

ANNEX

Biogeography and ecological distribution of environmental microbial protein families at a global scale

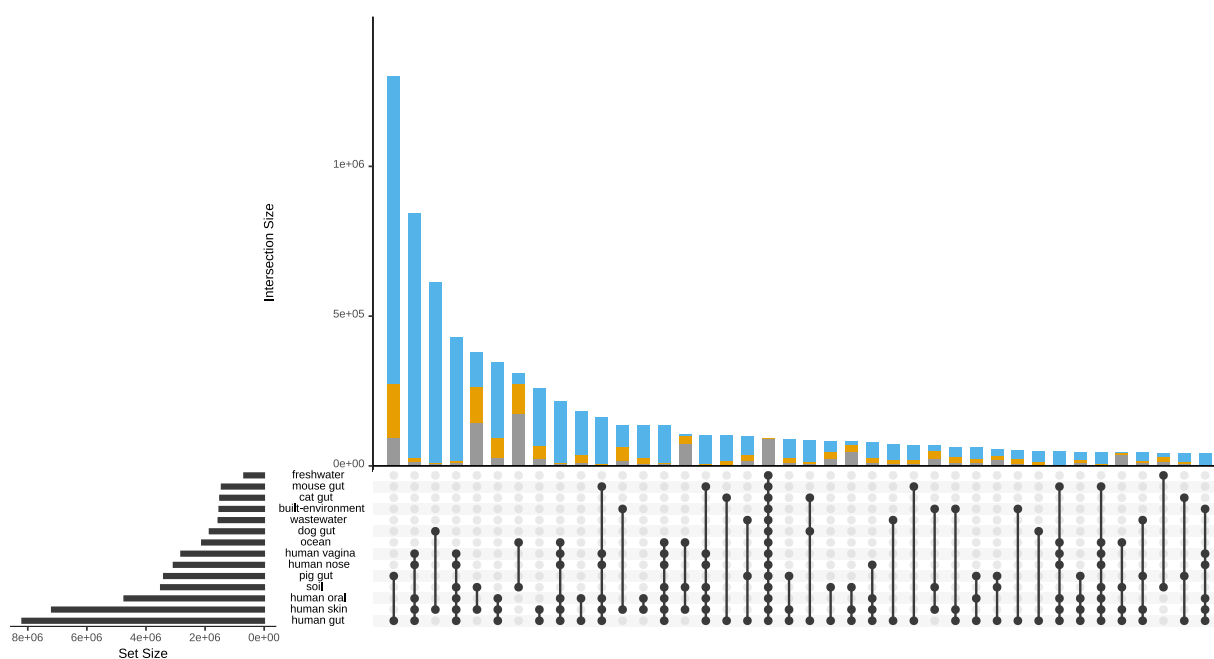


Figure 8.1: Distribution of known function (gray), unknown function (yellow), and novel families exclusive to uncultivated taxa (blue) across the GMGC, Barcharts indicate the number of families shared by the habitats indicated below. Horizontal bars represent the number of gene families detected in each of the 14 GMGC habitats. Only families detected on more than one habitat are shown

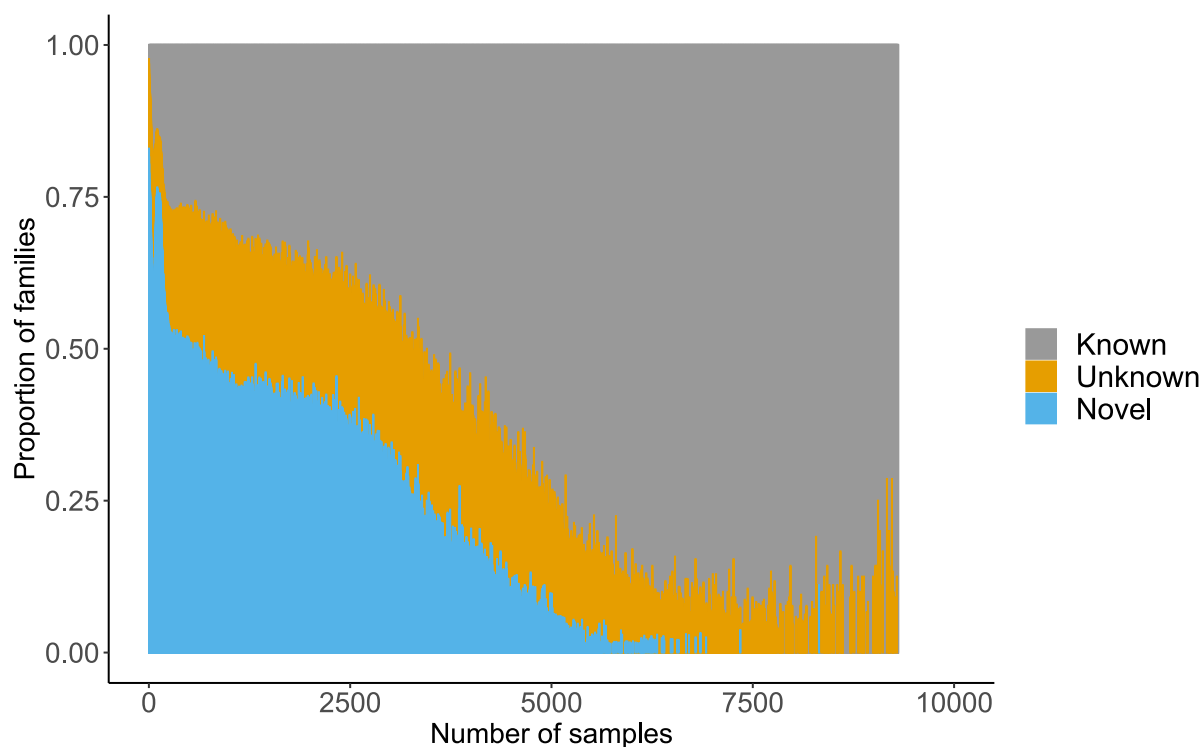


Figure 8.2: Proportion of known function, unknown function and novel gene families detected in the number of samples indicated

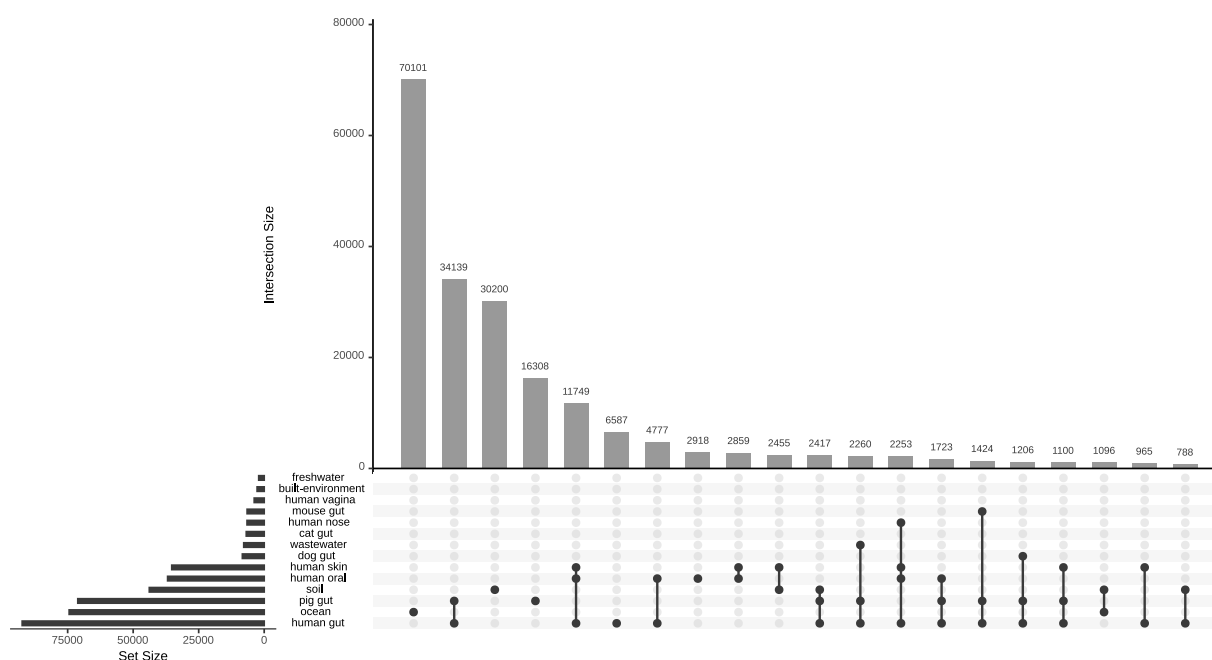


Figure 8.3: Distribution of filtered novel gene families in the GMGC habitats. Barcharts indicate the number of families shared by the habitats indicated below. Horizontal bars represent the number of curated novel gene families detected in each of the 14 GMGC habitats

Functional and evolutionary significance of unknown genes from uncultivated taxa

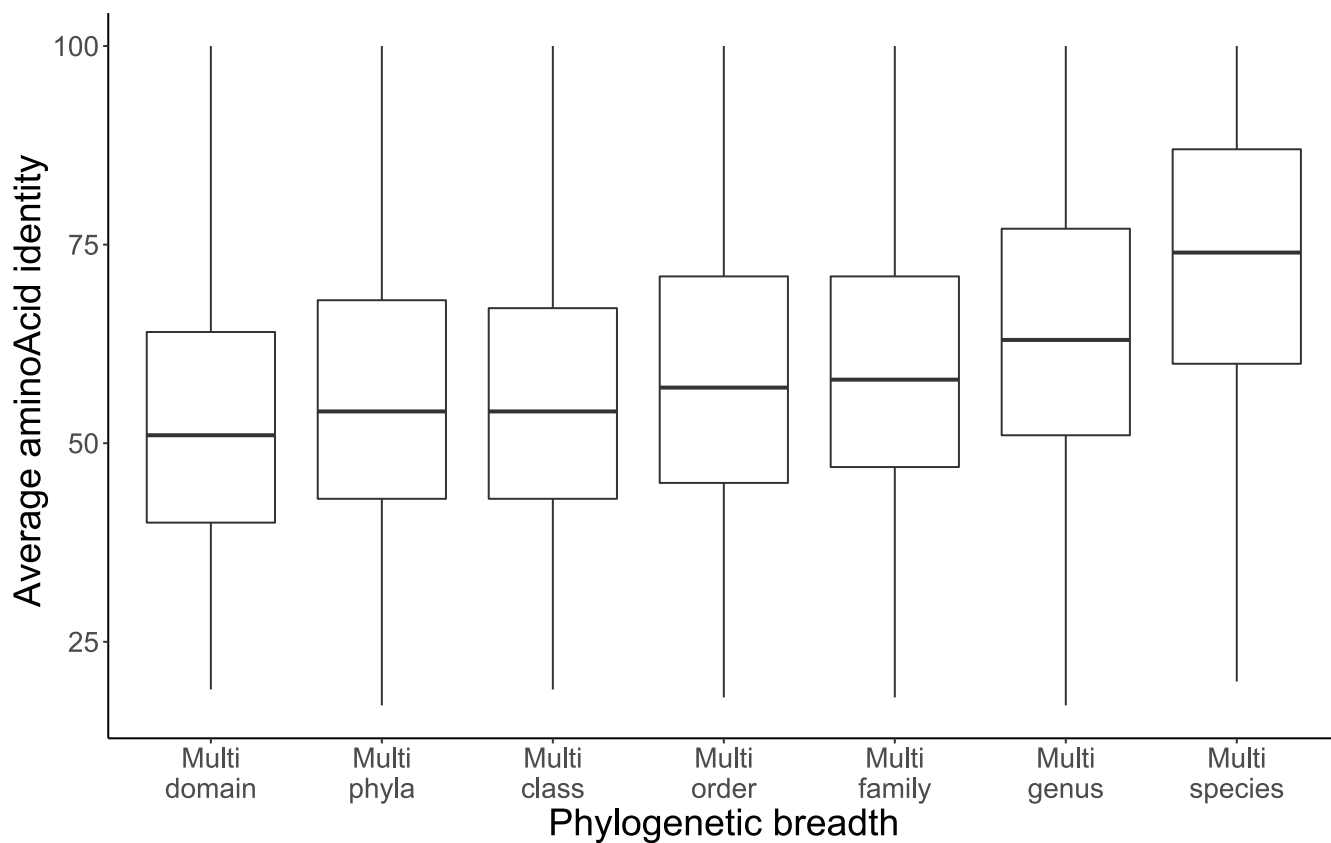


Figure 8.4: Average amino acid identity of the families, stratified by taxonomic breadth

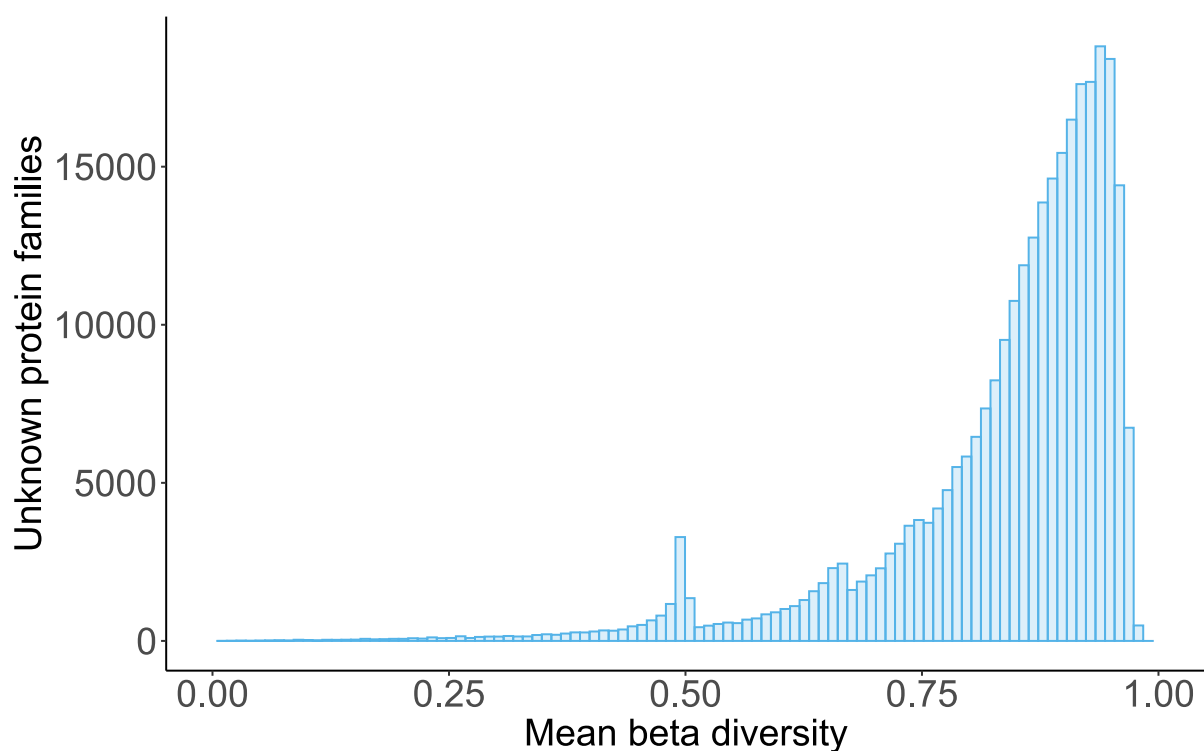


Figure 8.5: Distribution of beta diversity values calculated for the novel gene families

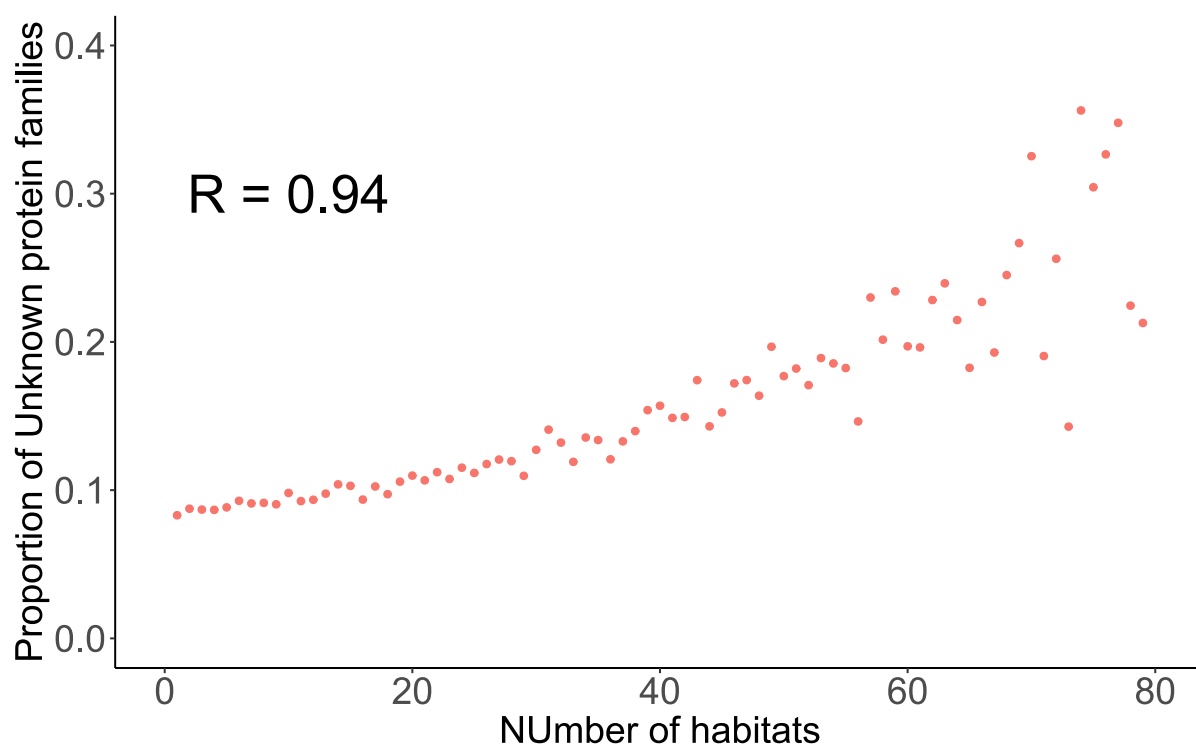


Figure 8.6: Proportion of protein families linked to plasmids or viral contigs with relation to the number of habitats they were detected in. R was calculated as the Spearman correlation coefficient

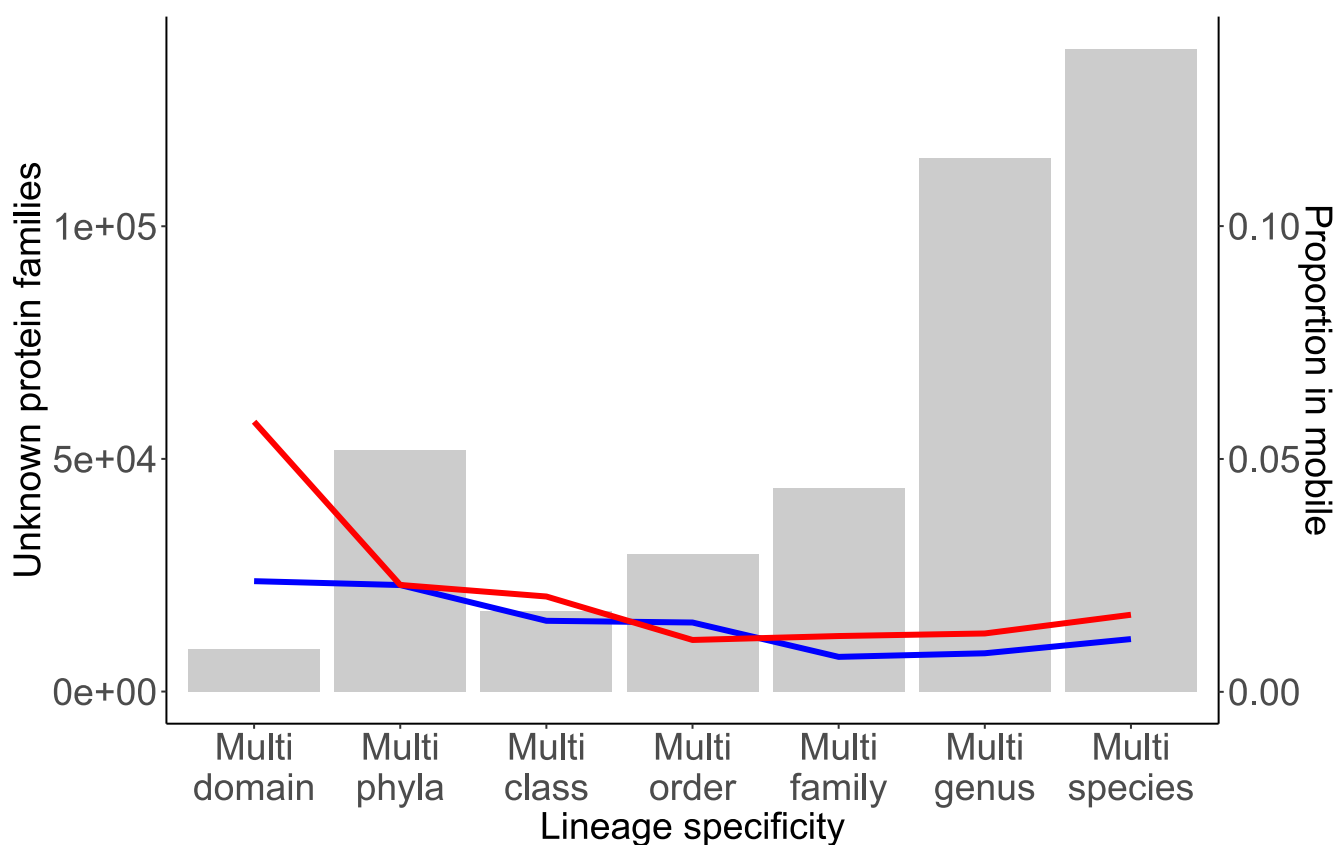


Figure 8.7: Number of unknown protein families confined to each taxonomic rank. The term *genus* in the x-axis indicates the number of protein families detected in multiple species from the same genus, while the *domain* bar indicates families spanning more than one phylum from the same domain. The blue and red lines indicate the proportion of protein families predicted as mobile in plasmids and viral contigs respectively. Equivalent to Figure 4.15C but requiring 30% of the members to be present in mobile elements)

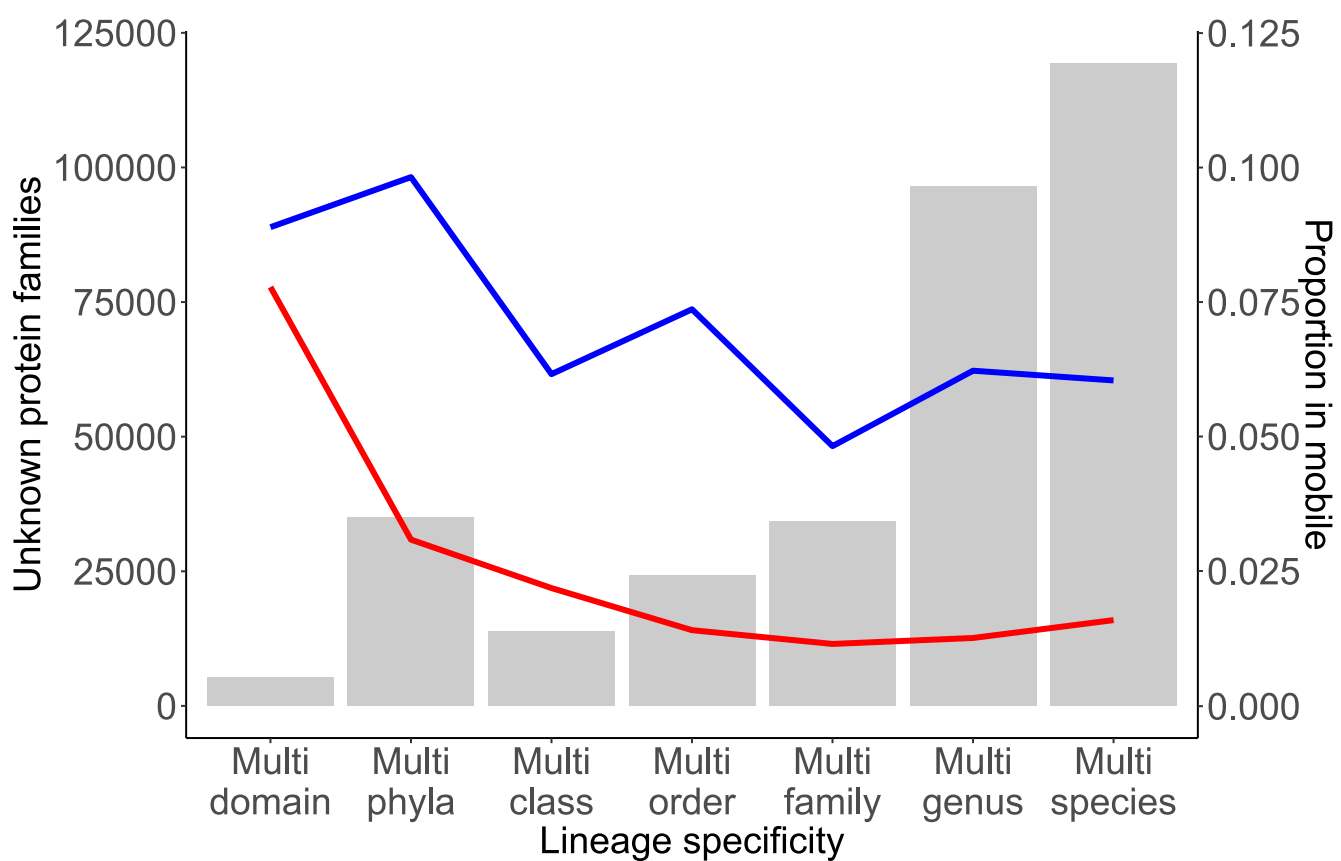


Figure 8.8: Equivalent to Fig. 3 4.15C but calculating LCAs as the most basal taxonomic group gathering 50% of the members of the family. Red: proportion of families in viral contigs. Blue: proportion of families in plasmids.

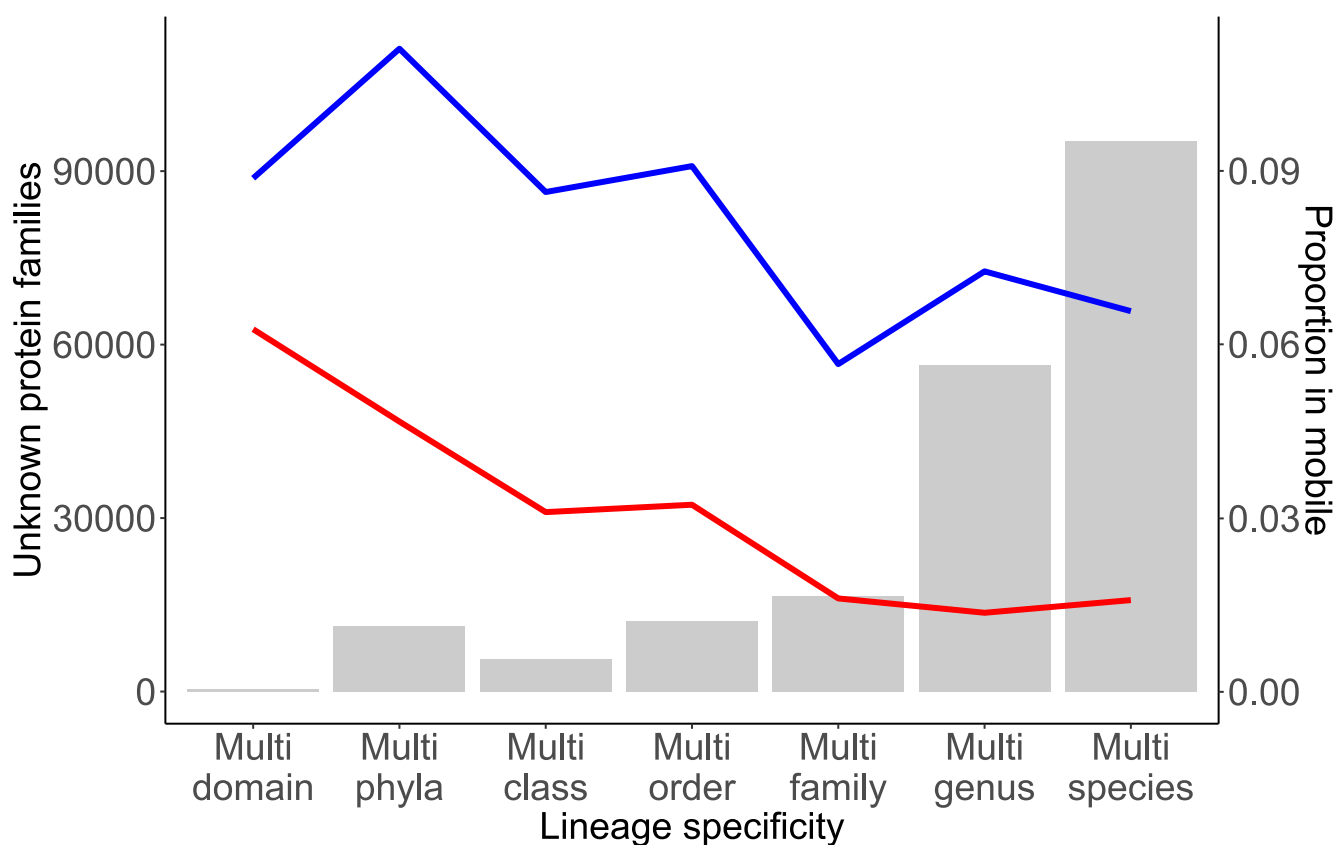


Figure 8.9: Equivalent to Fig. 3 4.15C but calculating LCAs as the most basal taxonomic group gathering 80% of the members of the family. Red: proportion of families in viral contigs. Blue: proportion of families in plasmids.

The code used for generating these results and the supplementary tables were uploaded to <https://github.com/AlvaroRodriguezDelRio/NovFamilies>.

