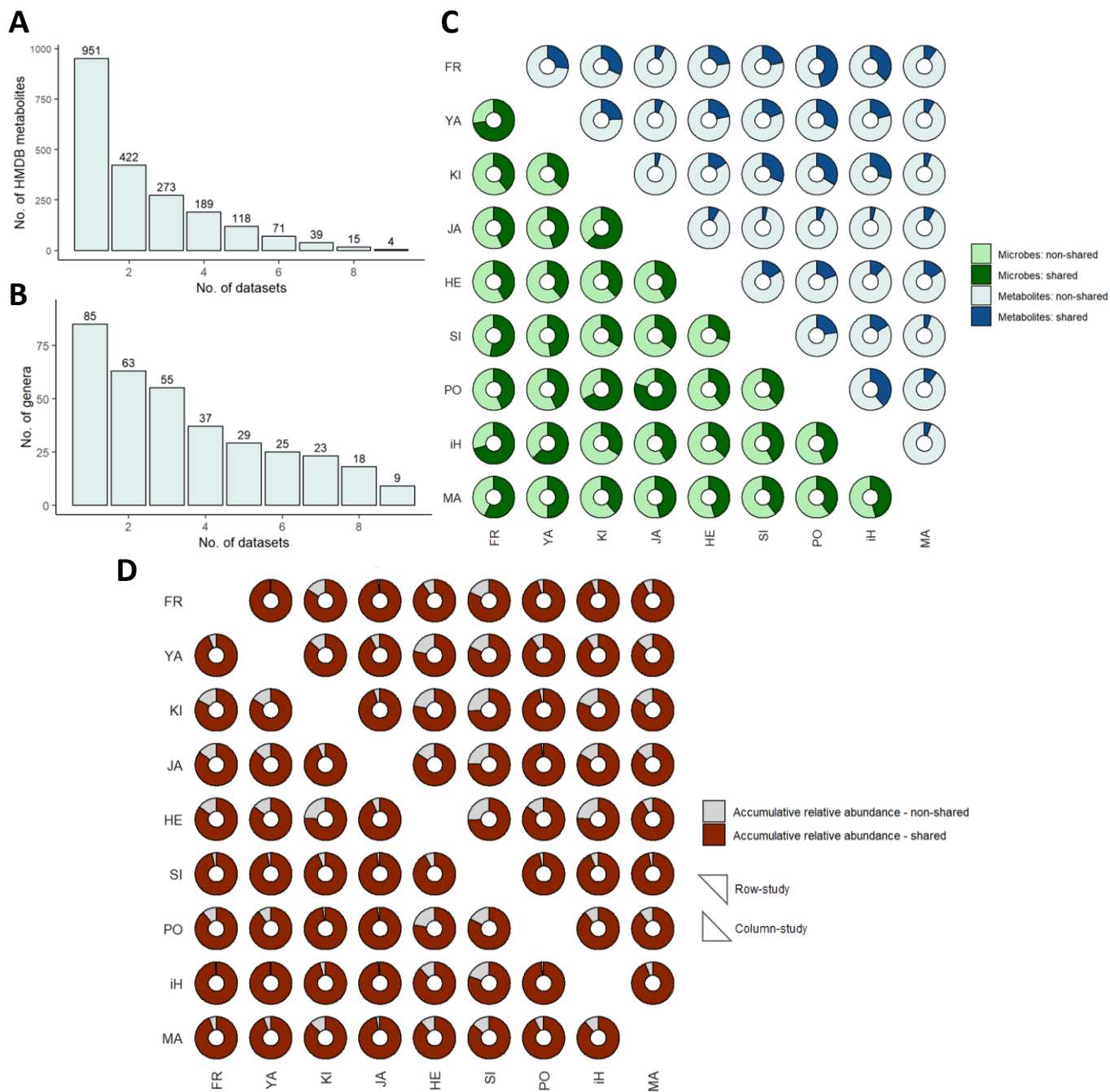
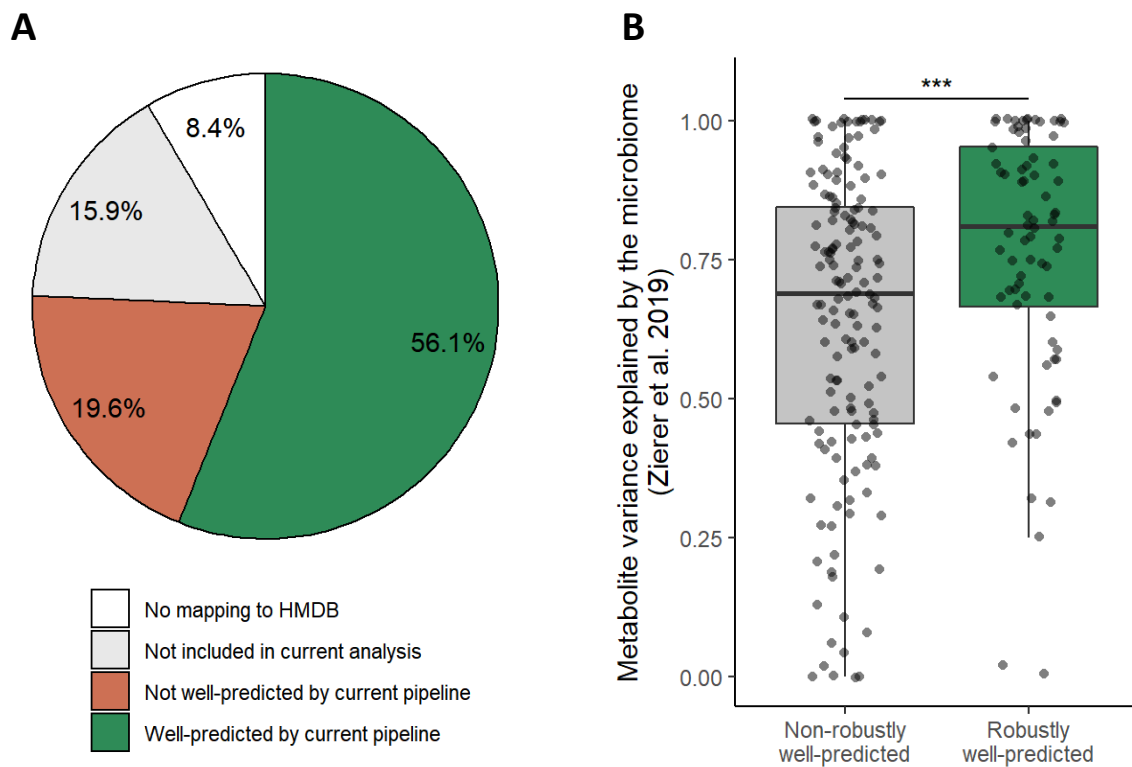


# Figure S1



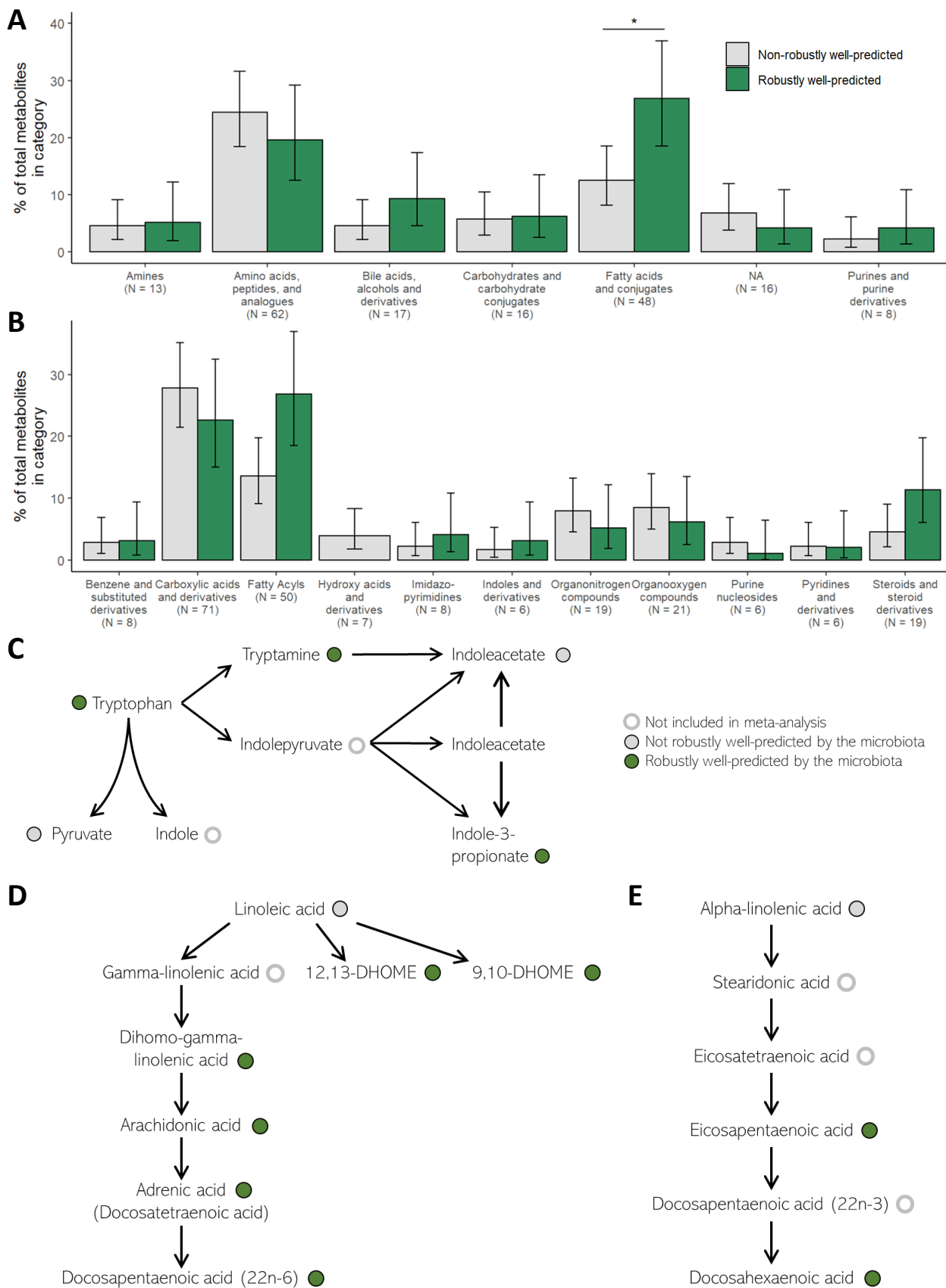
**Figure S1: Statistics about genera and metabolite features shared among datasets.** **(A)** Bar-plot illustrating the cumulative number of metabolites shared by any number of datasets, i.e. each bar represents the number of metabolites shared by at least the specified amount of datasets. Only non-rare metabolites that were successfully mapped to HMDB identifiers are considered (see Methods). **(B)** Bar-plot illustrating the cumulative number of genera shared by any number of studies. Only non-rare genera were considered (see Methods). **(C)** Pairwise comparisons in feature overlap between datasets. Each donut in the matrix represents a pair of studies (“row-study” and “column-study”). Donuts on the upper right triangle represent unique metabolites in both studies, where the dark blue portion represents unique metabolites shared between both studies and light blue represents metabolites available in one dataset but not the other. Similarly, the lower left triangle represents genera overlap between studies, with dark green indicating shared genera and light green indicating non-shared genera. **(D)** Pairwise comparison of genera overlap, considering genera relative abundances in each dataset. As in panel C, each donut represents a pair of studies. In the upper right donuts, the dark red slice indicates the accumulative relative abundance of genera in the row-study which are also present in the column-study (averaged across samples). The grey slice indicates the accumulative relative abundance of the non-shared genera. Lower left donuts should be interpreted similarly, replacing column- and row-studies in the above description. In all 4 plots, only the healthy samples are considered.

## Figure S2



**Figure S2: Comparisons to previous studies. (A)** Comparison of the 107 metabolites identified as well-predicted by Mallick *et al.* (2019), using the FRANZOSA\_IBD dataset, to predictability results obtained here on the same dataset. **(B)** Comparison of robustly well-predicted metabolites to estimates of the proportion of metabolite variance explained by the microbiome as calculated by Zierer *et al.* (2018). The association's significance was tested using a Mann-Whitney test. \*\*\*:  $P < 0.001$ .

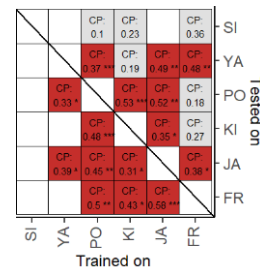
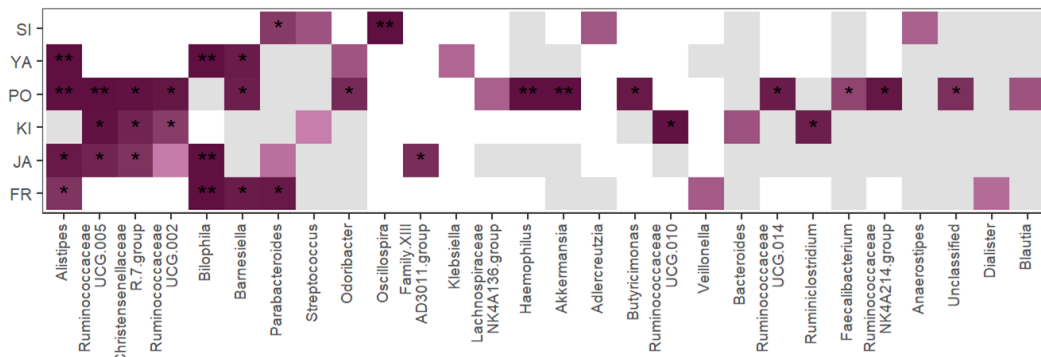
Figure S3



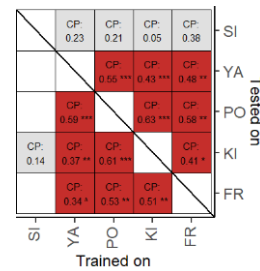
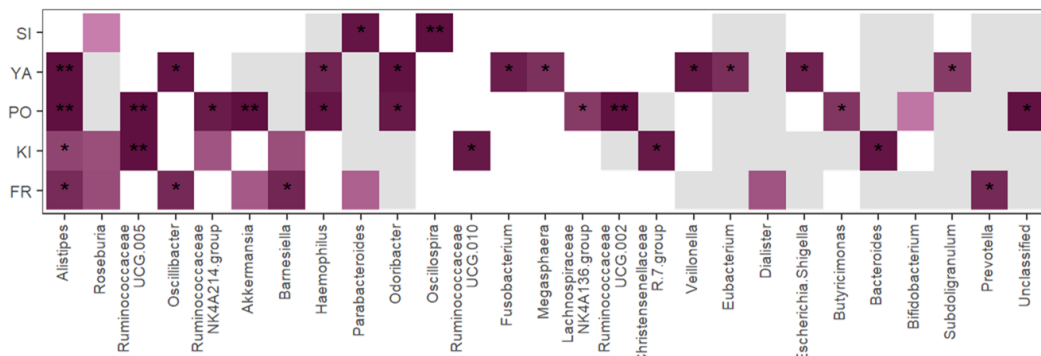
**Figure S3: Characteristics of robustly well-predicted metabolites.** **(A)** Comparison of metabolite distribution over metabolite sub-classes in the HMDB's classification system. Classes with less than 6 metabolites were discarded from this plot. \*: FDR-corrected P value < 0.05. **(B)** Comparison of metabolite distribution over metabolite classes in the HMDB's classification system. Classes with less than 6 metabolites were discarded. **(C)** An illustration of a part of the tryptophan metabolism pathway. Colored circles indicate whether the metabolite was robustly well-predicted by the gut microbiome (green), not robustly well-predicted (gray), or not included in the analysis (white). **(D)**, **(E)** An illustration of a part of the omega-3 and omega-6 metabolism pathways, respectively. Colored circles are as above.

# Figure S4

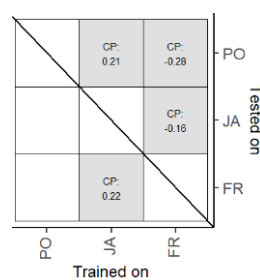
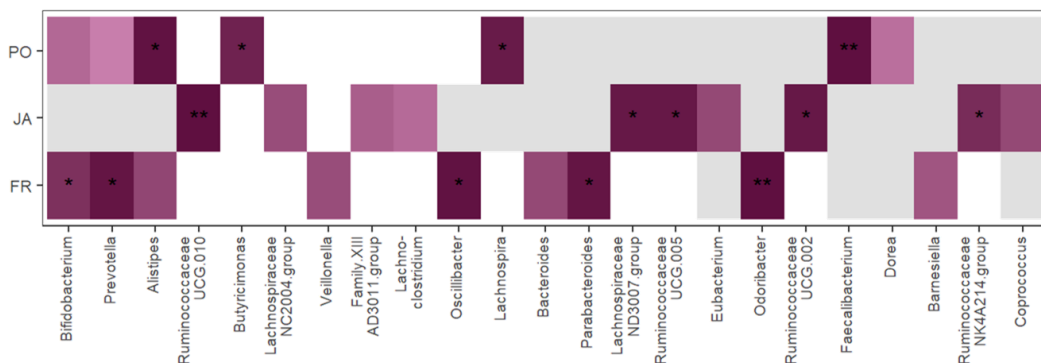
## A Taurine



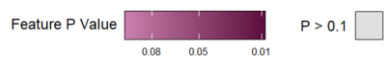
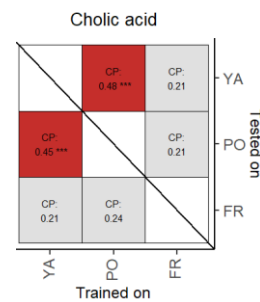
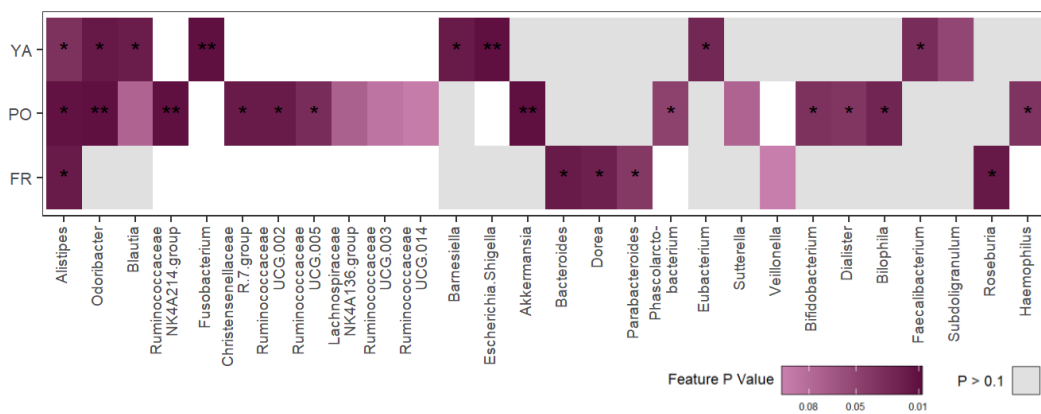
## B N-Acetylputrescine



## C Xanthine



## D Cholic acid



**Figure S4: Comparisons of metabolite models between datasets.** Detailed model comparison for taurine **(A)**, N-acetylputrescine **(B)**, xanthine **(C)** and cholic acid **(D)** metabolites. Panel explanations are as described in Figure 4.

Figure S5

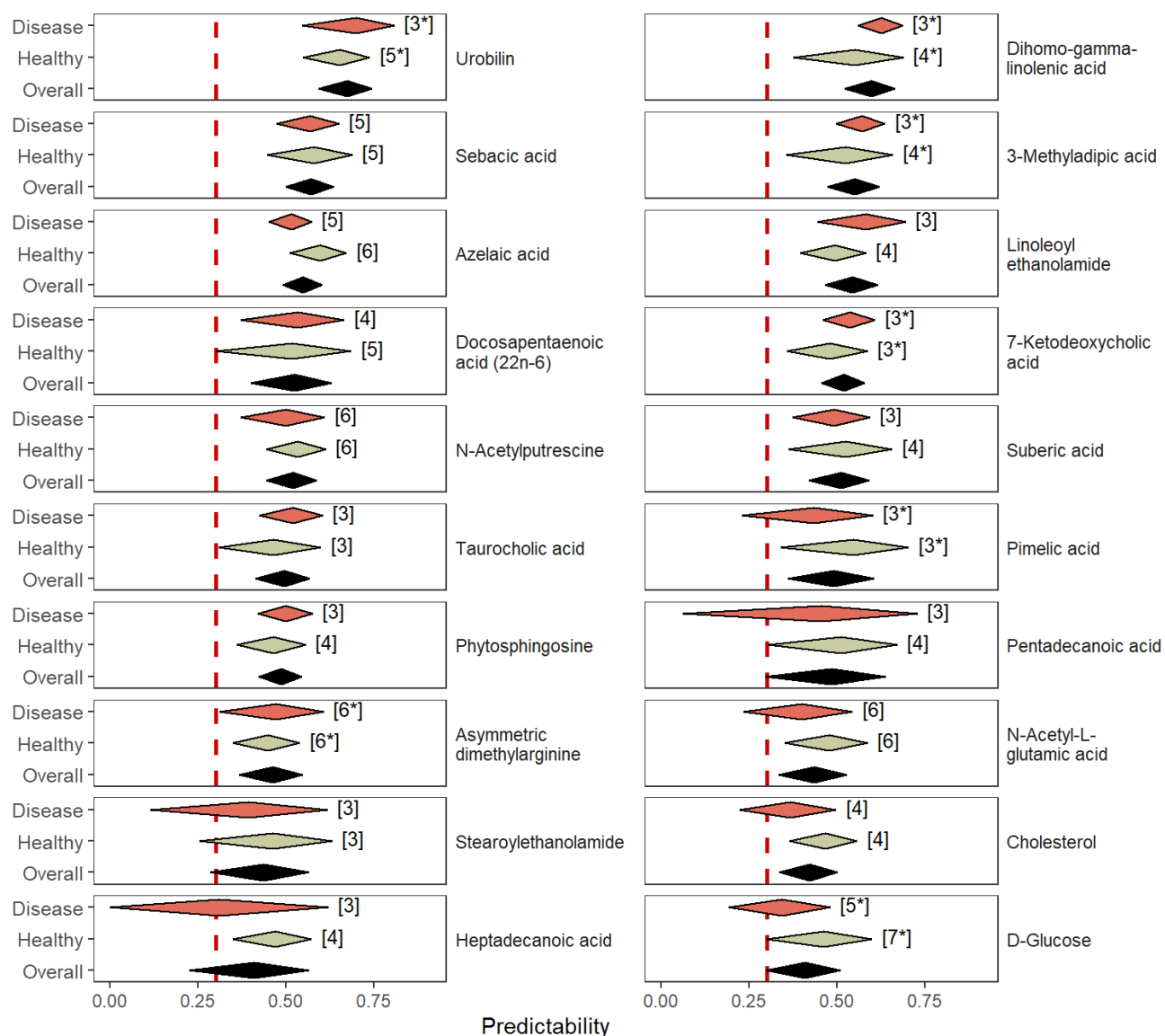
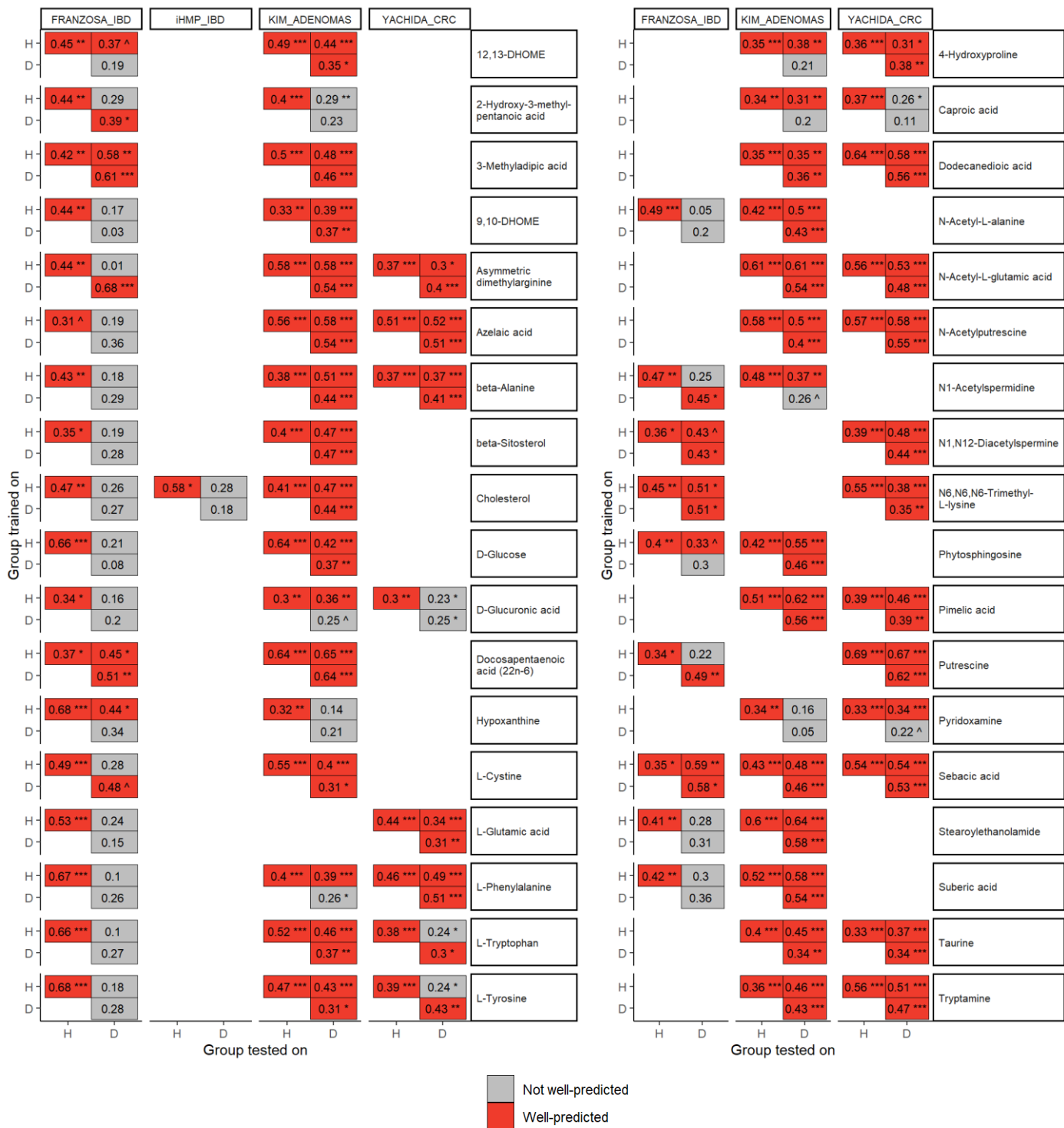


Figure S5: Examples of disease-independent robustly well-predicted metabolites. Top 20 robustly well-predicted metabolites which were also robustly well-predicted in disease. Diamonds and adjacent labels represent REM estimated mean predictability and the number of datasets, respectively, with a star indicating that this metabolite was annotated with low confidence for one or more studies. The red dashed line represents a Spearman’s correlation of 0.3, which we defined as the threshold for a “well-performing” predictive model. The 3 diamonds represent the estimated effects of disease datasets, healthy datasets, and all datasets, respectively.

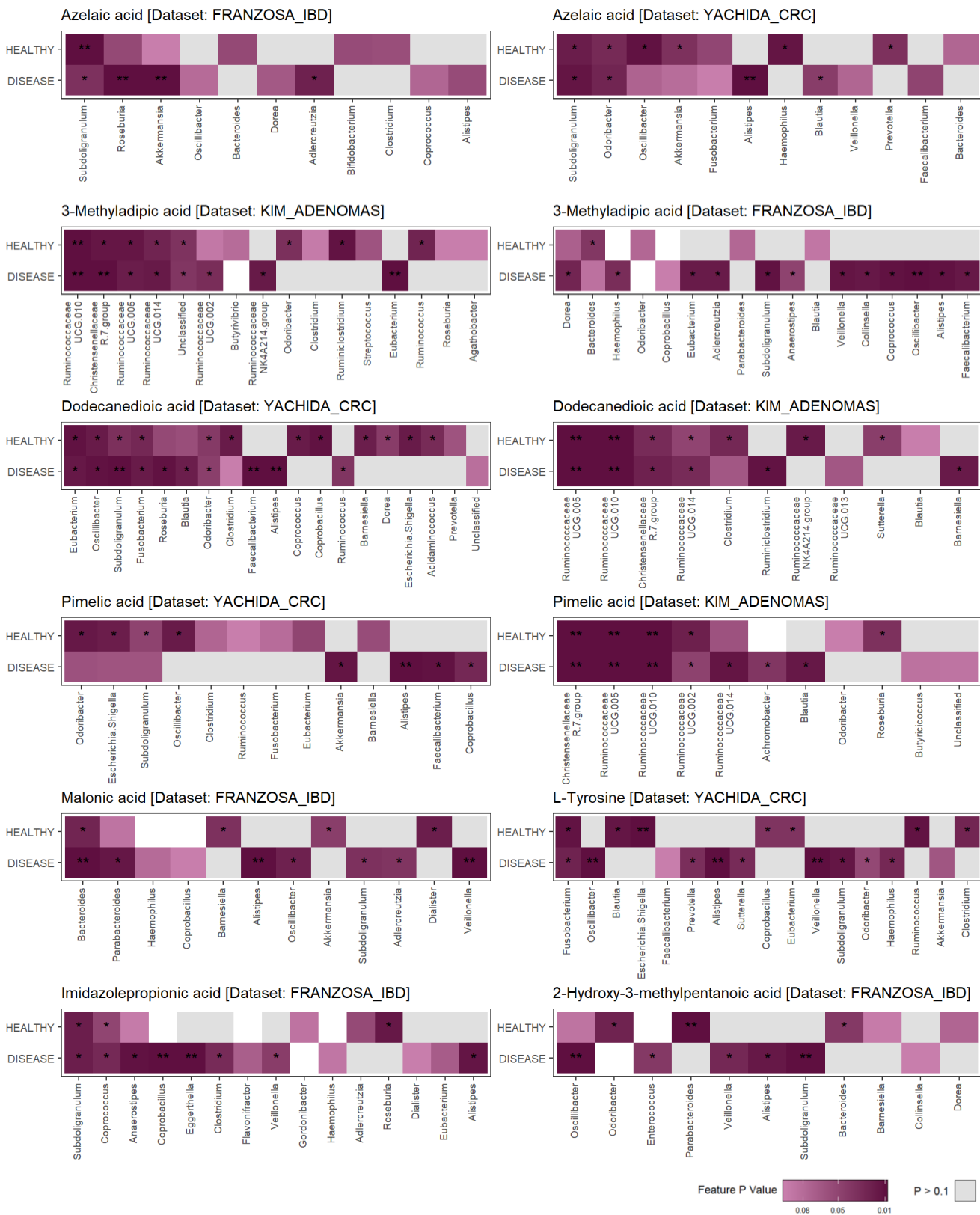


# Figure S6



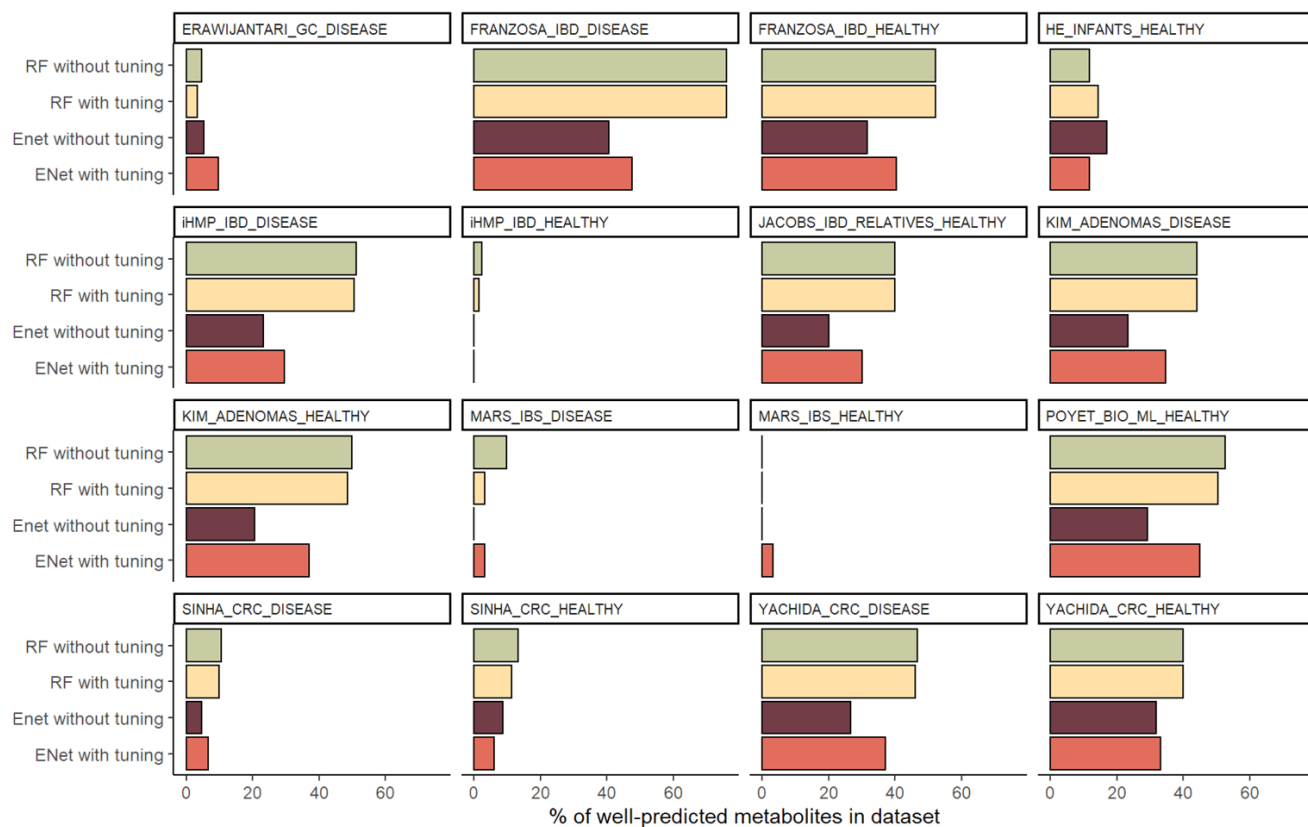
**Figure S6: Transferability of metabolite models from control to cases within the same study.** Cross-predictability analysis results are presented for each case-control study and each metabolite analyzed (see Methods). Only metabolites analyzed in at least 2 datasets are presented. In each panel, the upper left box represents model performance when trained and tested on the control group, the lower right box represents model performance when trained and tested on the cases group, and the upper right box represents model performance when trained on controls and tested on cases. Numbers in cells indicate Spearman's correlation between predicted and actual metabolite levels. Red cells indicate that the metabolite model was considered well-predicted. <sup>^</sup>: P < 0.1; \* : P < 0.05; \*\* : P < 0.01; \*\*\* : P < 0.001.

# Figure S7



**Figure S7: Examples of genus significant contributors' comparison between healthy and disease metabolite models.** Each panel describes significant contributors' comparison for a specific metabolite and dataset. Here upper rows and lower rows represent the control group and case group, respectively, from the same study. Colors and significance marks are as described in Figure 4.

# Figure S8



**Figure S8: Comparison of different machine learning pipelines (related to supplementary note 1).** Each panel describes the percent of well-predicted metabolites out of all metabolites analyzed generated by each machine learning pipeline, in a specific dataset. Both healthy and disease datasets are presented.