# Supplementary information

## Predicting the decision making chemicals used for bacterial growth

Kazuha Ashino[1], Kenta Sugano[2], Toshiyuki Amagasa[2,3], Bei-Wen Ying[1,*]

[1]Graduate School of Life and Environmental Sciences, University of Tsukuba, Ibaraki 305-8572, Japan
[2]Graduate School of Systems and Information Engineering, University of Tsukuba, Ibaraki 305-8573, Japan
[3]Center for Computational Sciences, University of Tsukuba, Ibaraki 305-8577, Japan
*Correspondence: ying.beiwen.gf@u.tsukuba.ac.jp

**Supplementary figures**

**A**

Step 1. Extract the exponential period (**B**, blue lines)

Step 2. Calculate the slopes of two neighbor records

Step 3. Remove the noise and/or error of the slopes

Step 4. Choose the maximal slope (**B**, red line), and average with its two neighbor slopes
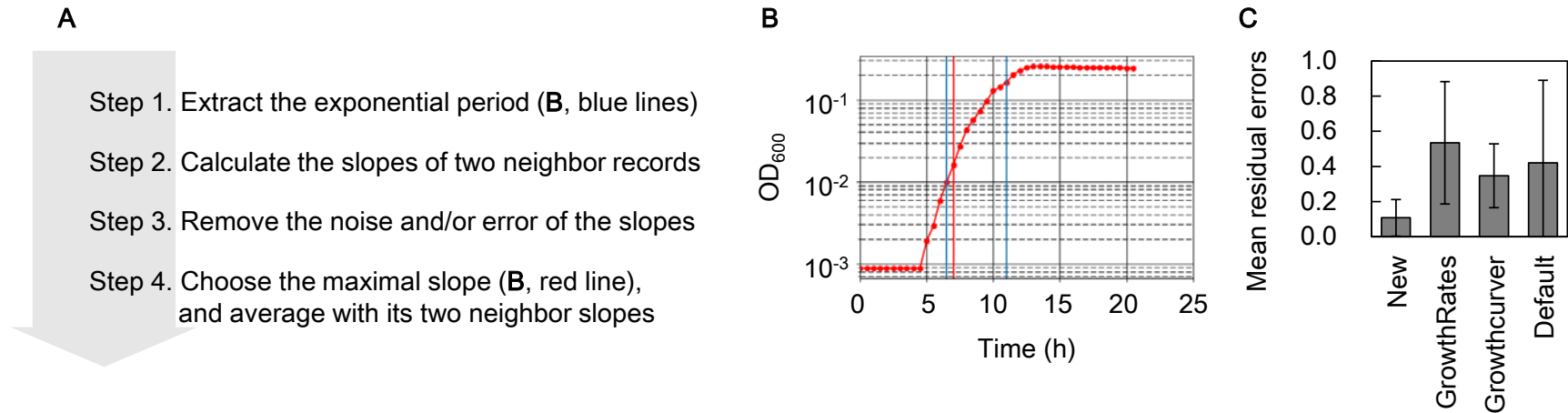
**B**



**C**



**Figure S1 Improved program for calculating the growth rate with Python. A.** Flowchart of data processing for calculating the growth rate. Four steps for estimating the growth rates from the growth curves are illustrated. **B.** A representative growth curve. Blue and red vertical lines represent the borders of the exponential phase and the time-point of the maximal growth rate, respectively. **C.** Comparison of varied tools for evaluating the growth rate. The mean residual errors, representing the precision of the evaluation, were calculated from 60 samples (growth curves) with four different tools. "New" and "Default" indicate the newly developed program in the present study and the tool installed on the plate reader, respectively. GrowthRate and Growthcurver are previously reported programs.
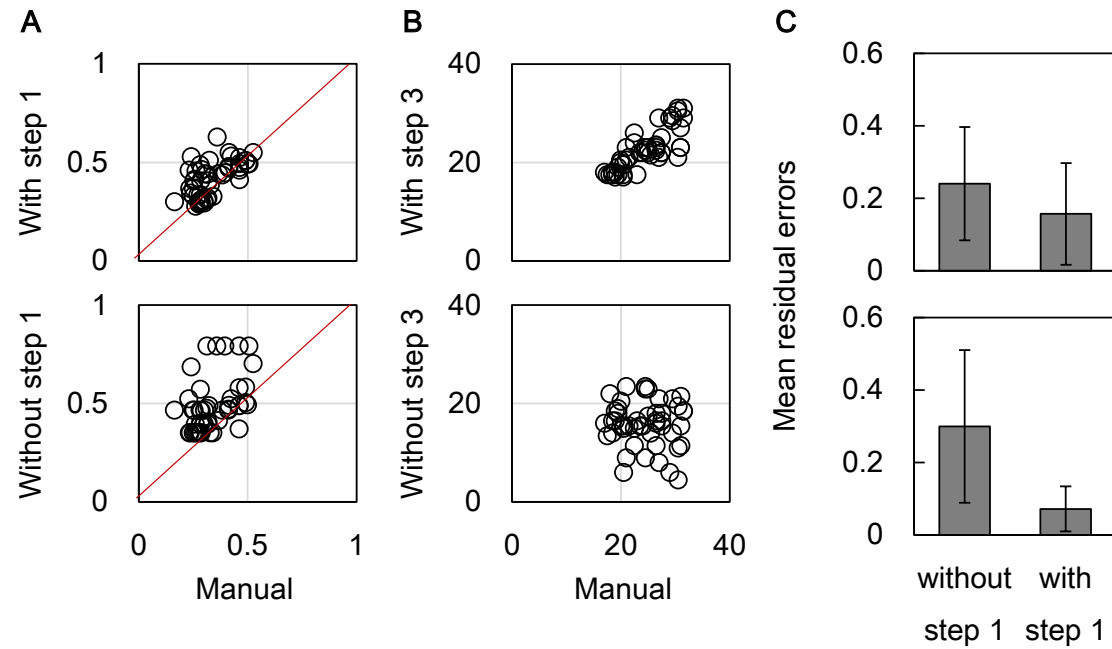
**Figure S2 Improved evaluation of the growth rate due to extracting the exponential period.** Whether step 1 of data processing (Fig. S1A) improved the calculation of the growth rate was examined. The growth rates calculated with or without step 1 (**A**, top and bottom panels, respectively) are plotted against the manually calculated growth rates, which were considered the true values. The red lines indicate equal values. For reference, the time points for the maximum growth rates (slopes) identified with or without step 1 (**B**, upper and bottom panels, respectively) are plotted against the manually identified time points. The mean residual errors of the growth rates and the time points (**C**, top and bottom panels, respectively) are calculated. The addition of step 1 to the program reduced the errors in the growth evaluation.
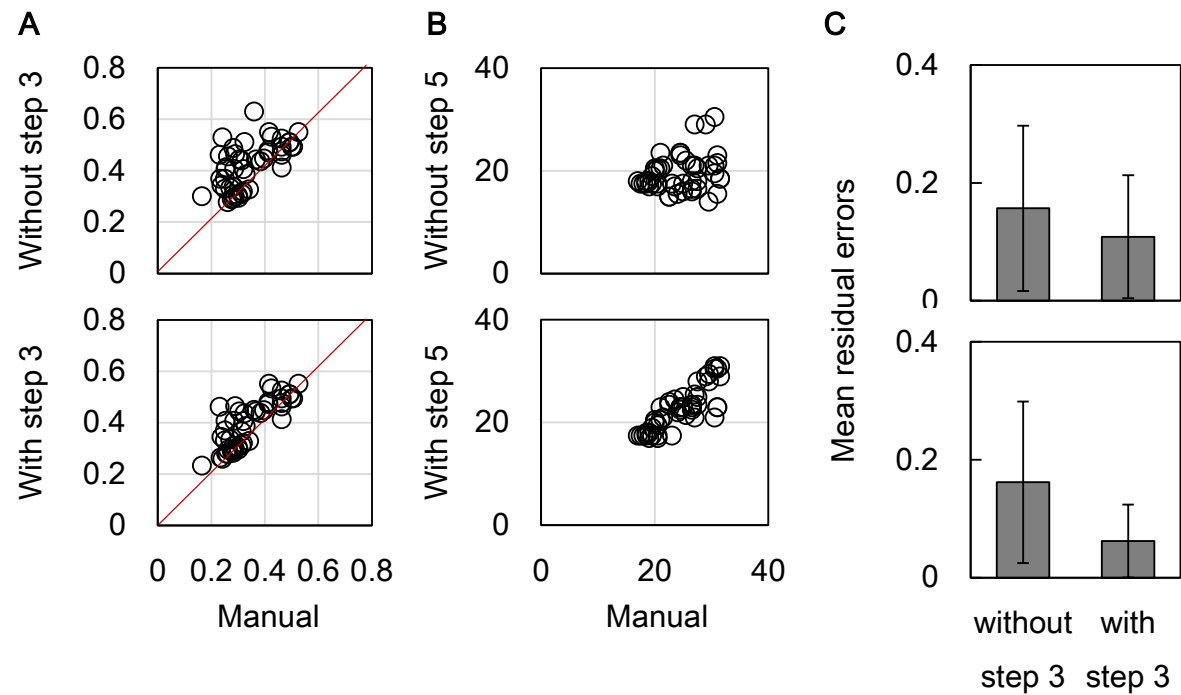
**Figure S3 Improved evaluation of growth rate by removing the noise.** Whether step 3 of data processing (Fig. S1A) improved the calculation of the growth rate was examined. The growth rates calculated with or without step 3 (**A**, top and bottom panels, respectively) are plotted against the manually calculated growth rates, which were considered the true values. The red lines indicate equal values. For reference, the time points of the maximum growth rates (slopes) identified with or without step 3 (**B**, upper and bottom panels, respectively) are plotted against the manually identified time points. The mean residual errors of the growth rates and the time points (**C**, upper and bottom panels, respectively) are calculated. The addition of step 3 to the program reduced the errors in the growth evaluation.
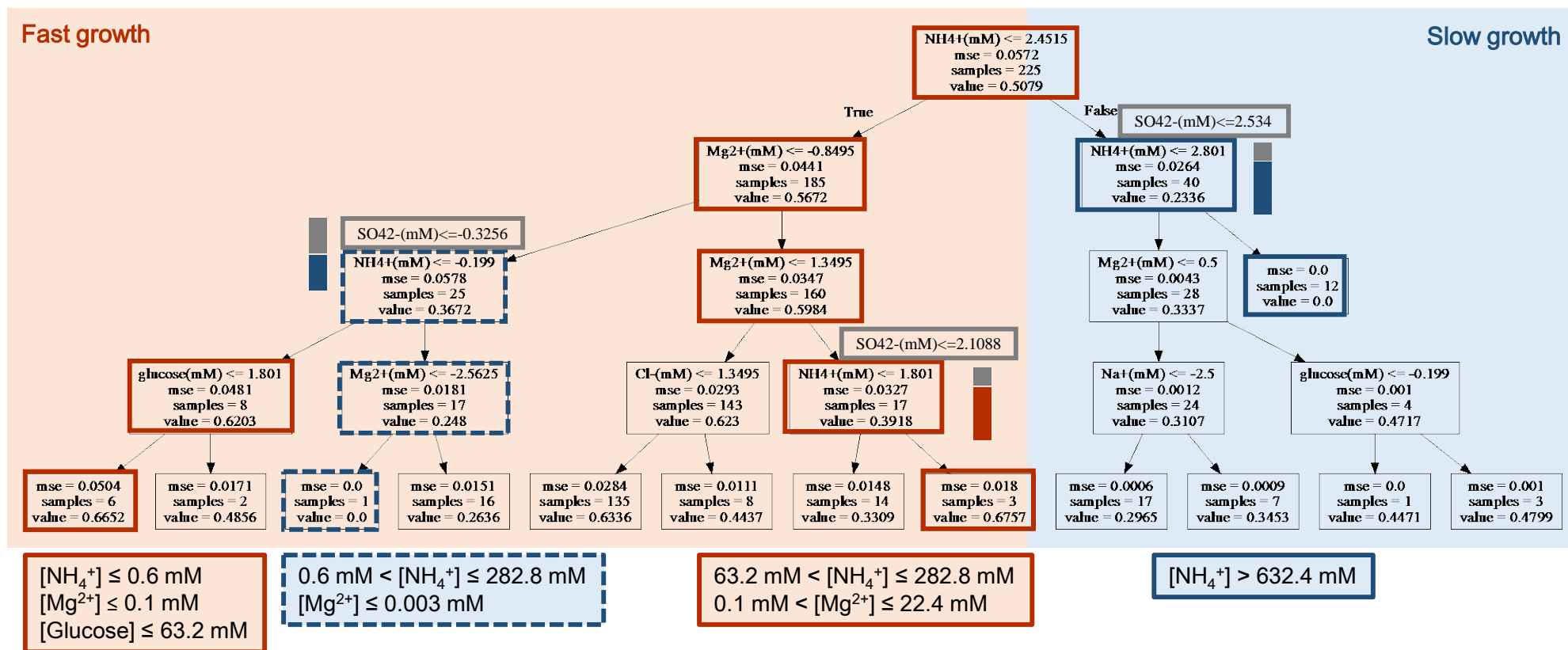
**Figure S4 Decision tree for the growth rate.** The data set for C225 was applied to decision tree learning to predict the growth rate. The resulting tree with a maximal depth of four is shown, and its accuracy is indicated with cross validation (mse). The chemicals predicted as decision elements in the growth rate appear in the tree. Orange and blue illustrate the branches of fast and slow growth, respectively. The name of the selected chemical and its concentration in a logarithmic scale for bifurcation, the value of cross validation for this selection, the number of data used for this selection, and the mean growth rate of the data used for this selection are summarized in the squares, from top to bottom. The squares involved in the paths of the best and worst chemical combinations for growth are highlighted with bold orange and blue lines, respectively. The grey boxes represent the alternative chemicals predicted with the decision tree. The dual-colour bars beside the boxes indicate the frequencies of the alternatives at the same levels, revealing the stability of the tree/prediction. The chemical combinations predicted to cause either fast or zero growth are summarized in the large squares, in which the ranges of the chemical concentrations are shown over linear scales.
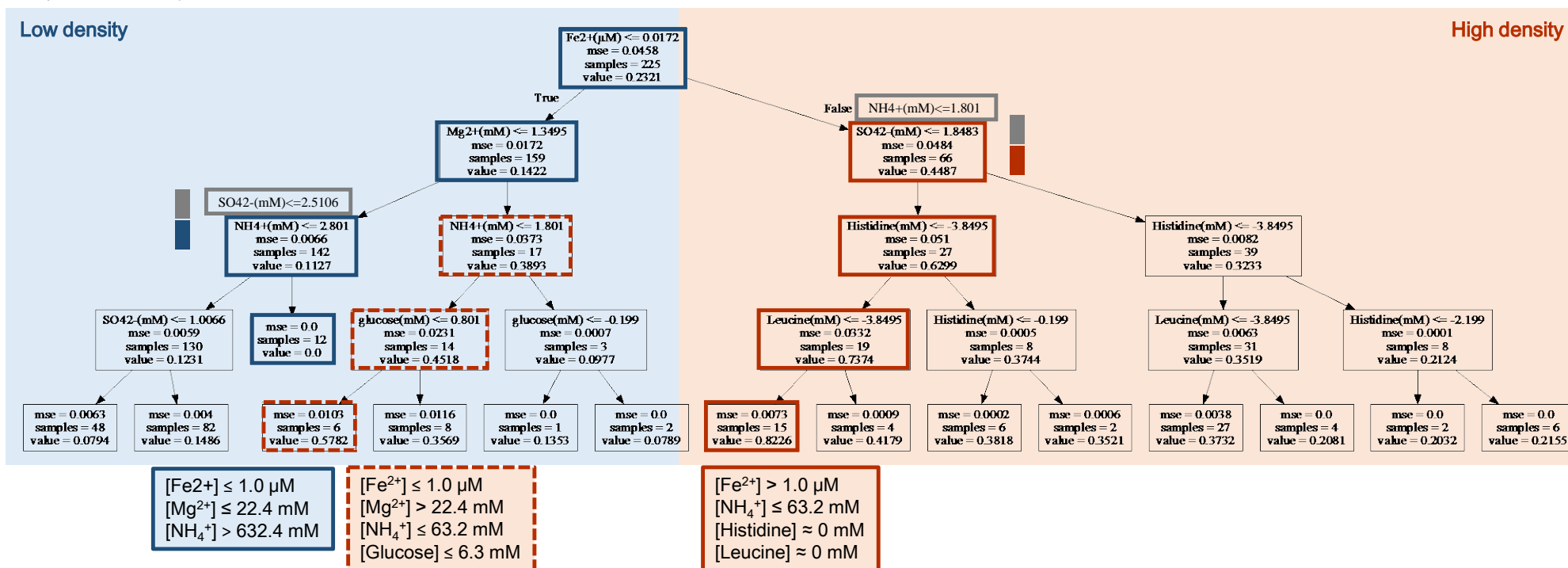
5

Population density, C225, mse=0.01



**Figure S5 Decision tree of the saturated density.** The C225 data set was applied to the decision tree learning to predict the saturated population density. The resulting trees with a maximal depth of four are shown, and its accuracy is indicated with cross validation (mse). The chemicals predicted to be the decision elements for the saturated density appear in the tree. Orange and blue illustrate the branches of high and low density, respectively. The name of the selected chemical and its concentration over a logarithmic scale for bifurcation, the value of cross validation of this selection, the number of data used for this selection, and the mean saturated density of the data used for this selection are summarized in the squares, from top to bottom. The squares involved in the paths of the best and worst chemical combinations for population density are highlighted with bold lines in orange and blue, respectively. The grey boxes represent the alternative chemicals predicted with the decision tree. The dual-colour bars beside the boxes indicate the frequencies of the alternatives at the same levels, revealing the stability of the tree/prediction. The chemical combinations predicted to cause either high or zero density are summarized in the large squares, for which the ranges of chemical concentrations are shown in the linear scales.
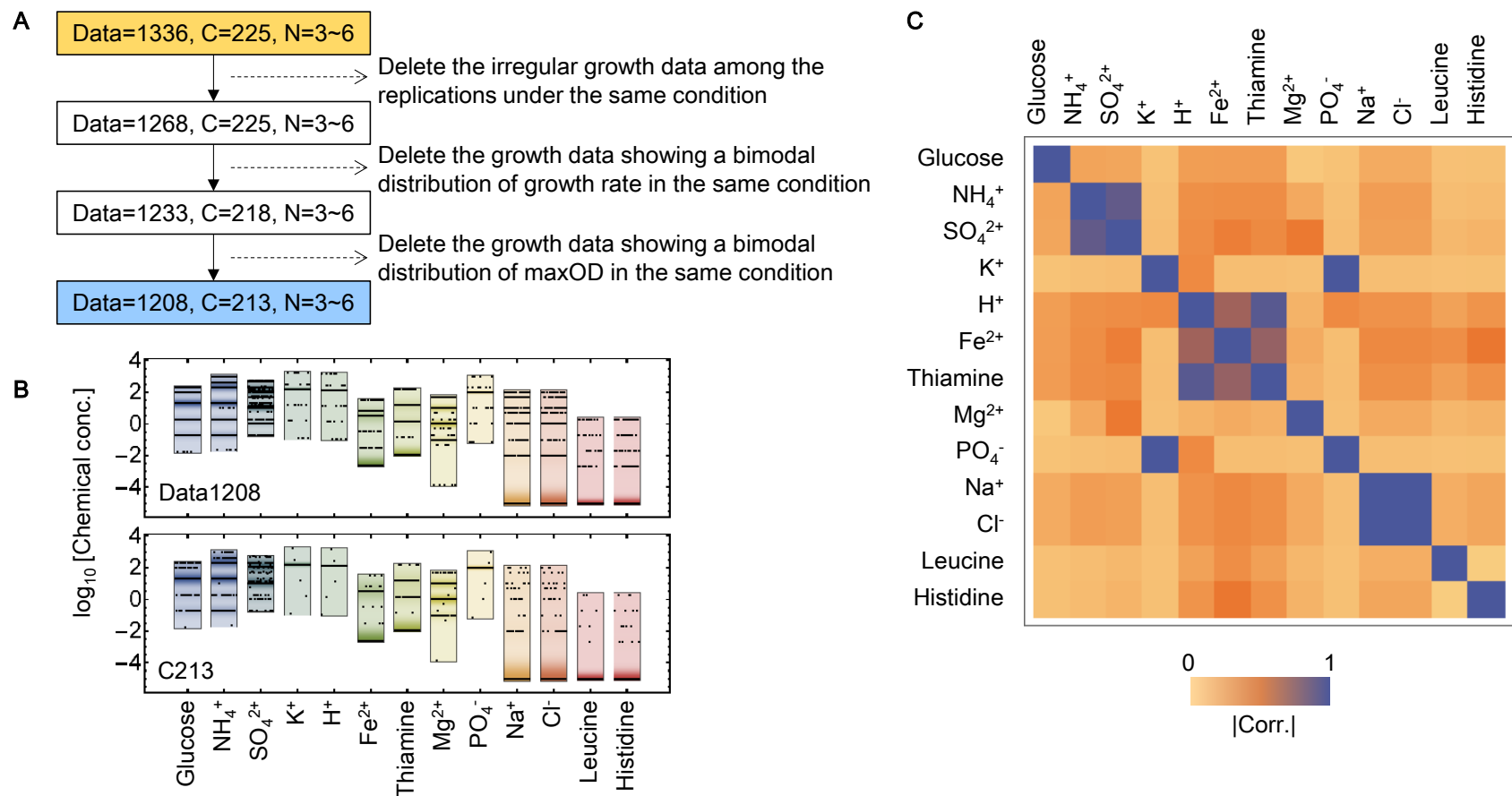
**Figure S6 Data distribution and the chemical combinations of the processed data sets. A.** Flowchart of the data refinement. Data, C, and N represent the numbers of the individual growth curves acquired in the assay, the tested combinations from the ten chemical compounds, and the biological replicates per combination, respectively. The data set highlighted in orange was the purified data set used in the following analyses. **B.** Data distributions at the single chemical level. A total of 13 chemicals, which are encompassed within the ten compounds, are indicated. The data distributions at individual chemical levels are shown as coloured bars. The tested concentrations of these 13 chemicals are spotted in black, within the corresponding distribution bars. Those condensed spots appear as black lines. The top and bottom panels indicate the distributions of 1208 individual growth curves and the 213 combinations, respectively. **C.** Relationships between the concentration changes of 13 individual chemicals. The matrix represents the correlations of the changes in the concentrations of any two chemicals. The gradation from light orange to dark blue indicates the correlation coefficients from low to high.
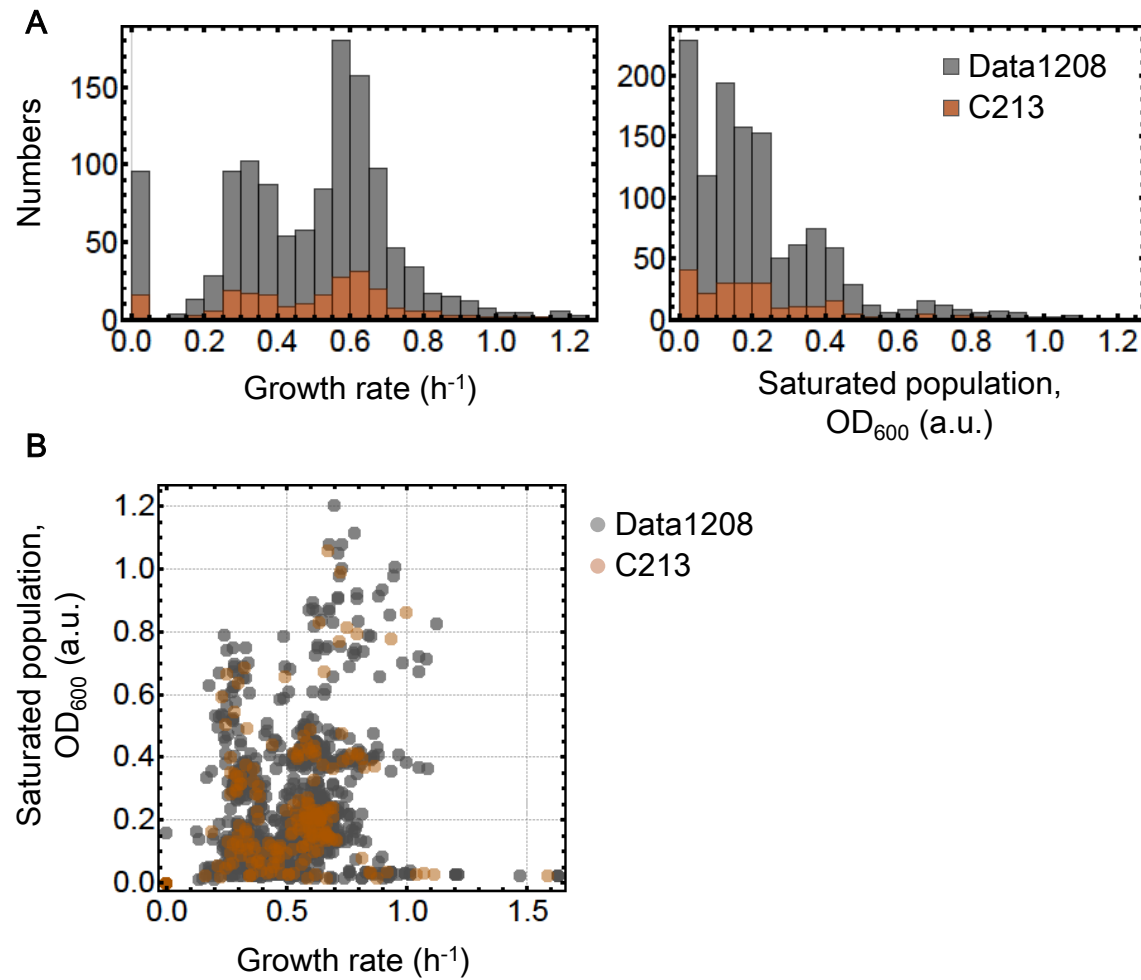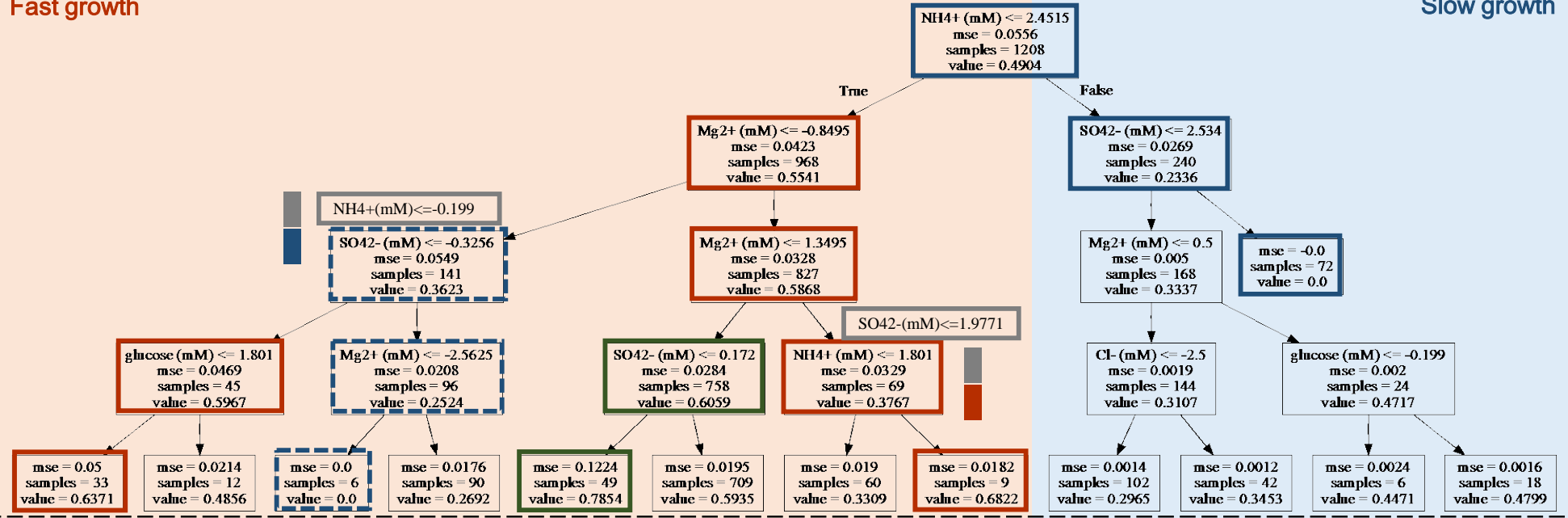
**Figure S7 Growth rates and saturated population densities of the processed data sets. A.** Histograms of the growth rate and the saturated density. The left and right panels show the histograms of the growth rate and the saturated density, respectively. **B.** Relation between the growth rate and the saturated density. The Spearman rank correlation coefficients of Data1208 and C213 are 0.33 ($p$=5e-33) and 0.61 ($p$=7e-23), respectively. The Data1208 and C213 data sets are indicated in grey and bronze, respectively.
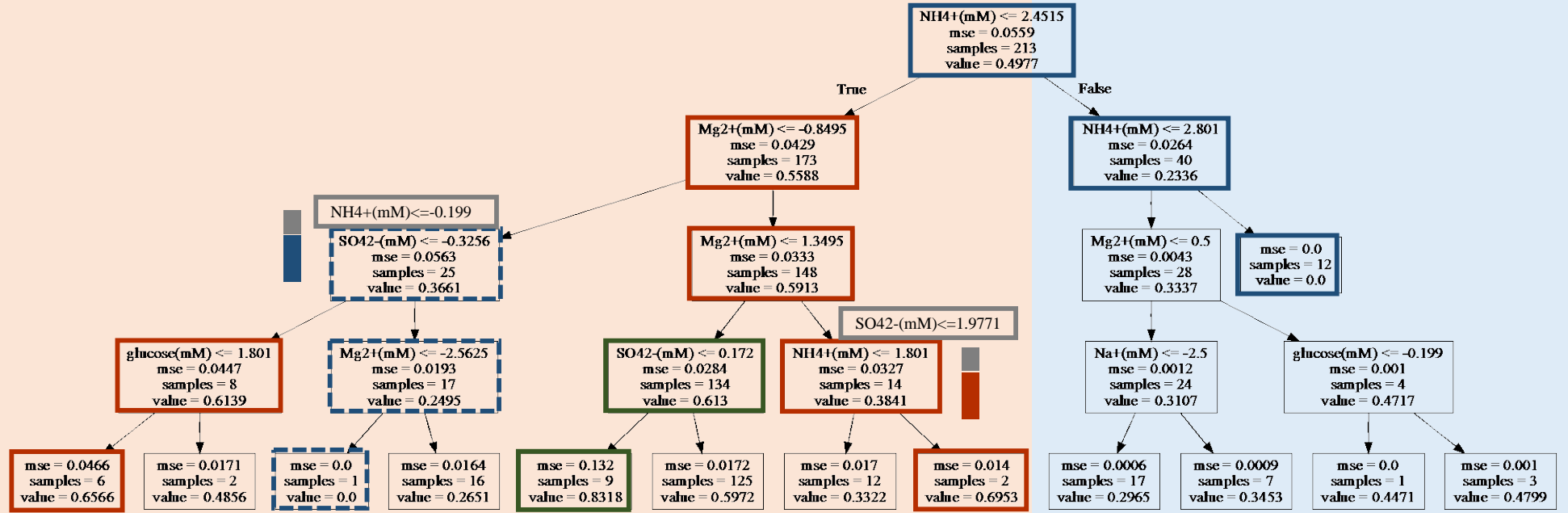
**Figure S8 Decision tree of the growth rate predicted using the clean data sets.** The purified Data1208 and C213 data sets (upper and bottom panels, respectively) were applied to the decision tree learning to predict the growth rate. The resulting trees with a maximal depth of four are shown, and the accuracy is indicated with cross validation (mse). The chemicals predicted to be the decision elements for the growth rate appear in the tree. Orange and blue illustrate the branches of fast and slow growth, respectively. The name of the selected chemical and its concentration over a logarithmic scale for bifurcation, the value of cross validation of this selection, the number of data used for this selection, and the mean growth rate of the data used for this selection are summarized in the squares, from top to bottom. The squares involved in the paths directed towards fast growth are highlighted with bold lines in orange and green, and those with zero growth are in blue. The grey boxes represent the alternative chemicals predicted with the decision tree. The dual-colour bars beside the boxes indicate the frequencies of the alternatives at the same levels, revealing the stability of the tree/prediction. The decision trees were highly identical to those acquired with the noisy data sets (Data1336, C225). An alternative path for fast growth (highlighted in green) was newly determined.
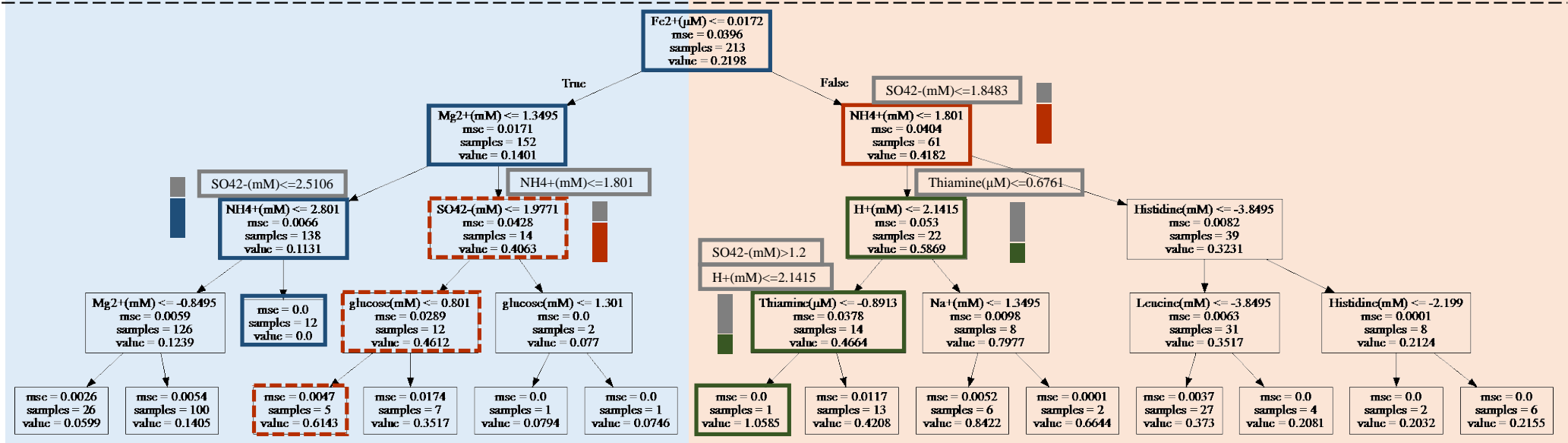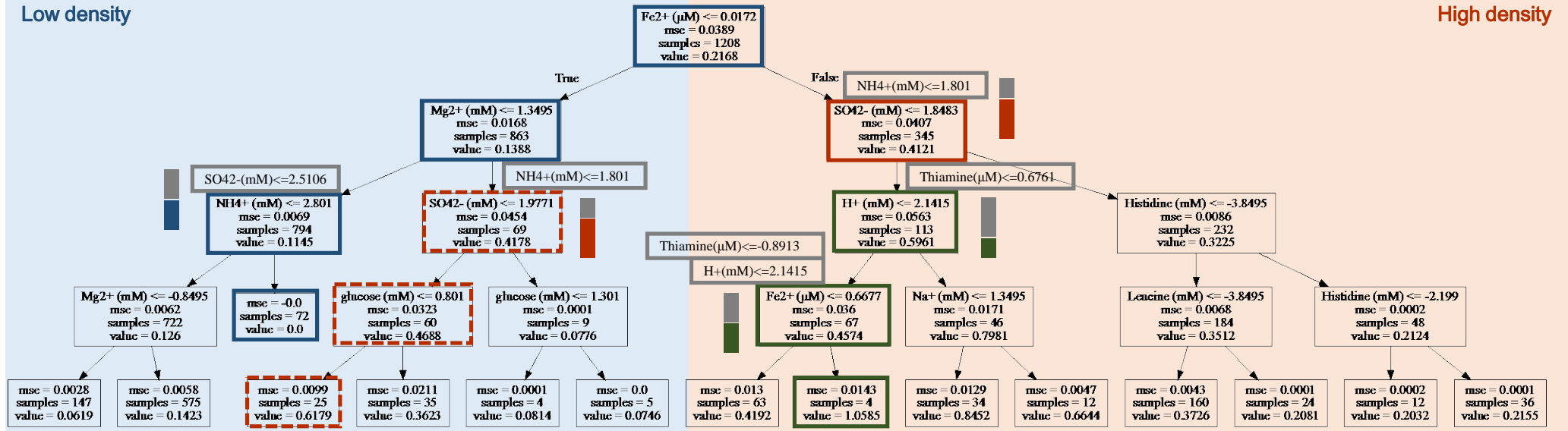
**Figure S9 Decision tree of the saturated density predicted with the clean data sets.** The purified Data1208 and C213 data sets (the top and bottom panels, respectively) were applied to the decision tree learning to predict the saturated population density. The resulting trees with a maximal depth of four are shown, and the accuracy is indicated with cross validation (mse). The chemicals predicted to be the decision elements for saturated density appear in the tree. Orange and blue illustrate the branches of high and low density, respectively. The name of the selected chemical and its concentration in a logarithmic scale for bifurcation, the value of cross validation for this selection, the number of data used for this selection, and the mean saturated density of the data used for this selection are summarized in the squares, from top to bottom. The squares involved in the paths directed towards the high density are highlighted with bold lines in orange and green, and those moving towards zero density are in blue. The grey boxes represent the alternative chemicals predicted with the decision tree. The dual-colour bars beside the boxes indicate the frequencies of the alternatives at the same levels, revealing the stability of the tree/prediction. Low-density branches were highly identical to those based on the noisy data sets (Data1336, C225). An alternative path for the highest density was newly determined (highlighted in green).