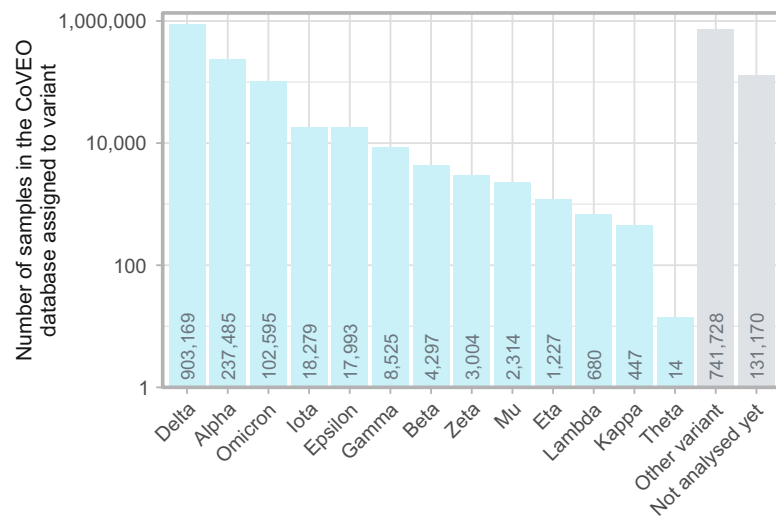
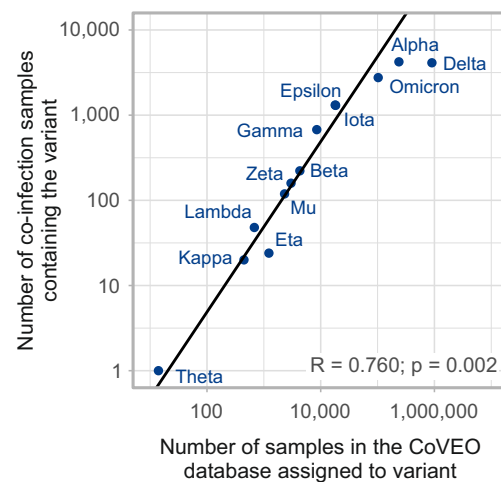


Supplementary Figure 1. The number of samples in the CoVEO database. **a.** The number of good-quality SARS-CoV-2 samples with a human host assigned to different variants in the CoVEO database. **b.** The relationship between the total number of samples assigned to a specific variant and the number of co-infection samples containing the variant. The straight black line represents a linear dependence with a slope of 1 on a log-log graph. R represents the Pearson-correlation coefficient, and the corresponding p-value is derived from a two-sided t-test (n=13 variants).

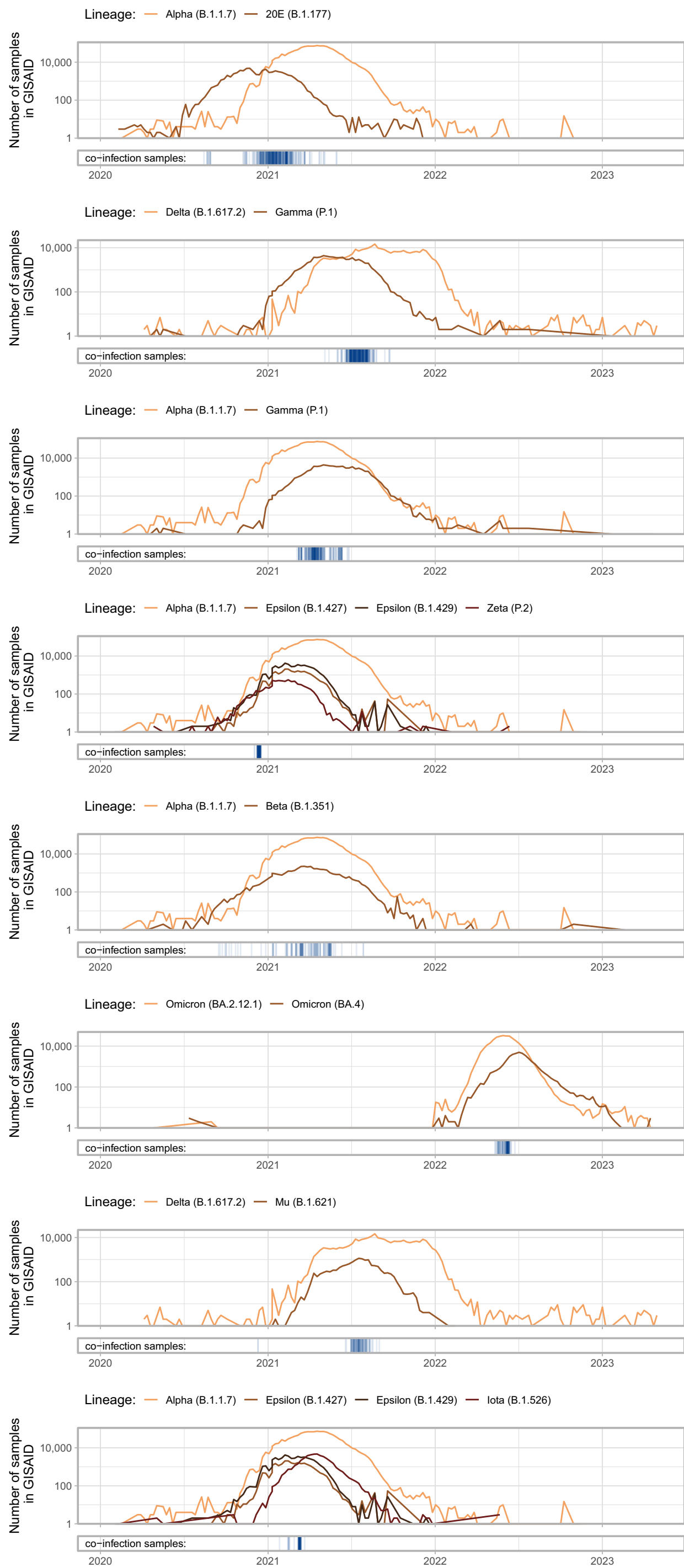
a

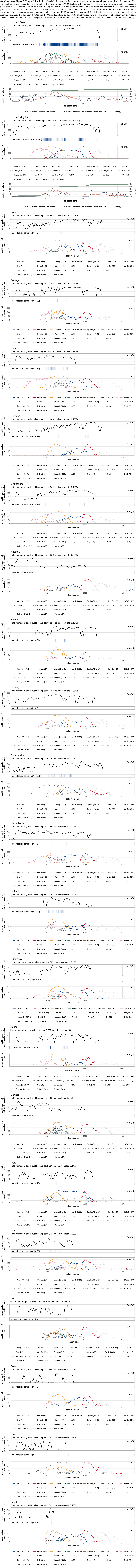


b

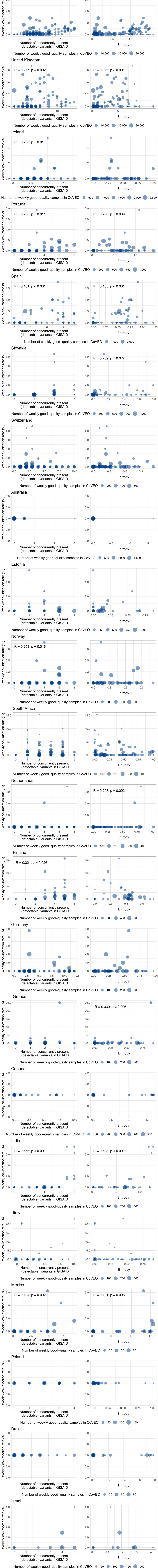


Supplementary Figure 2. Temporal distribution of co-infection samples for variant combinations with more than 50 samples. (The same figures for the top 4 most abundant combinations are shown in Figure 1c of the main manuscript.) Prevalence curves indicate the number of GISAID samples assigned to the respective variants (binned weekly). Blue vertical lines on the bottom panels mark the collection date of co-infection samples of the given variants.



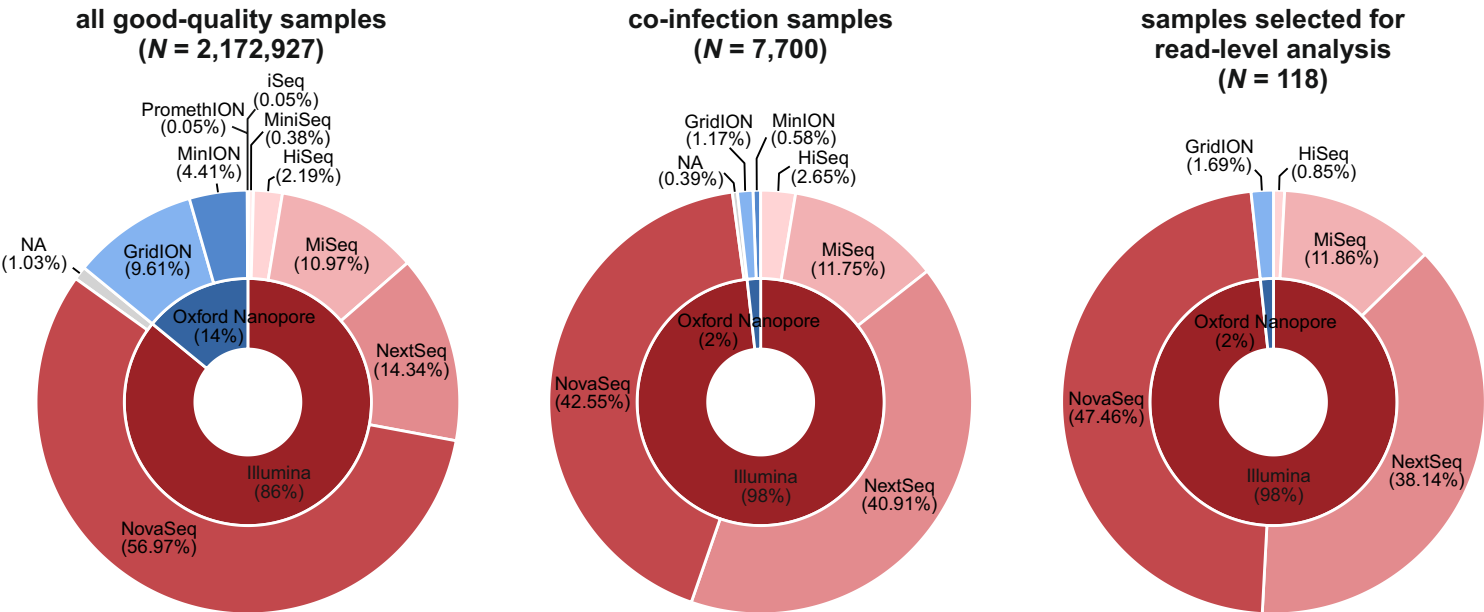


Supplementary Figure 4. Weekly co-infection rate in the function of genetic diversity. Weekly co-infection rate was calculated as the percentage of co-infection samples out of all good-quality samples in the given country, in the given week. Genetic diversity was either defined as the number of lineages concurrently present in the given country in the given week in GISAID data (left panels) or the information entropy (right panels, see Methods). Only lineages investigated by the present study were considered. Weeks for which the number of good-quality samples in the CoVEO database did not reach 10 were discarded. Each marker represents the data for a single week. Marker size corresponds to the number of good-quality samples available. “R” and “p” indicate Pearson-correlation coefficients and respective two-sided t-test p-values (non-significant ($p \geq 0.05$) results are not displayed).



mutational status

Supplementary Figure 6. Distribution of sequencing platforms and instruments within all good-quality samples included in the study, for co-infection samples and for samples selected for downstream read-level analysis.



Supplementary Figure 7. Schematic diagram of the workflow used to produce the data analyzed in the study. Only those steps of the VEO variant calling pipeline (github.com/enasequence/covid-sequence-analysis-workflow) and features of the COVID-19 Data Portal (https://www.covid19dataportal.org/) are shown that are relevant to the current analysis. More details on both can be found in Rahman et al. (2023). For the CoVEO database, only those tables and fields are displayed that were queried during data processing. Queries, codes, data files and visualizations are uploaded to the csabaiBio/SARSCoV2-coinf github repository. Postgres, PostgreSQL and the Slonik Logo are trademarks or registered trademarks of the PostgreSQL Community Association of Canada and used with their permission.

