# Using Machine Learning to Predict Antimicrobial Minimum Inhibitory Concentrations and Associated Genomic Features for Nontyphoidal *Salmonella*

## Supplemental Information

Marcus Nguyen, S. Wesley Long, Patrick F. McDermott, Randall J. Olsen, Robert Olson, Rick L. Stevens, Gregory H. Tyson, Shaohua Zhao and James J. Davis

**Supplemental Figures**

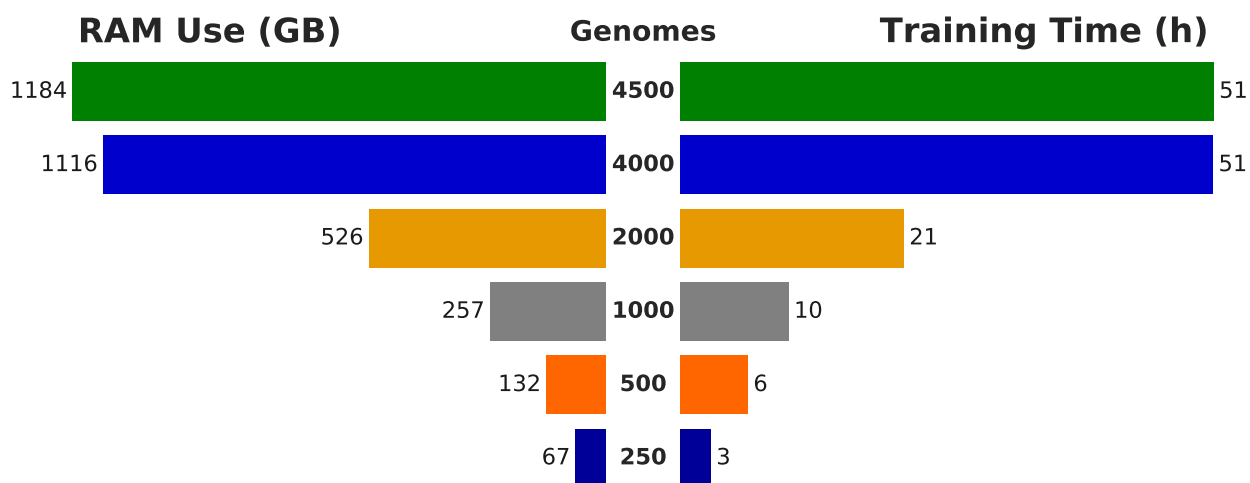| RAM Use (GB) | | Genomes | | Training Time (h) |
|---|---|---|---|---|
| 1184 | | **4500** | | 51 |
| 1116 | | **4000** | | 51 |
| 526 | | **2000** | | 21 |
| 257 | | **1000** | | 10 |
| 132 | | **500** | | 6 |
| 67 | | **250** | | 3 |

**Figure S1.** Tornado plot of the peak memory use (GB) and time required to train a 5-fold cross-validated model (hours), using a large server (4 Intel E5-4669v4 CPUs @ 2.20GHz, 1,585,224,876 kB RAM) and the data sets from Figure 1. Training was done using 170 of the 174 available logical cores. From top to bottom, bars in the plot correspond to training sets containing 4500, 4000, 2000, 1000, 500 and 250 genomes respectively. Values on the left depict peak RAM usage, and values on the right depict training time.

Tree scale: 0.01

Model Accuracy
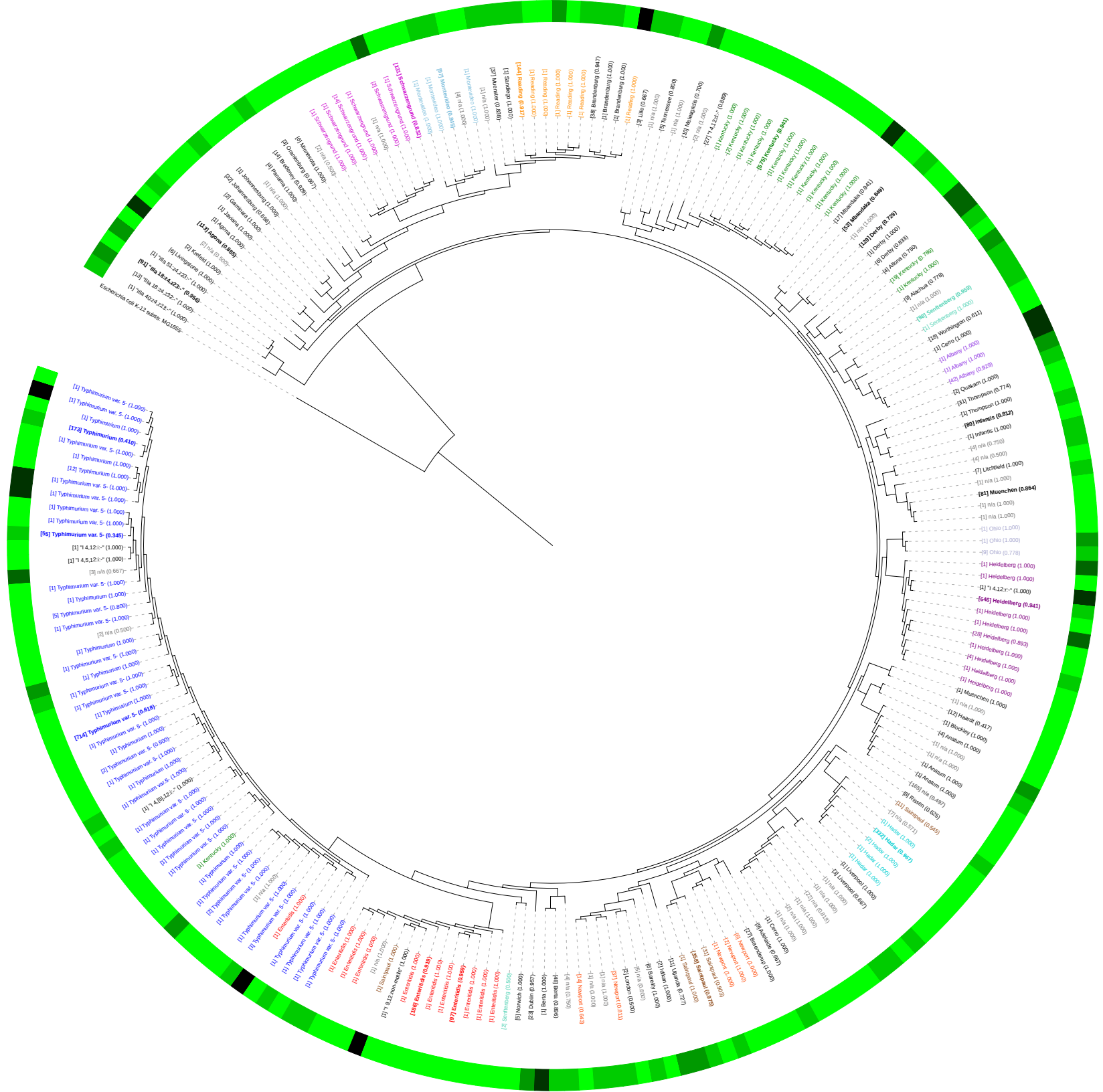95-100%
90-95%
85-90%
80-85%
75-80%
≤70%

**Figure S2.** Phylogenetic tree of the *Salmonella* strains used in this study. The tree is based on a concatenation of alignments for the *rpoB* and *rpoC* genes and is rooted on *E. coli* K-12. Zero-length branches are collapsed and are depicted by a single tip and labeled according to the most common serotype occurring at that tip. The number in square brackets is the number of genomes depicted by a tip, and tips with ≥ 50 taxa are shown in bold. The number in parentheses is the fraction of genomes at a given tip that belong to the listed serotype. Tips representing common serotypes are colored for ease of identification. The color bar surrounding the tree represents the average accuracy (within ± 1 two-fold dilution step) of the 4,500 genome model for the genomes represented by each tip, with black being the least accurate (≤ 70%) and bright green being the most accurate (95-100%).
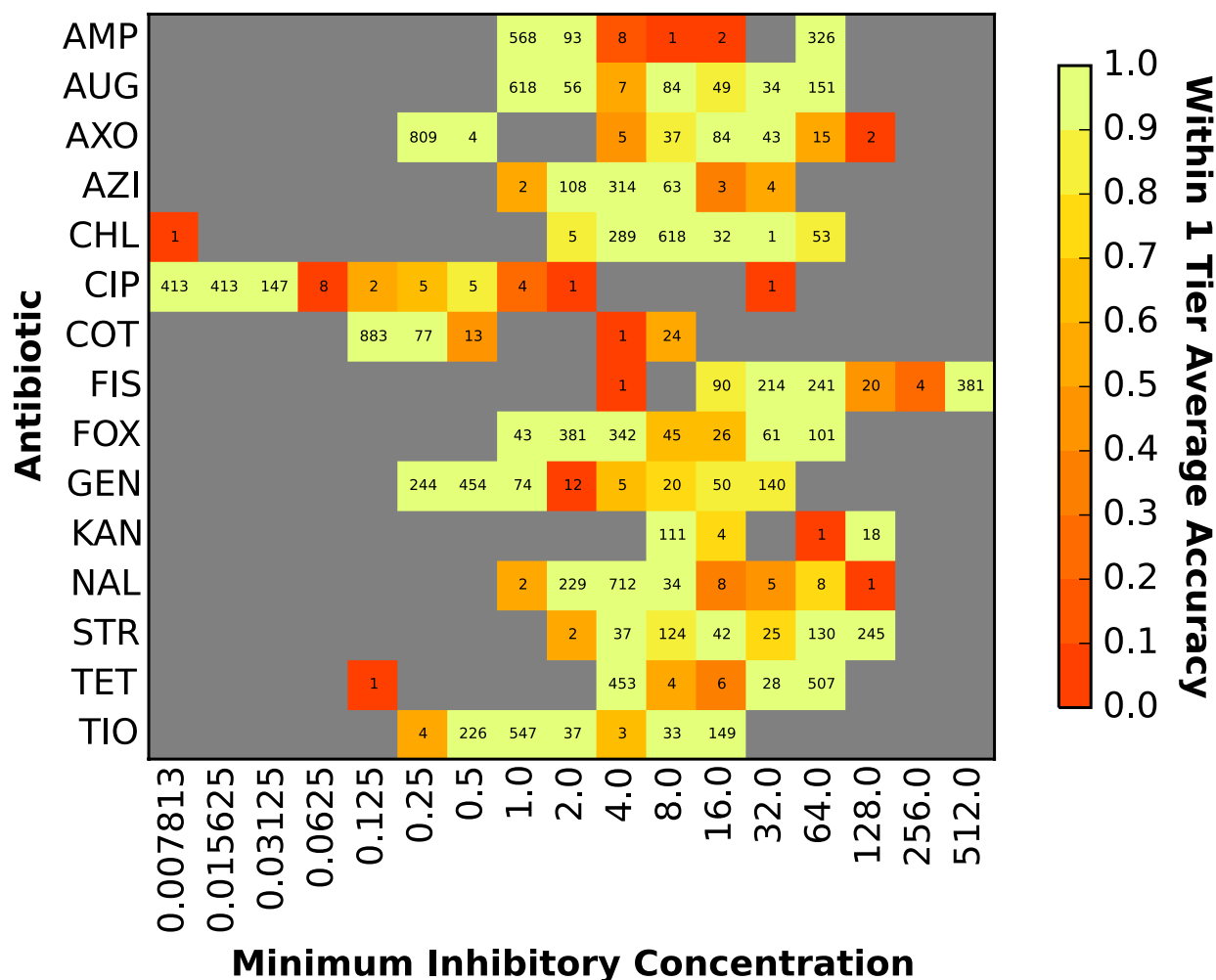
**Figure S3.** The accuracy of each single antibiotic model based on the set of ≤ 1000 diverse genomes. The heat map depicts the accuracy within ±1 two-fold dilution step of the laboratory-derived MIC. The X-axis shows the MIC (µg/ml) and each antibiotic is shown on the Y-axis. The accuracy for each antibiotic-MIC combination is depicted by color with bright yellow/green being the most accurate and red being the least accurate. The values shown in each cell are the number of genomes with that MIC for a given antibiotic.
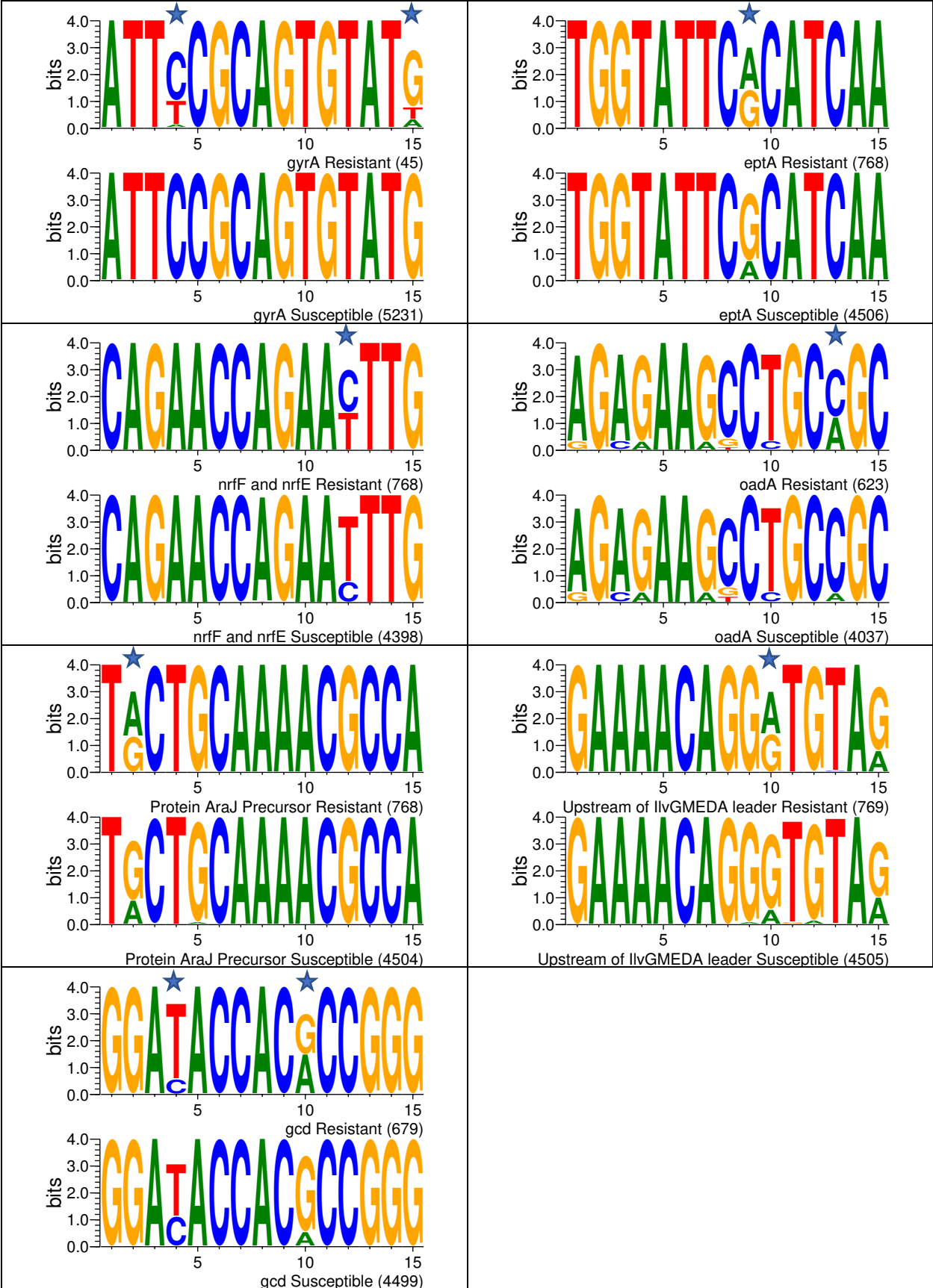
gyrA Resistant (45)

gyrA Susceptible (5231)

eptA Resistant (768)

eptA Susceptible (4506)

nrfF and nrfE Resistant (768)

nrfF and nrfE Susceptible (4398)

oadA Resistant (623)

oadA Susceptible (4037)

Protein AraJ Precursor Resistant (768)

Protein AraJ Precursor Susceptible (4504)

Upstream of IlvGMEDA leader Resistant (769)

Upstream of IlvGMEDA leader Susceptible (4505)

gcd Resistant (679)

gcd Susceptible (4499)

**Figure S4.** WebLogo plots for the susceptibility-related k-mer regions containing significant SNPs from Table 7 (P < 0.001, based on a Chi-square test). SNPs are significant in both the set of 1,000 diverse genomes and in the full set of 5,278 genomes (data are depicted for the full set of 5,278 genomes). Logos for the k-mer region in the resistant genomes is depicted above the logos for the k-mer region from the susceptible genomes. Significant SNPs are depicted with a star. Gene names and short hand descriptions correspond to the following PATRIC annotations: gyrA, DNA gyrase subunit A (EC 5.99.1.3); eptA, Phosphoethanolamine transferase EptA; nrfE and nrfF, Formate-dependent nitrite reductase complex subunit NrfF and Cytochrome c-type heme lyase subunit nrfE, nitrite reductase complex assembly (the k-mer occurs in both genes); oadA, Oxaloacetate decarboxylase alpha chain (EC 4.1.1.3); gcd, Glucose dehydrogenase, PQQ-dependent (EC 1.1.5.2).