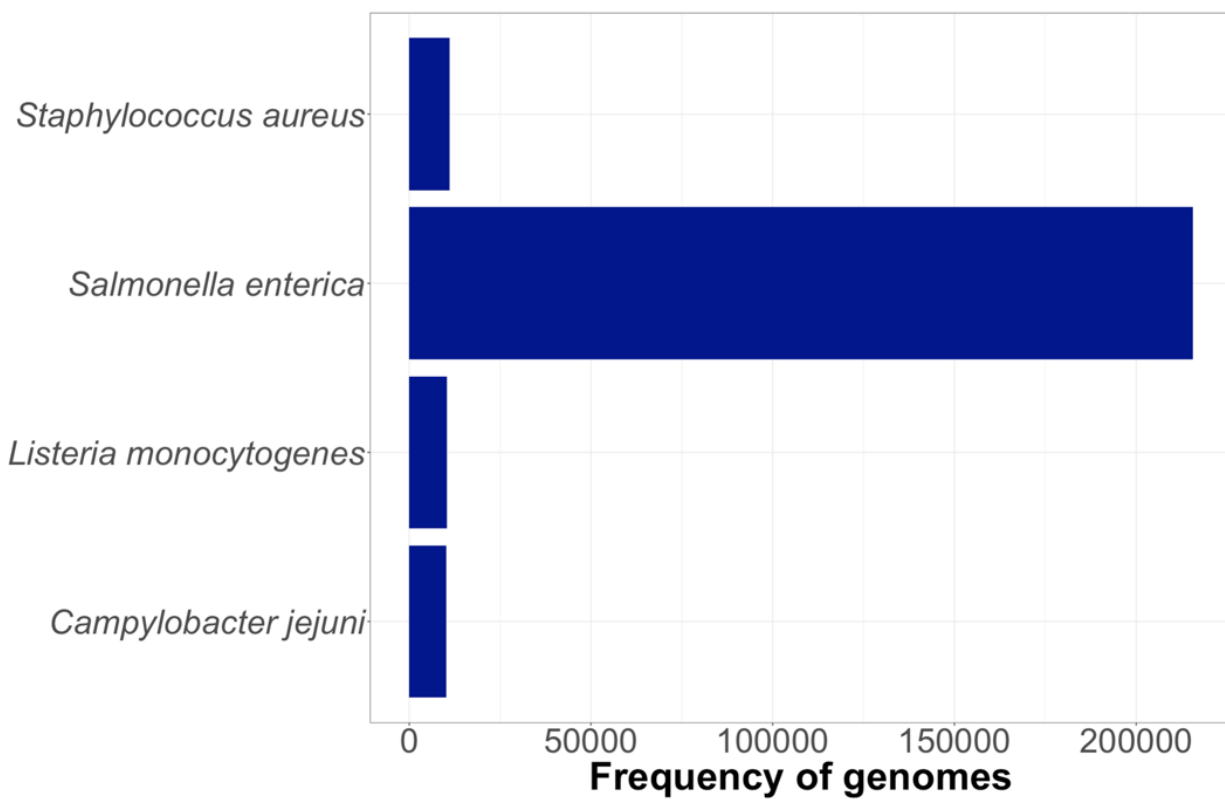
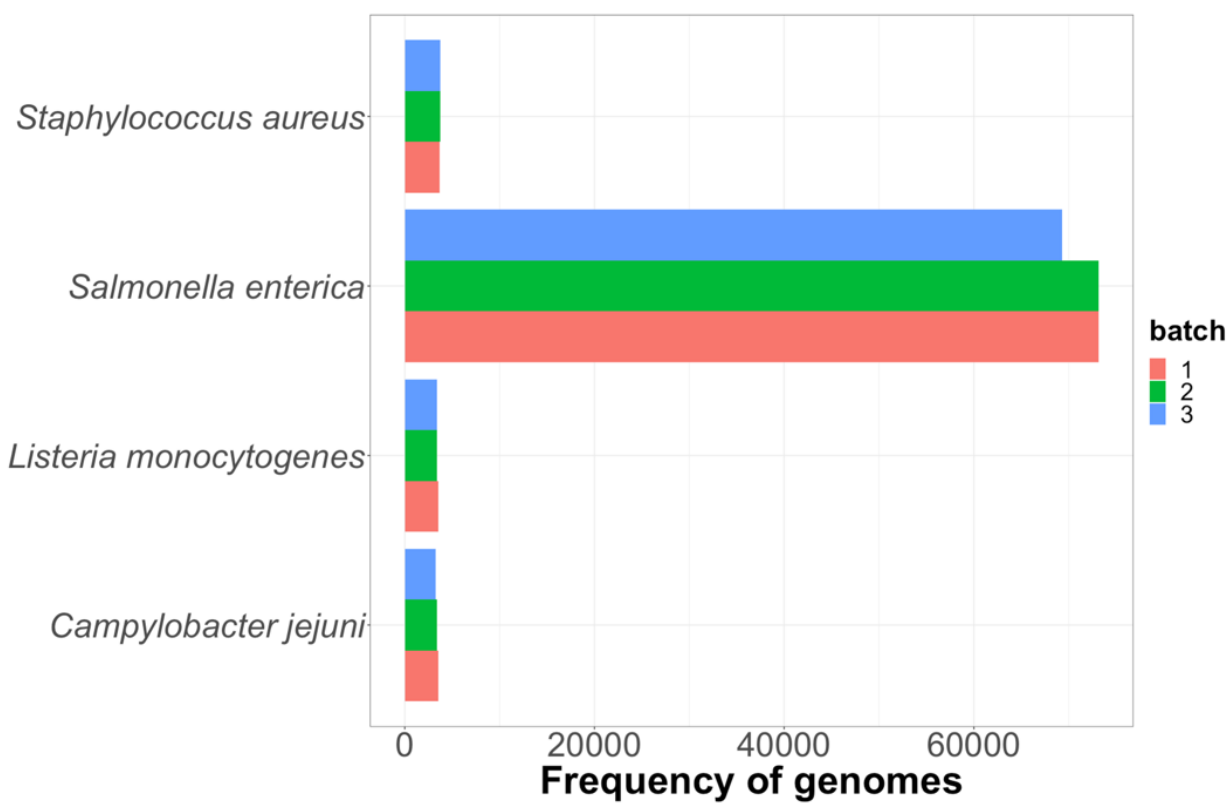
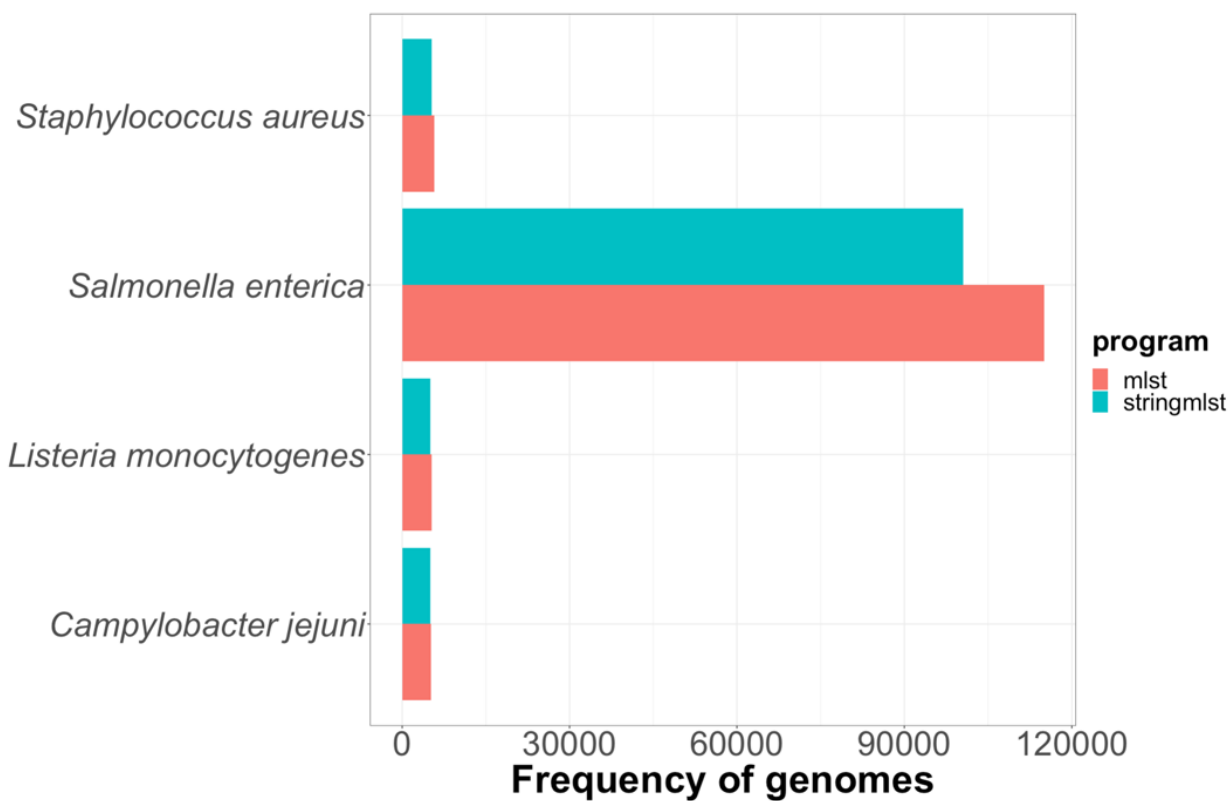
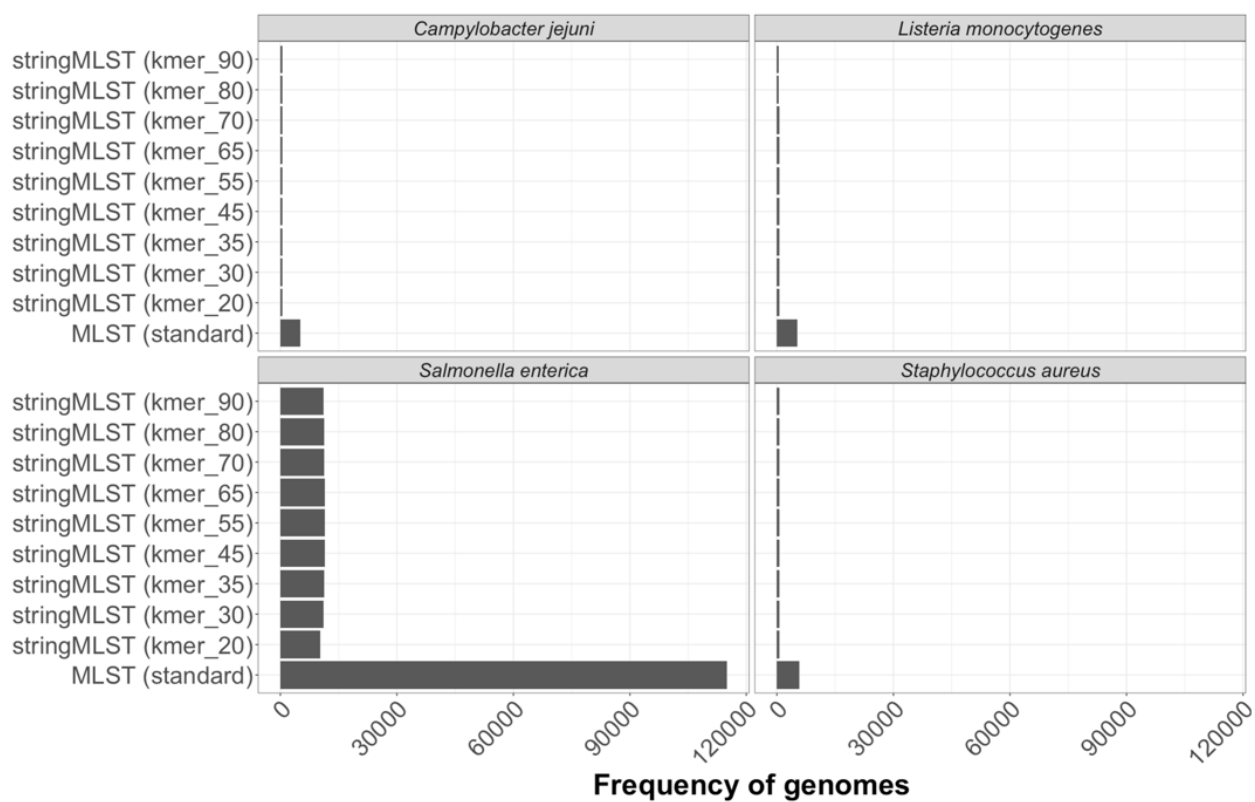


**A.****B.**

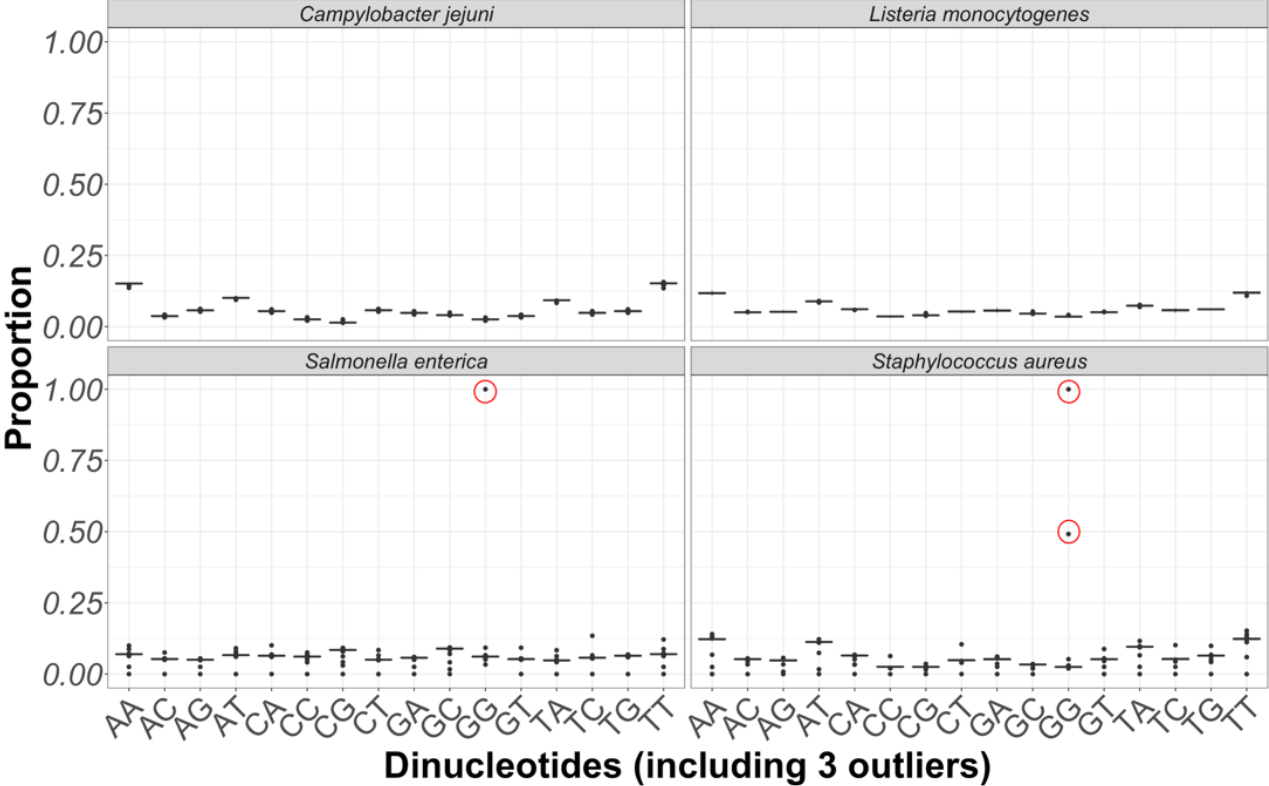
**C.**



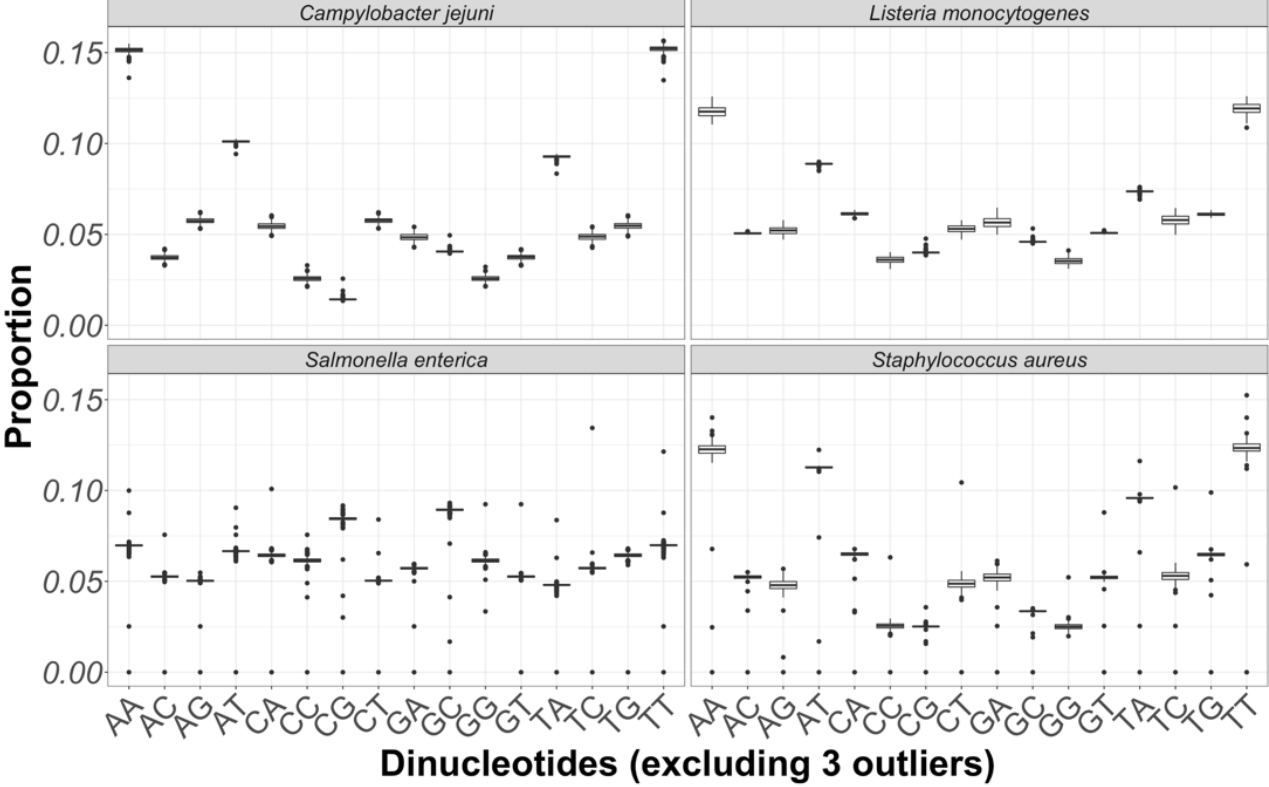
**D.**

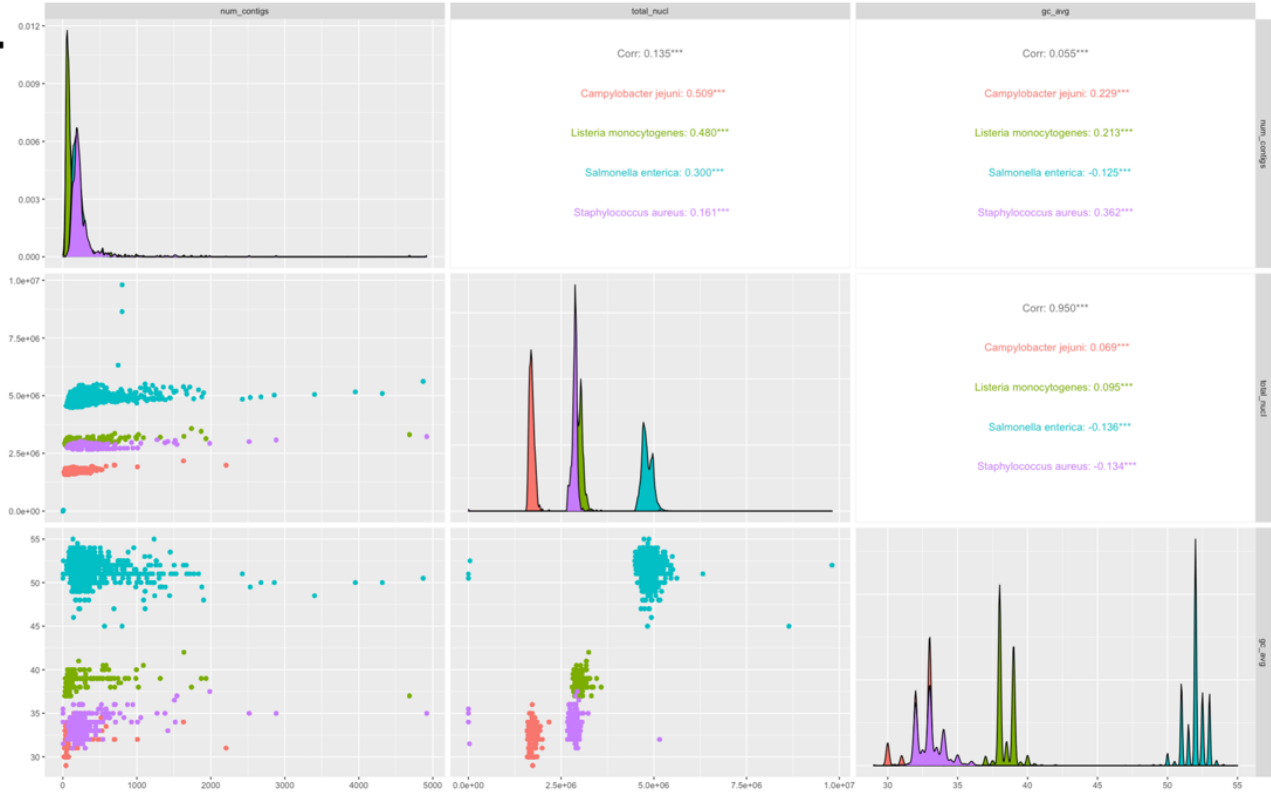


E.



F.



**G.**

**Figure S2.** Summary statistics of the frequency of genomes, including the distribution of dinucleotides and bivariate associations between genome-intrinsic variables across all four bacterial species.

Frequency-based distribution of randomly selected genomes across bacterial species (A), including a stratification by batch (B), program (C), and further differentiation by k-mer length used by stringMLST (D). Proportion of all sixteen pairs of dinucleotides present in a bacterial genome, across species, with (E) or without outliers (red-circled data points) (F). (G) Bivariate association between genome-intrinsic variables across species with statistical significance measured by the Pearson correlation coefficient (Corr). Genome-intrinsic variables used were the total number of contigs (num\_contigs), the total number of nucleotides per genome (assembly), and the GC% content per genome (gc\_avg). (G) Asterisks refer to the degree of significance for the correlation coefficient, with  $p$ -value thresholds being:  $*p < 0.05$ ,  $**p \leq 0.01$ ,  $***p \leq 0.001$ ,  $****p \leq 0.0001$  and NS = not significant at  $p \geq 0.05$ . Across all figures A-G, data for *S. enterica* Subspecies *enterica* (*S. enterica*) contained an even proportion of genomes across twenty serovars.