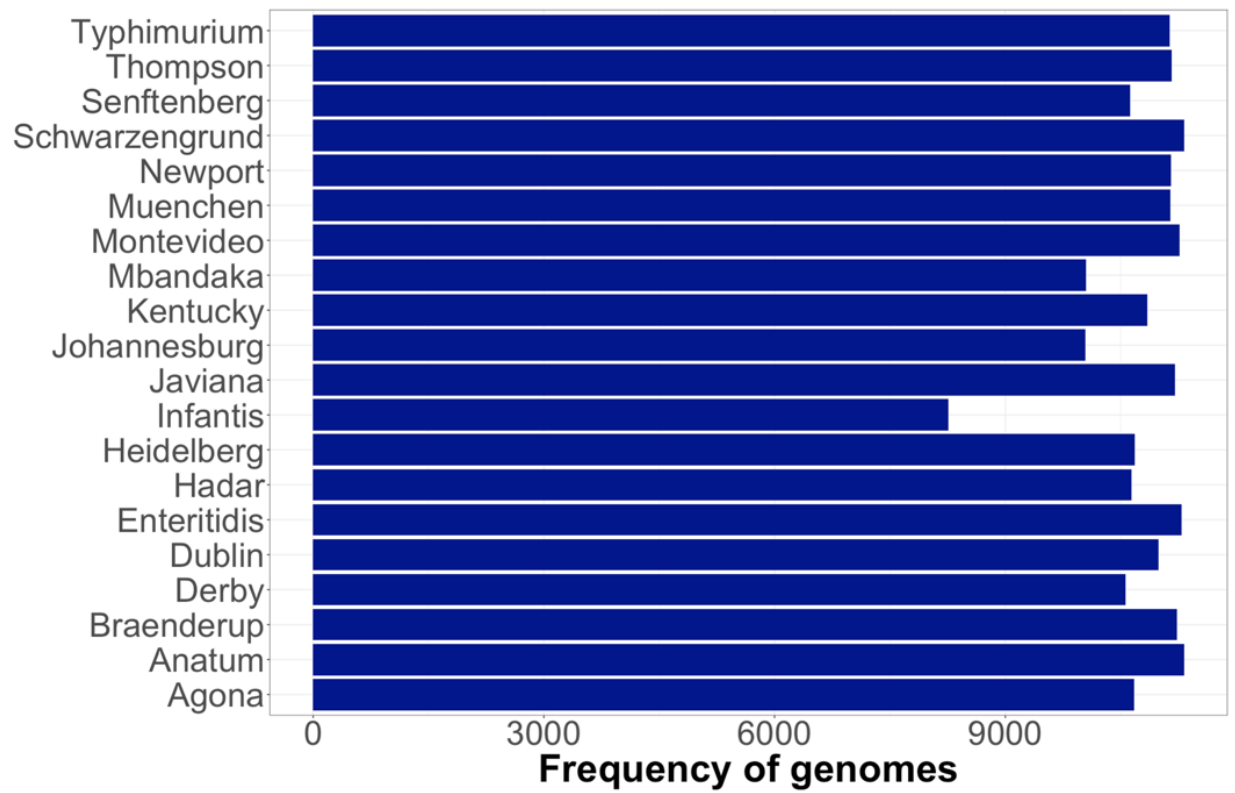
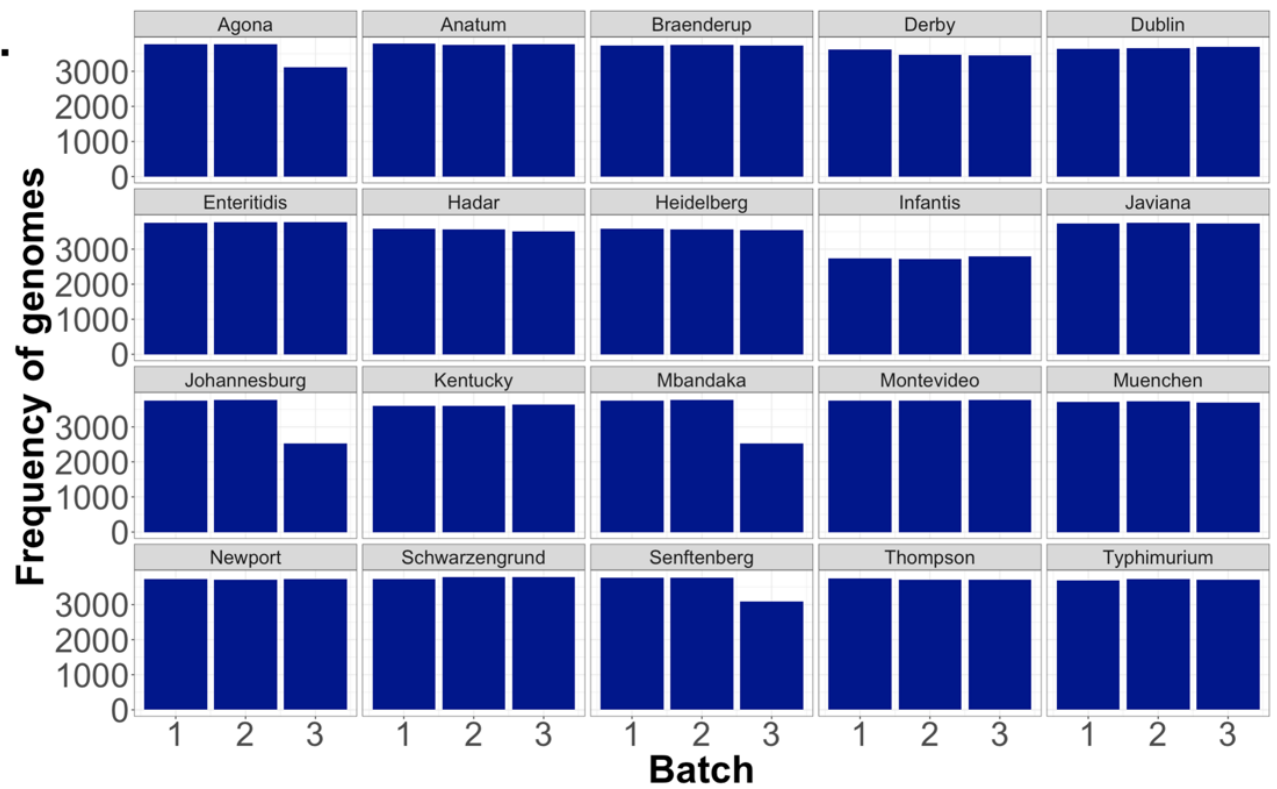
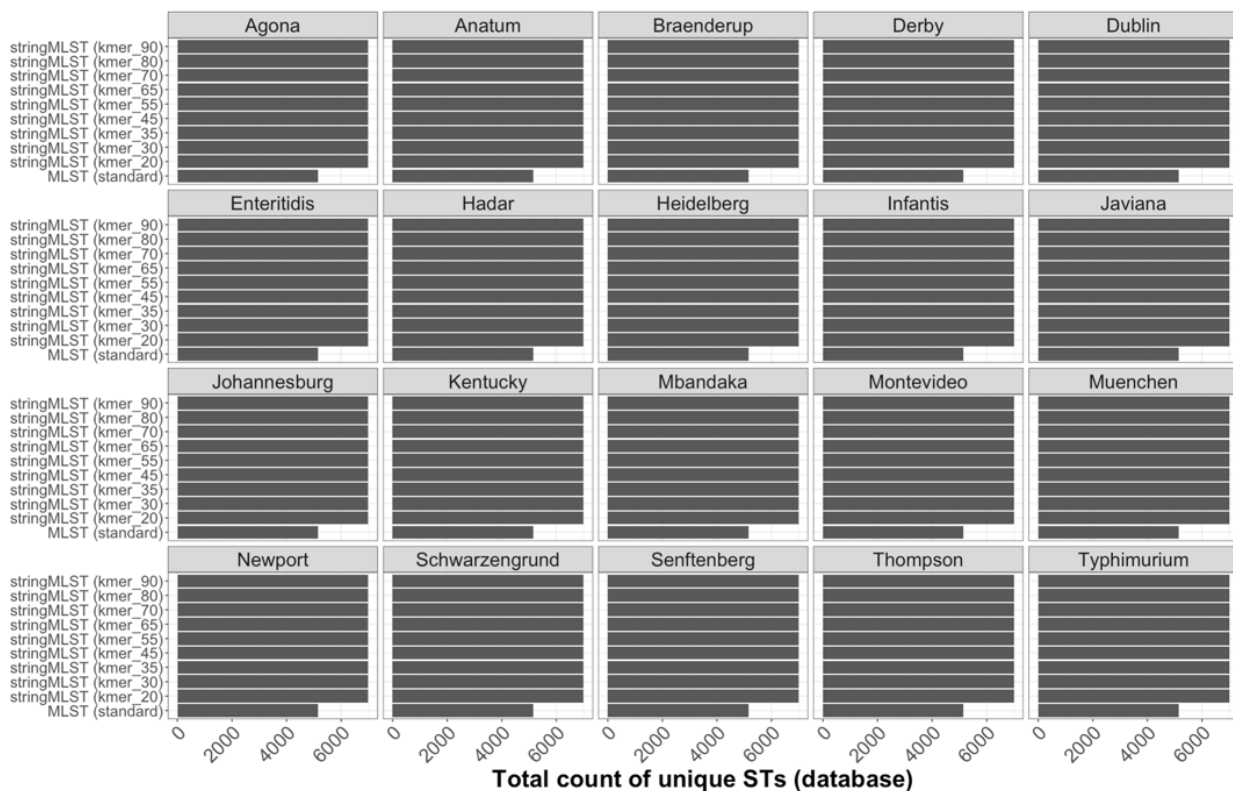
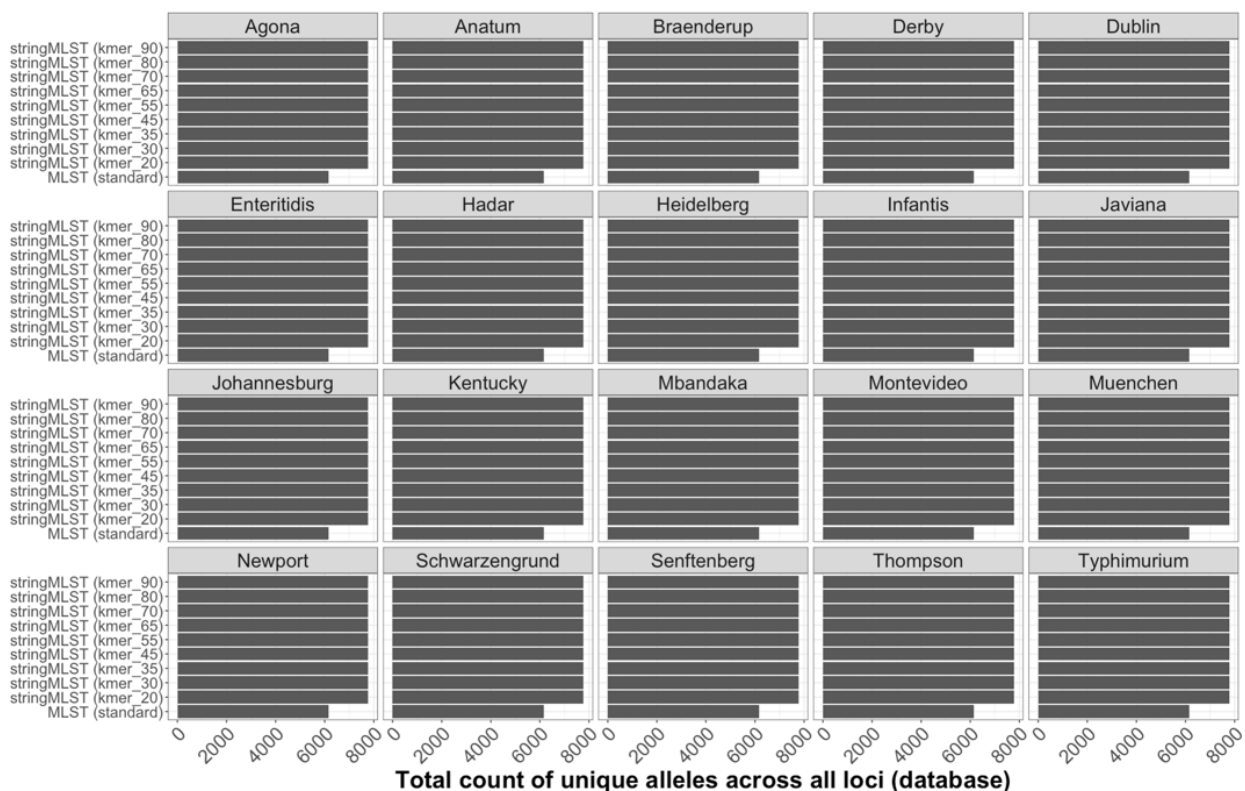


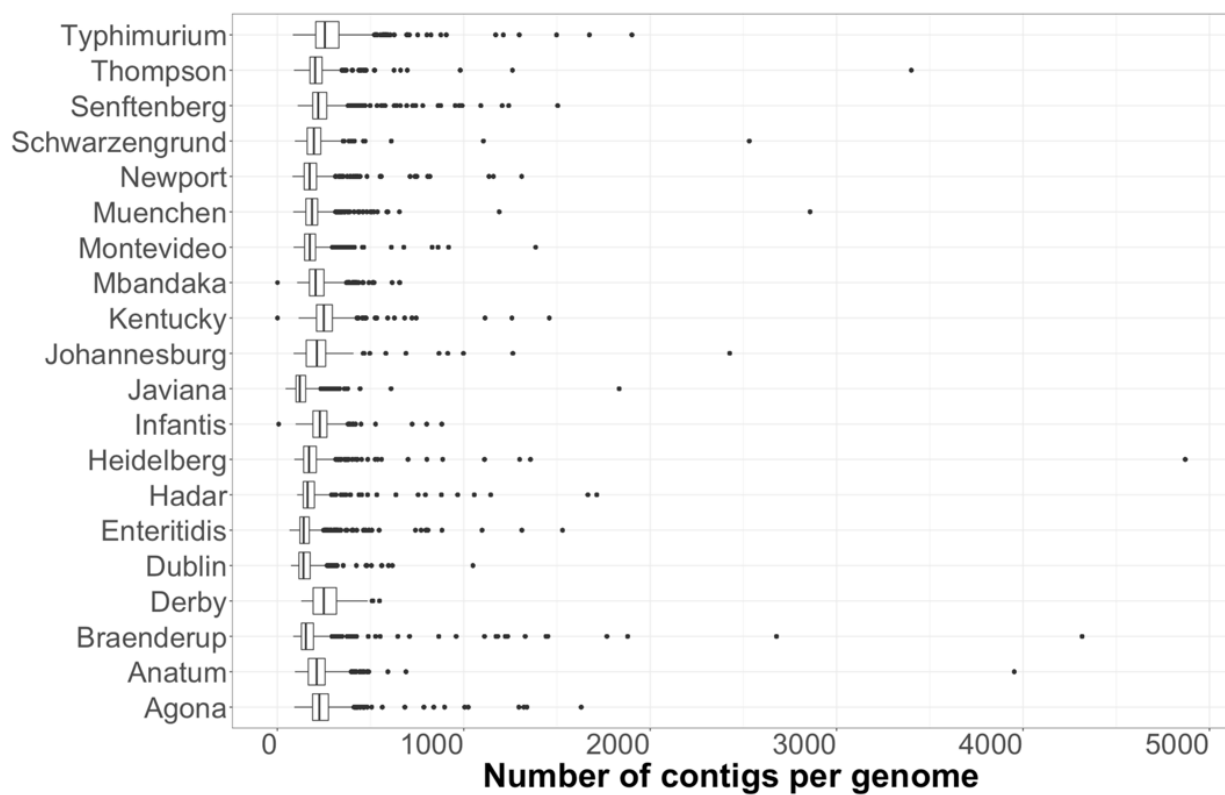
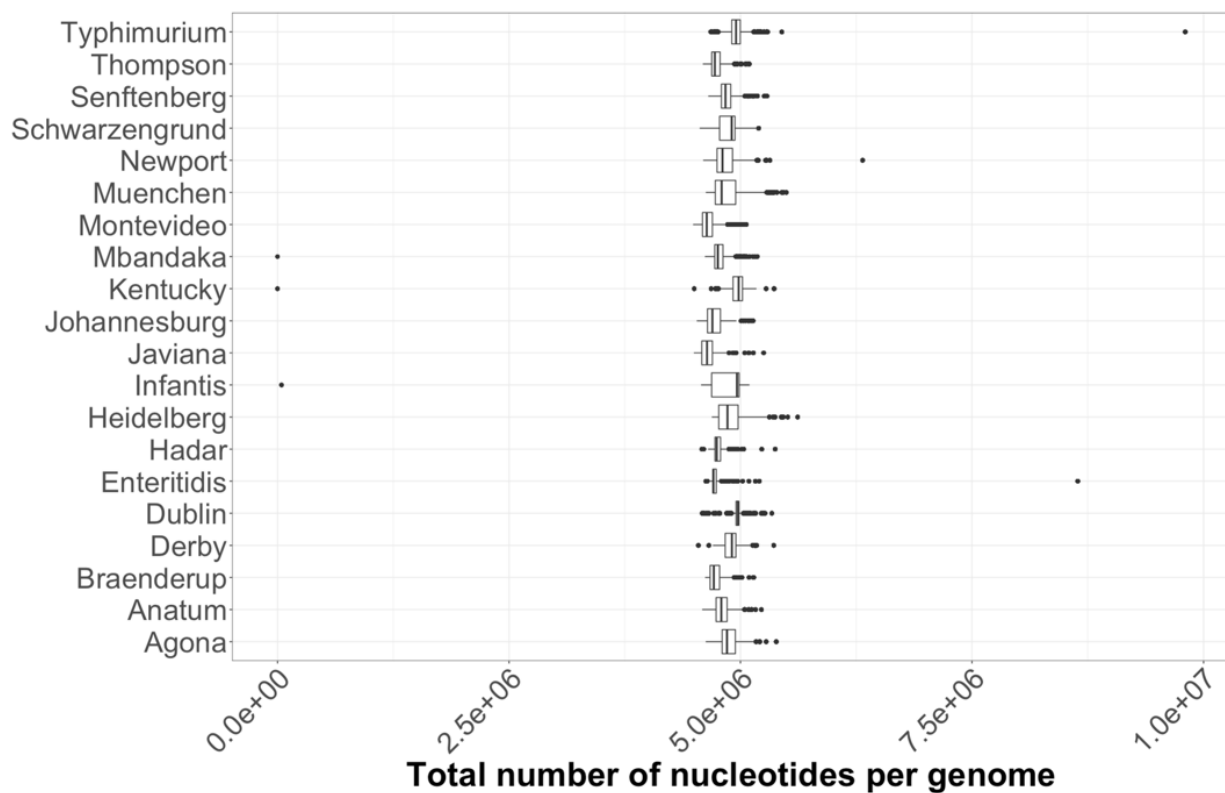
A.**B.**

E.

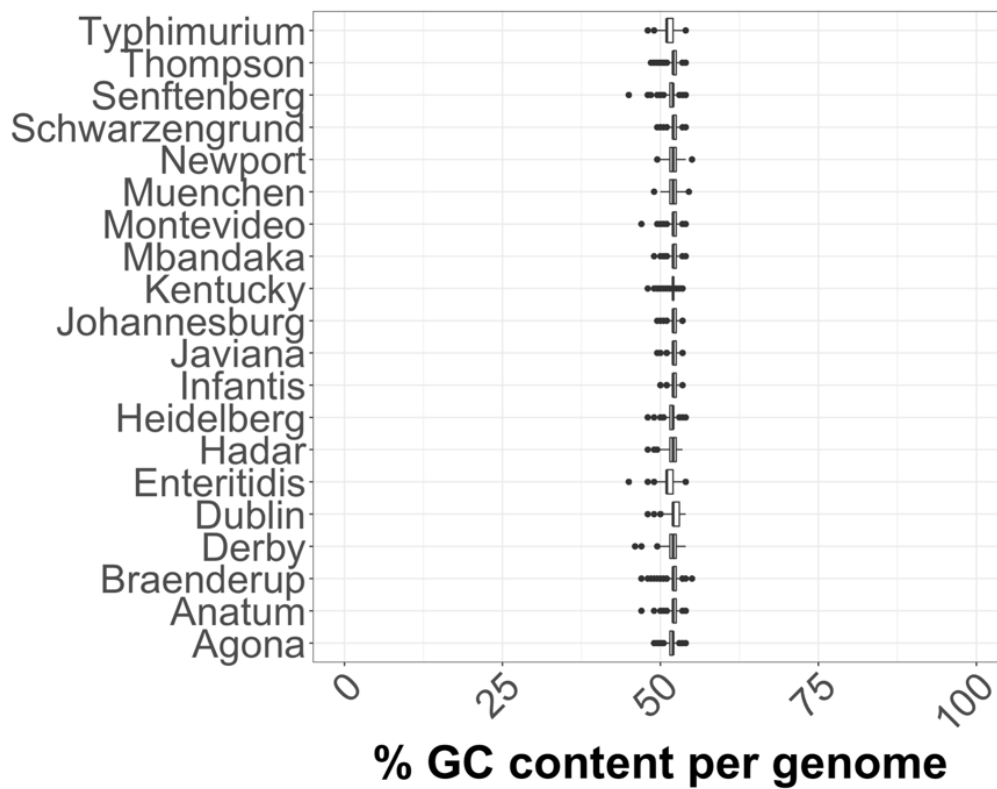


F.

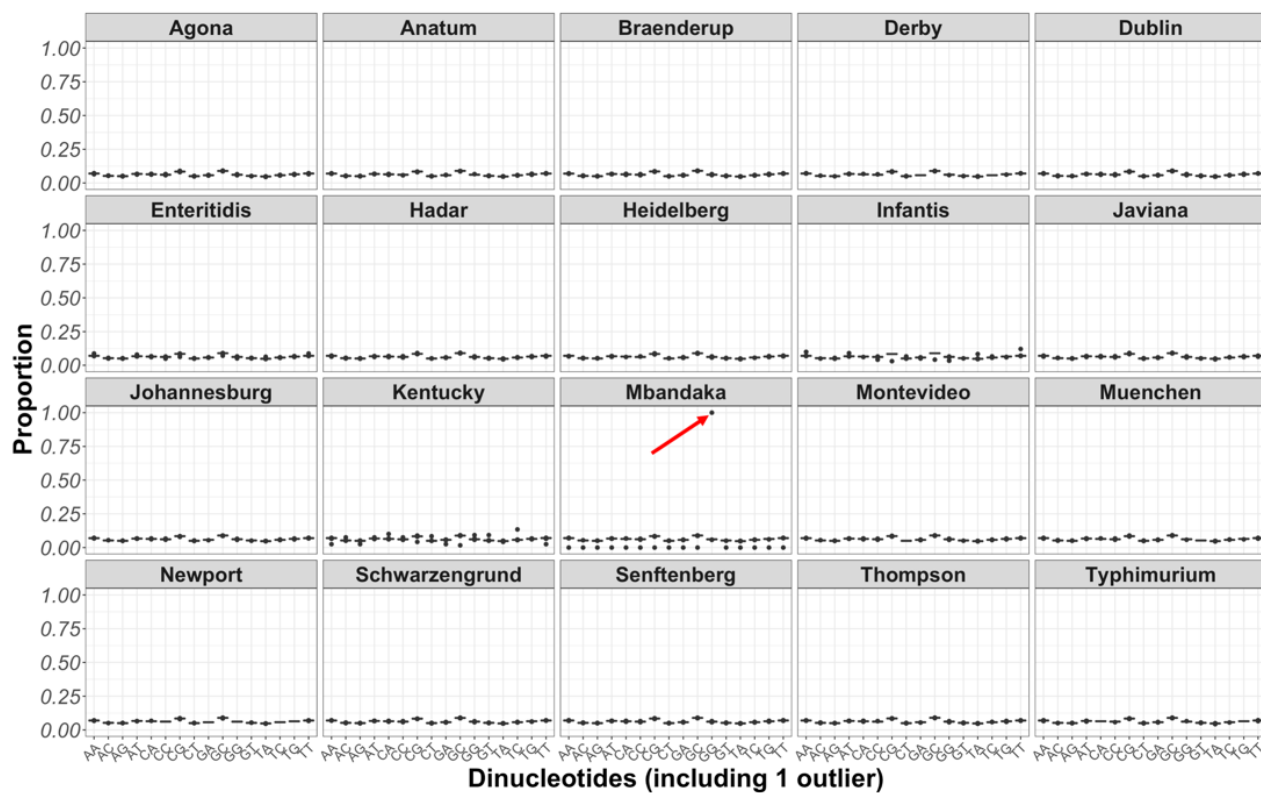


G.**H.**

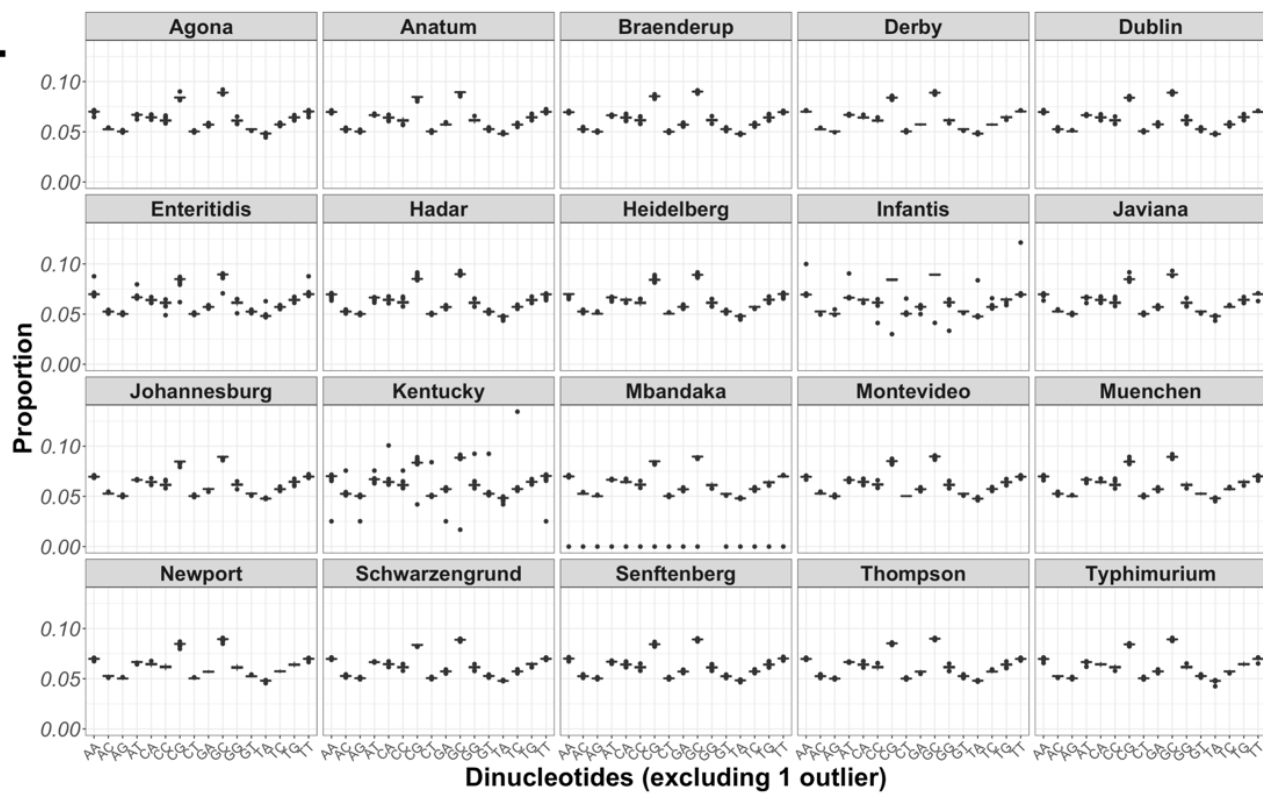
I.



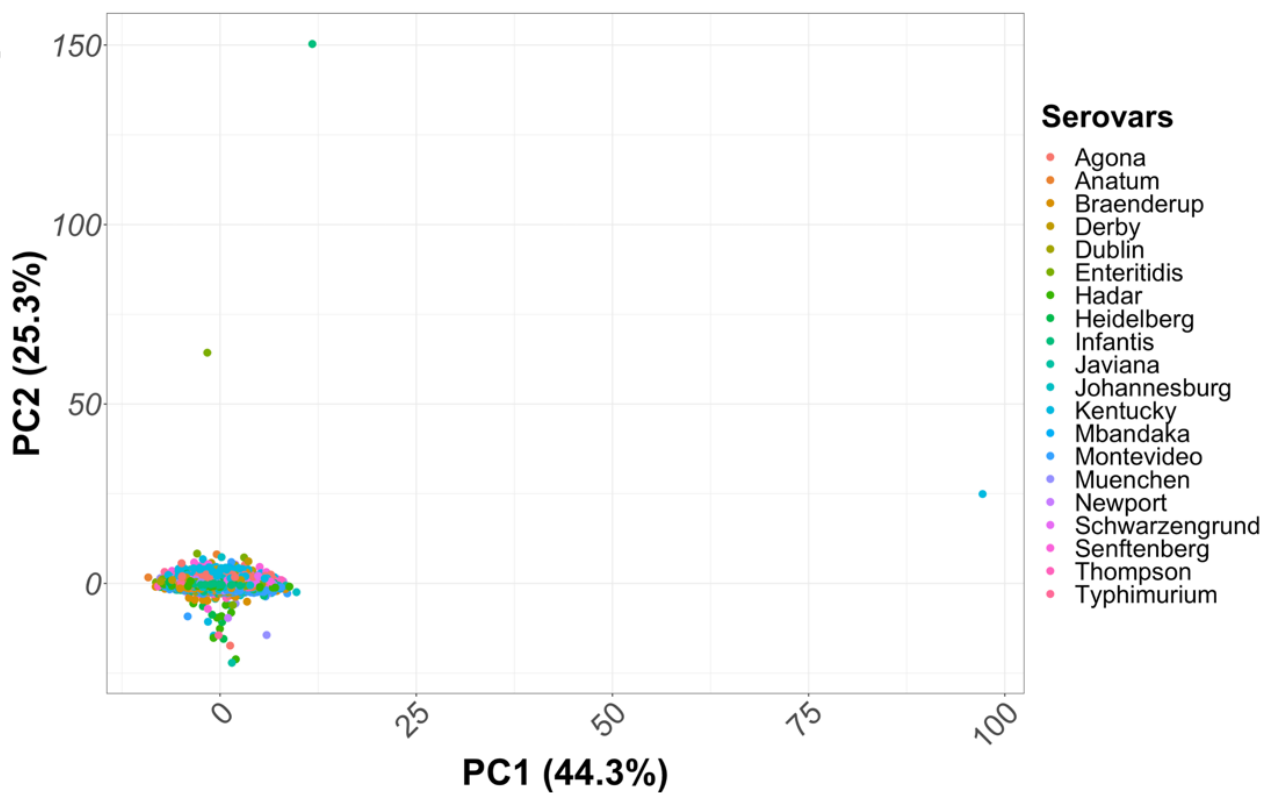
J.



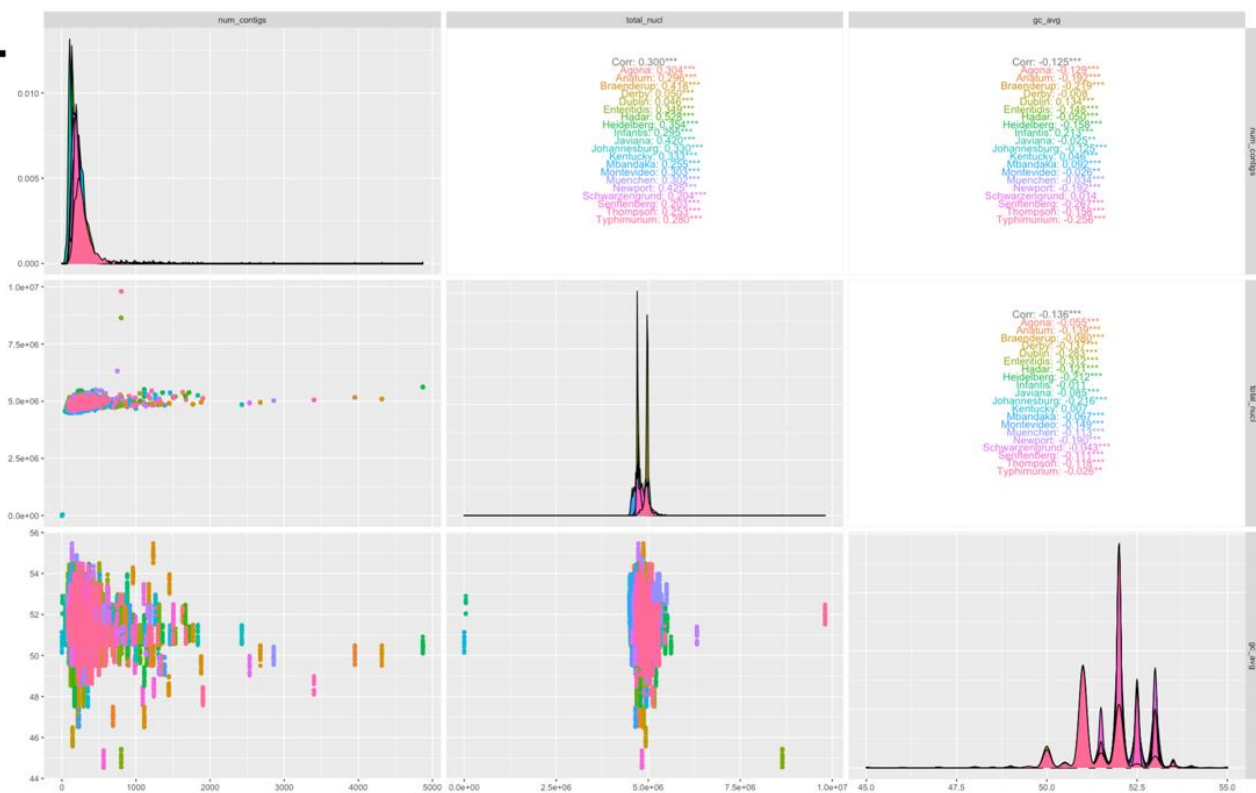
K.



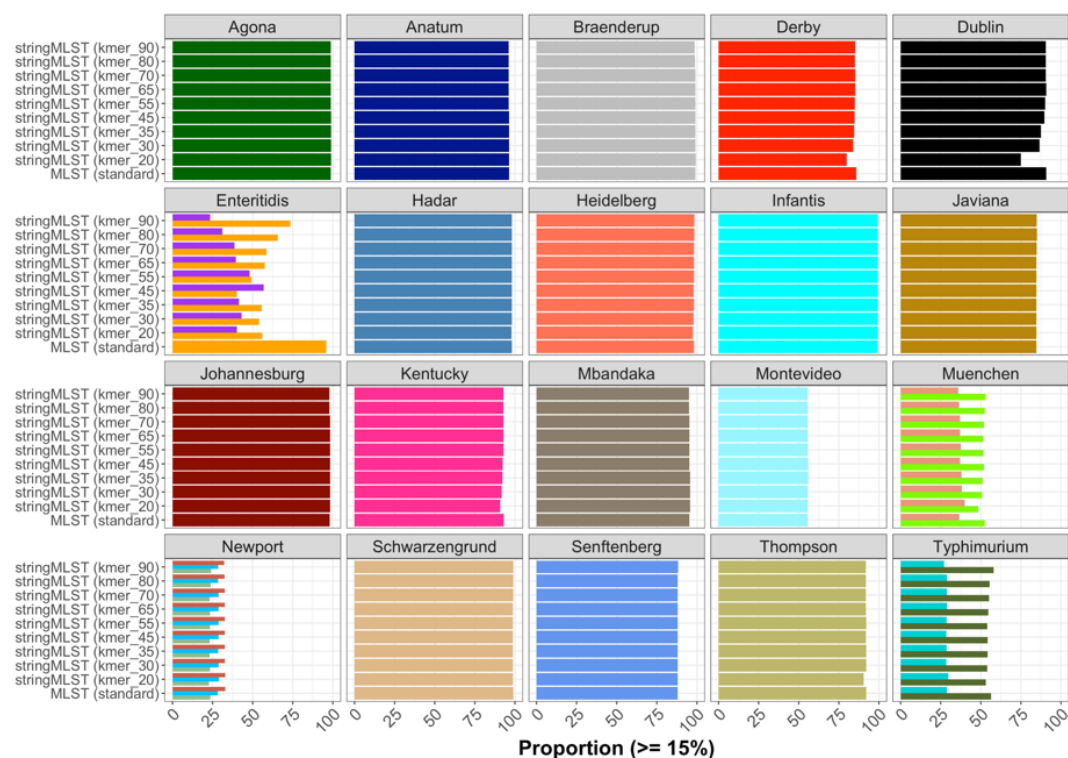
L.



M.

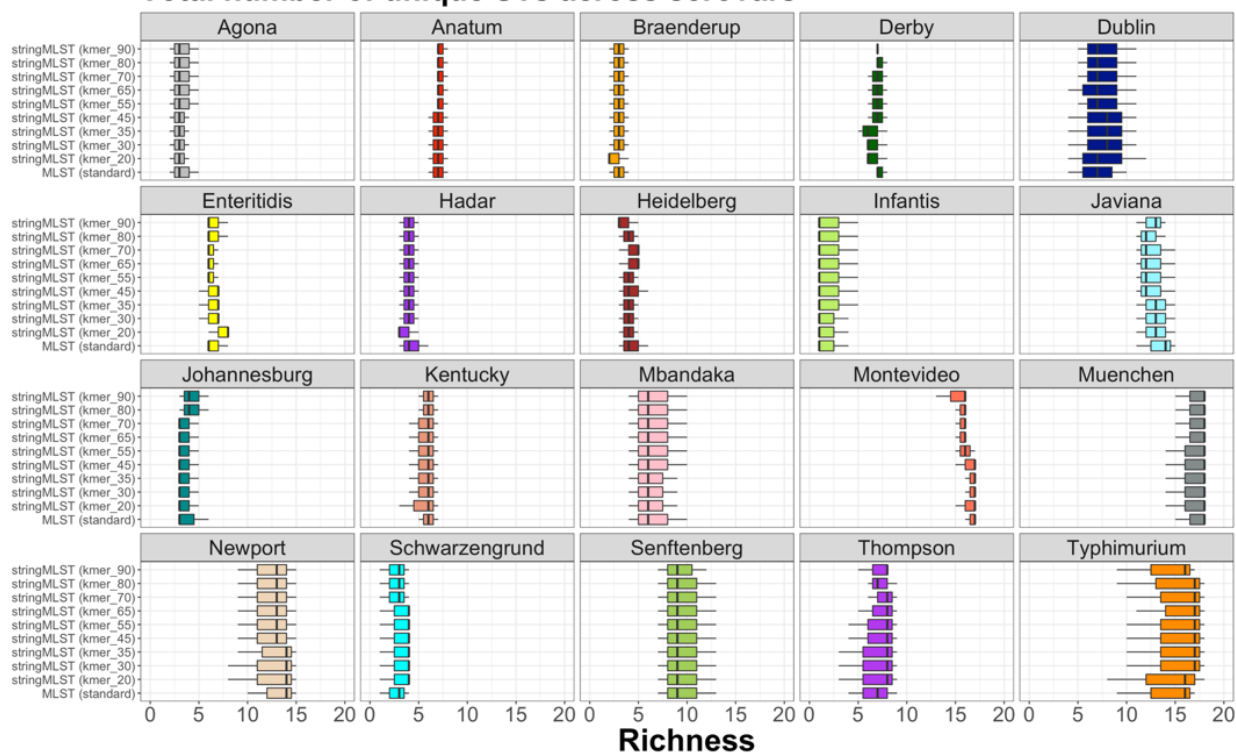


N.



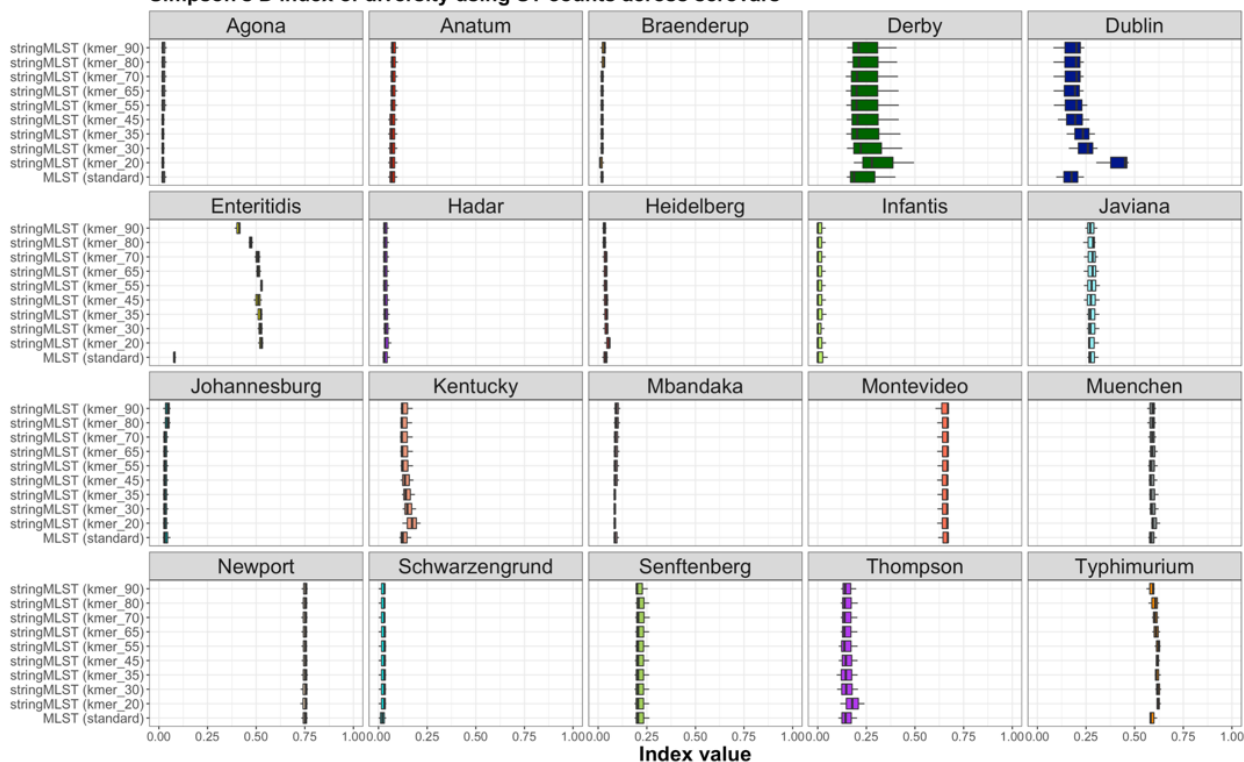
O.

Total number of unique STs across serovars



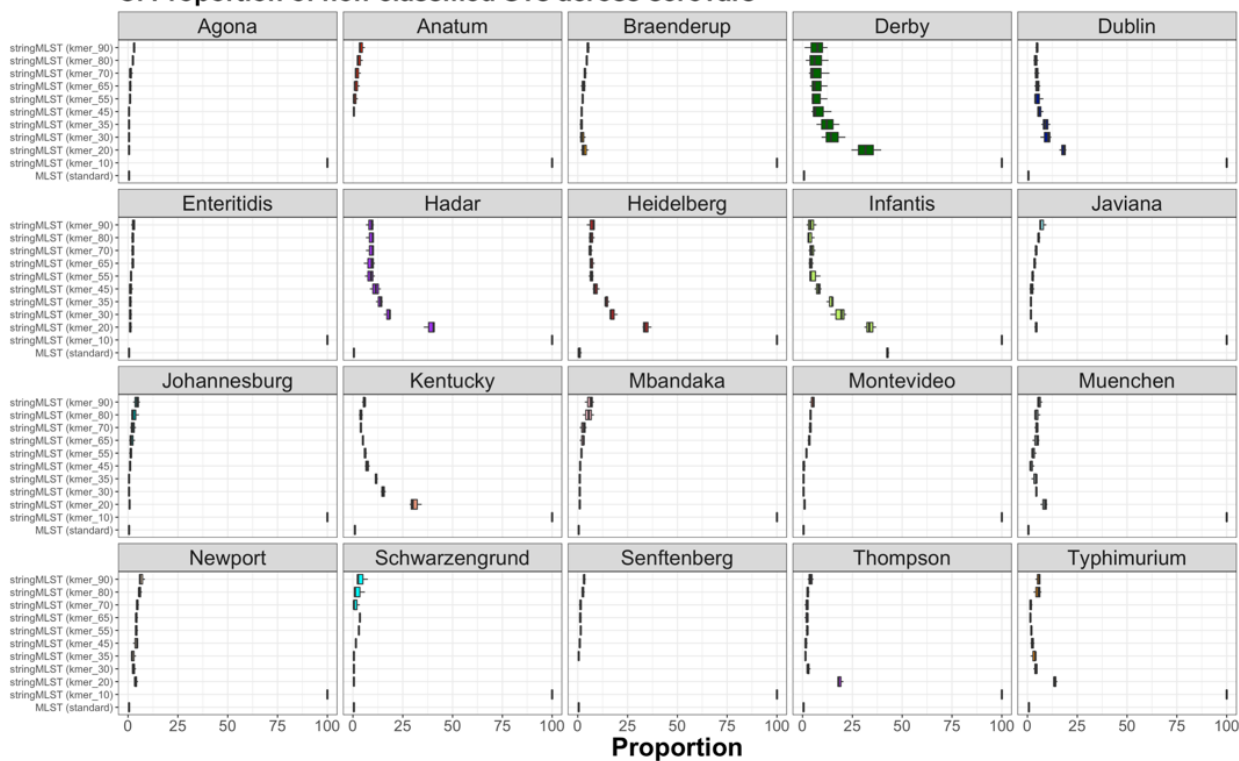
P.

Simpson's D index of diversity using ST counts across serovars



Q.

C. Proportion of non-classified STs across serovars



R.

D. Standard deviation for the proportion of non-classified STs across serovars

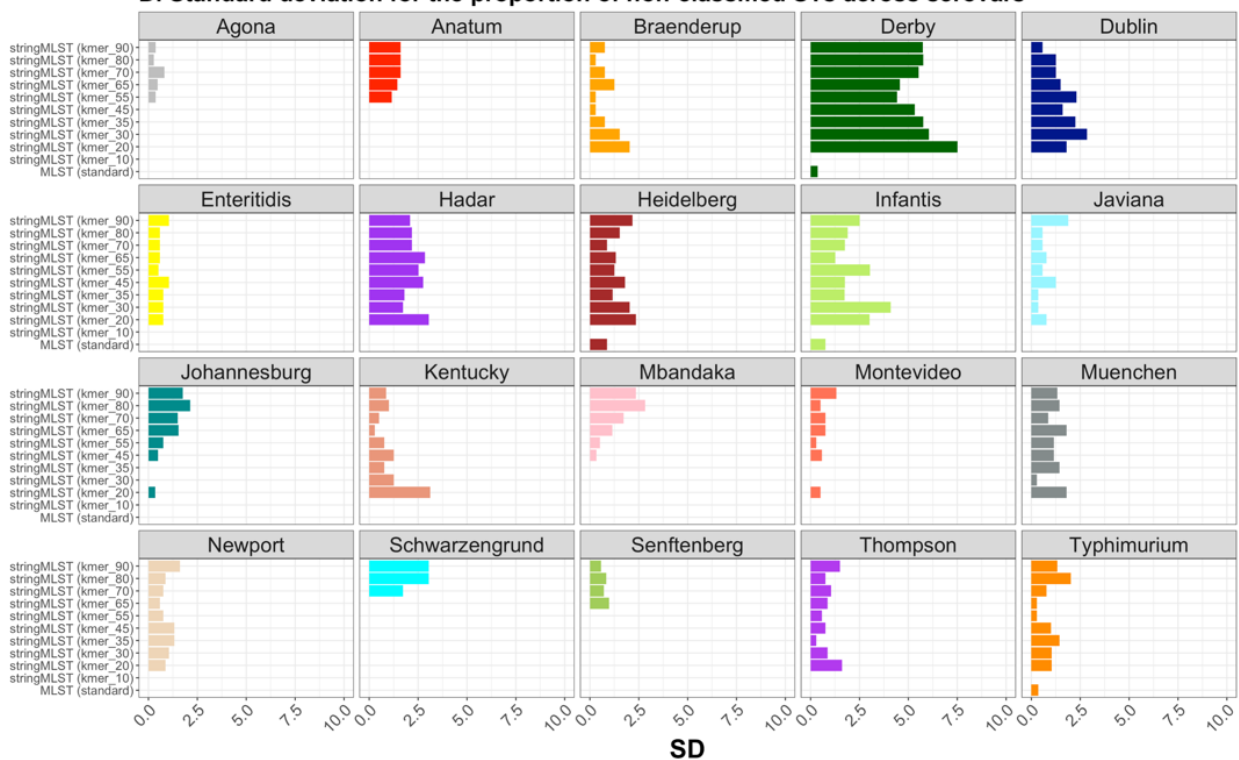


Figure S3. Summary statistics of the frequency of genomes, including the distribution of dinucleotides and bivariate associations between genome-intrinsic variables across all twenty *S. enterica* serovars (serovars).

Frequency-based distribution of randomly selected genomes across serovars (A), including a stratification by batch (B), program (C), and further discrimination by k-mer length used by stringMLST (D). Total count of unique STs (E), and alleles across all loci (F) by serovar and program. (G) Total number of contigs per genome across all serovars. (H) Total number of nucleotides per genome across all serovars. (I) Percent of GC% content per genome across all serovars. Proportion of all sixteen pairs of dinucleotides present, across serovars, with (J) or without outliers (red-circled data points) (K). (L) Principal component analysis plot depicting two PCs, along with variance explained, for the distribution of serovars using the dinucleotide data as input. (M) 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). (M) 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 $p \geq 0.05$. (N) Proportion of dominant STs (proportion $\geq 15\%$) across programs and serovars. Distribution of ST richness (O), Simpson's D index of diversity ($1 - D$) (P), proportion of non-classified STs (Q), and the standard deviation of the proportion of non-classified STs (R) across programs and serovars.