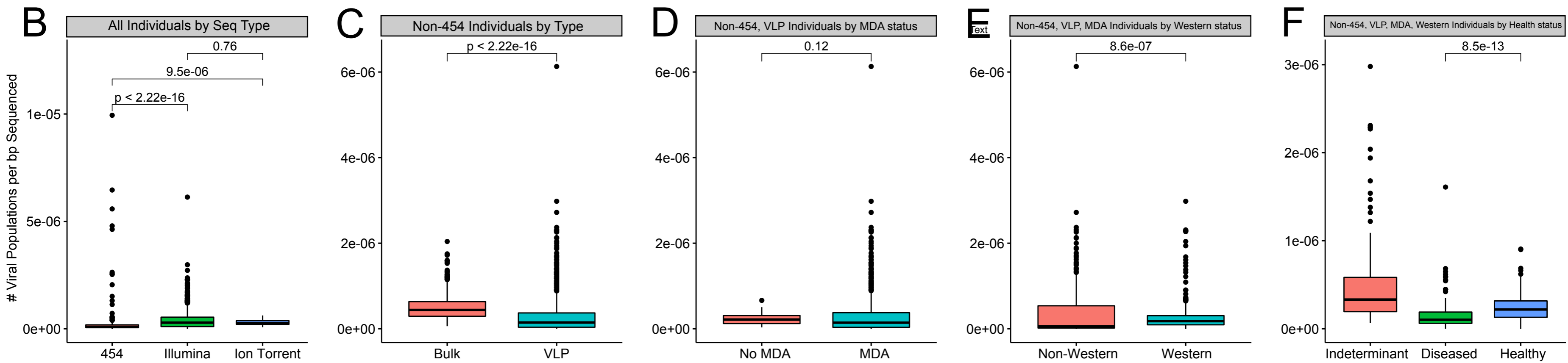
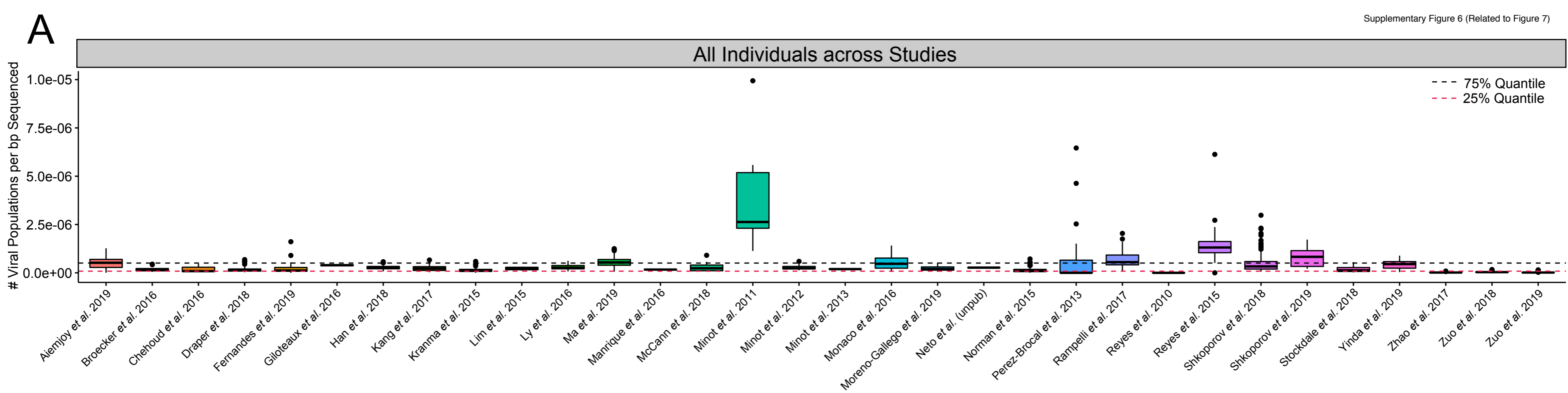


B









Supplementary Figure 1 (Related to Figure 2). Pipeline for the selection, processing and virus validation of human gut datasets. (a) Pipeline workflow showing that datasets were processed individually. Reads were filtered for quality, trimmed, and reads that mapped to Φ x174 and the human genome were removed. The remaining reads were assembled into scaffolds, filtered for lengths ≥ 1.5 kb, and run through tools that collectively utilize homology to viral reference databases, probabilistic models on viral genomic features, and viral k -mer signatures to identify putative viral genomes. The putative set of viral genomes were decontaminated using the BUSCO and VPF hmm models. Viral genomes were then deduplicated to get a total of 33,242 viral populations. (b) Density plot showing the number of BUSCO hits per total number of genes (BUSCO ratio) for all viruses in Viral Refseq v96. The **inset** dot plot showing the distribution of the BUSCO ratio for bacteriophages, with the highest value being 0.067. This value was the max value allowed for all GVD bacteriophages. (c) Histogram showing the number of GVD viral populations with different numbers of viral protein family (VFP) hits.

Supplementary Figure 2 (Related to Figure 2). (a) Barplot showing the number of assembled viral genomes versus the number of deduplicated viral populations per study. (b&c) Box plots showing median and quartiles of proportion of the total reads sequenced that mapped to the GVD viral contigs (b) per study and (c) and between VLP and bulk metagenomes.

Supplementary Figure 3 (Related to Figure 3). GVD improves read recruitment.

Boxplots showing median and quartiles of (a) the number of viral populations detected in a study not included in GVD, and (b) the number of reads mapped per study to the individual virome, Viral Refseq v96, JGI IMG/VR, or GVdb databases. All pairwise comparisons were performed using Mann-Whitney U-tests.

Supplementary Figure 4 (Related to Figure 4). There are no core viral populations or viral clusters across GVD samples. (a & c) Histogram showing the number of (a) viral populations and (c) viral clusters (VCs) present in different percentages of GVD samples. The vast majority of viral populations and VCs are found in $<10\%$ of the individuals. (b & d) Hive plot showing the percentage of GVD samples each (b) viral population and (d) VC is detected within. The dots on the x-axis represent each GVD viral population or VC in ascending order of the percentage of GVD samples that they are found within. The y-axis is the percentage of GVD samples that each viral population or VC is detected within. CrAssphage viral populations are highlighted in yellow in plot (b).

Supplementary Figure 5 (Related to Figure 5). (a) Barplots showing the proportion of those viruses that are bacteriophages, archaeal viruses, or eukaryotic viruses (top) and the proportion of those viruses that are dsDNA, ssDNA, or RNA viruses (bottom). The total number of assembled viral contigs and viral populations per study are available in **Supplementary Fig. 2a** and further details in **Supplementary Table 6**. Studies that used VLP-enriched or bulk metagenomes or used multiple displacement amplification (MDA) were indicated by marking in between the barplots. No viral contigs ≥ 1.5 kb were assembled from the Aiemjoy *et al.* 2019 and Reyes *et al.* 2010 studies. Boxplots showing median and quartiles of the median contig length per study (b) of VLP and bulk metagenomes and (c) of the different VLP-enrichment methodologies across the studies. Scatter plots with linear regressions lines showing the number of assembled viral contigs per base paired sequenced per study in (d) only VLP and (e) only bulk metagenome studies.

Supplementary Figure 6 (Related to Figure 7). Removing confounding variables for cross-study viral diversity analyses. (a) Box plots showing median and quartiles of the number of viral populations per base pair sequenced across the different studies. The dashed black and red lines represent the 75% and 25% quantiles, respectively, of the number of viral populations per base pair sequenced across all individuals in each study. (b-f) Boxplots showing median and quartiles of the number of viral populations per base pair sequenced comparing sequencing platform, enrichment type (VLP vs. bulk), MDA status, geographic origin (Western vs. non-Western), and health status all sequentially and additively tested. All pairwise comparisons were performed using Mann-Whitney U-tests. (g) Boxplots showing median and quartiles of the number of viral populations per base pair sequenced per study for the remaining non-454, VLP-enriched, MDA, healthy Western individuals. The dashed black and red lines represent the 75% and 25% quantiles, respectively, of the number of viral populations per base pair sequenced across all remaining individuals in each study.

Supplementary Figure 7 (Related to Figure 7). Statistical differences across viral groups assessed across the life stages. (a & b) Box plots showing median and quartiles of the number of viral populations per base pair sequenced for all GVD viruses, bacteriophages, eukaryotic viruses, and different viral families including crAssphage across the life stages across healthy Western individuals for (a) all viral populations and (b) with putative contaminant viral populations removed. All pairwise comparisons were performed using Mann-Whitney U-tests.