

APPENDIX

A ANALYSIS OF THE uGLAD ARCHITECTURE

A.1 uGLAD PARAMETER DETAILS

For the sake of completeness of understanding the GLAD architecture, we graciously borrow the algorithm (see Alg.2) and neural network design details from Shrivastava et al. (2020b) here. GLAD parameter settings details are: ρ_{nn} was a 4 layer neural network and Λ_{nn} was a 2 layer neural network. Both used 3 hidden units in each layer. The non-linearity used for hidden layers was \tanh , while the final layer had sigmoid (σ) as the non-linearity for both, ρ_{nn} and Λ_{nn} (refer to Figure 2). The learnable offset parameter of initial Θ_0 was set to $t = 1$. It was unrolled for $L = 30$ iterations. The optimizer used was adam with the learning rates were chosen between $[0.001, 0.005]$.

B CASE STUDY: ANAEROBIC DIGESTION

Our algorithm development was inspired by a practical problem of domain exploration in anaerobic digestion. Anaerobic digestion is a growing field addressing waste management with both environmental benefits (reduced odor and pathogens, improved soil health, reduction in methane emissions) and economic value from use of captured methane gas. Despite numerous studies, the dynamics of organisms' growth in digesters, their dependence on conditions (temperature, pH, feedstock mix, nitrogen to carbon ratio, etc.) and their impact on methane yield are not well understood.

We present findings based on a public dataset from a study of anaerobic digesters at Danish wastewater plants Jiang et al. (2021). Data is available at NCBI under bioproject accession number PRJNA637463.

Data comes from a 6-year study of 46 digesters located at 22 Danish treatment plants. We have three types of digesters, operating at different temperatures (mesophilic, mesophilic with thermal hydrolysis pre-treatment, thermophilic). Digesters operate continuously with sludge retention rate of 15.8 to 35.6 days. Samples were taken at 3-month and 6-month intervals, so they can be treated as i.i.d. We have a total of 1,010 sludge samples, 418 used to sequence archaea and 592 bacteria, performed using 16S rRNA gene amplicon sequencing. Analysis resulted in identification of 33,047 bacterial and 878 archaeal unique amplicon sequence variants (ASVs). 70% of genera and 93% of the species were determined to be novel or previously unclassified. This presents problems for all approaches attempting to utilize existing databases to determine organisms' function for the purpose of grouping and feature selection. In fact, one of the best ways to determine an organism's function is based on checking properties of organisms whose abundance numbers in the digester best correlate with the given organism's numbers.

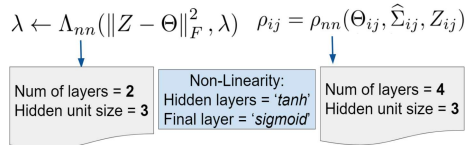
Algorithm 2: GLAD

Function GLADcell($\widehat{\Sigma}, \Theta, Z, \lambda$):

- $\lambda \leftarrow \Lambda_{nn}(\|Z - \Theta\|_F^2, \lambda)$
- $Y \leftarrow \lambda^{-1} \widehat{\Sigma} - Z$
- $\Theta \leftarrow \frac{1}{2} \left(-Y + \sqrt{Y^\top Y + \frac{4}{\lambda} I} \right)$
- For all** i, j **do**
 - $\rho_{ij} = \rho_{nn}(\Theta_{ij}, \widehat{\Sigma}_{ij}, Z_{ij})$
 - $Z_{ij} \leftarrow \eta_{\rho_{ij}}(\Theta_{ij})$
- return** Θ, Z, λ

Function GLAD($\widehat{\Sigma}$):

- $\Theta_0 \leftarrow (\widehat{\Sigma} + tI)^{-1}, \lambda_0 \leftarrow 1$
- For** $k = 0$ **to** $K - 1$ **do**
 - $\Theta_{k+1}, Z_{k+1}, \lambda_{k+1} \leftarrow \text{GLADcell}(\widehat{\Sigma}, \Theta_k, Z_k, \lambda_k)$
- return** Θ_K, Z_K

Figure 2: **Minimalist** neural network architectures designed for GLAD

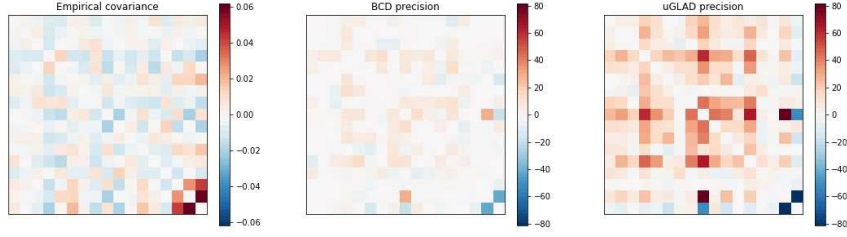


Figure 3: uGLAD recovered precision matrix compared to empirical covariance and precision matrix recovered by the BCD algorithm for archaea at **family** level

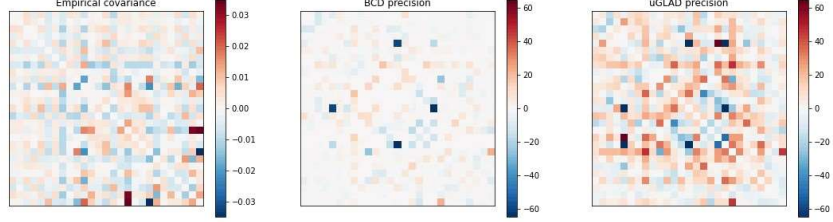


Figure 4: uGLAD recovered precision matrix compared to empirical covariance and precision matrix recovered by the BCD algorithm for archaea at **species** level

Our algorithm works with any input, including: ASVs filtered by frequency, ASVs rolled up to higher taxonomy levels (species, genus, family), ASVs abundance normalized in various ways Badri et al. (2020). We calculate the partial correlation matrix from the precision matrix. Each entry of the partial correlation matrix P_{ij} shows the correlation of the feature x_i, x_j given the values of the other features are observed. This helps us obtain the direct dependence of the features.

$$P_{ij} = -\frac{\Theta_{ij}}{\sqrt{\Theta_{ii}\Theta_{jj}}} \quad (9)$$

We use networkx package to visualize the graphs, presenting positive correlations in green and negative in red, with edge weights corresponding to the strength of the correlation.

Figures 4 and 3 present precision matrices recovered by our algorithm and BCD with empirical covariance shown for comparison. Figures 5 and 6 show corresponding graphs for archaea at family and species level.

Figures 8 and 7 show a result of multitask learning based on digester type: mesophilic (operating at temperature 38°C) and thermophilic (operated at temperature 53.6°C). The two graphs' edges are filtered to show only edges common to both graphs, which is a small fraction of all edges. Note that in some cases, the sign of the correlation (and the color of the edge) changes depending on digester's type.

Our model is being used by domain experts to gain insight into the domain of anaerobic digestion. One use case is to understand properties of newly discovered bacteria and archaea by analyzing which known organisms their abundance in digester samples correlates with (positively or negatively). That can lead to focusing attention on a smaller organism set. Another use case centers around understanding the role of digester conditions and feedstock mix on organisms' growth and methane yield. The results presented in recovered graphs lead to new hypotheses and new experiments being designed to test them.

C CASE STUDY: INFANT MORTALITY

We used uGLAD to recover the graph for the domain of infant mortality. The dataset is based on CDC Birth Cohort Linked Birth – Infant Death Data Files of Health et al.. It describes pregnancy and birth variables for all live births in the U.S. together with an indication of an infant's death before the first birthday. We used the data for 2015 (latest available), which includes information about 3,988,733 live births in the US during 2015 calendar year.

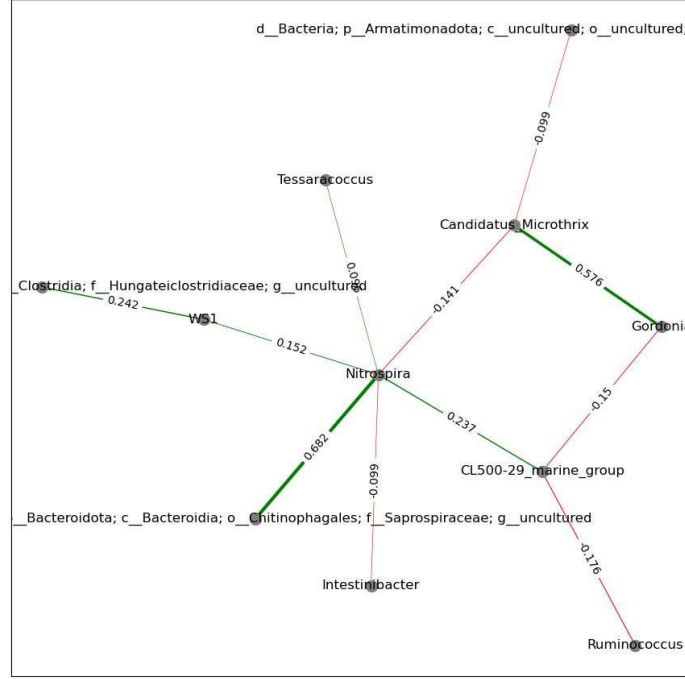


Figure 7: Example of **multitask** learning applied to different digester types. uGLAD graph for bacteria at genus level for thermophilic digesters showing only edges common to all digester types. Note that edge colors (correlation signs) are different from the corresponding graph for mesophilic digesters. Edge color indicates the sign of the correlation: green - positive, red - negative, edge weight corresponds to correlation's strength.

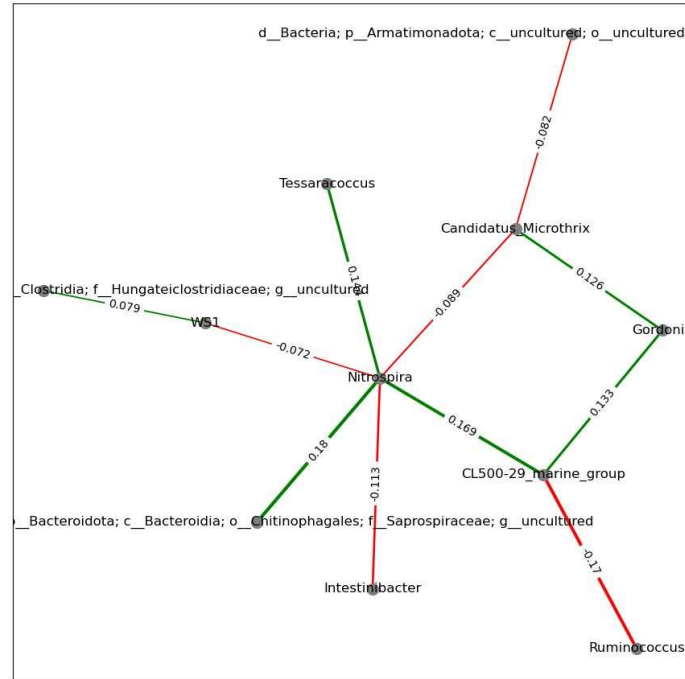


Figure 8: Example of **multitask** learning applied to different digester types. uGLAD graph for bacteria at genus level for mesophilic digesters showing only edges common to all digester types. Edge color indicates the sign of the correlation: green - positive, red - negative, edge weight corresponds to correlation's strength.

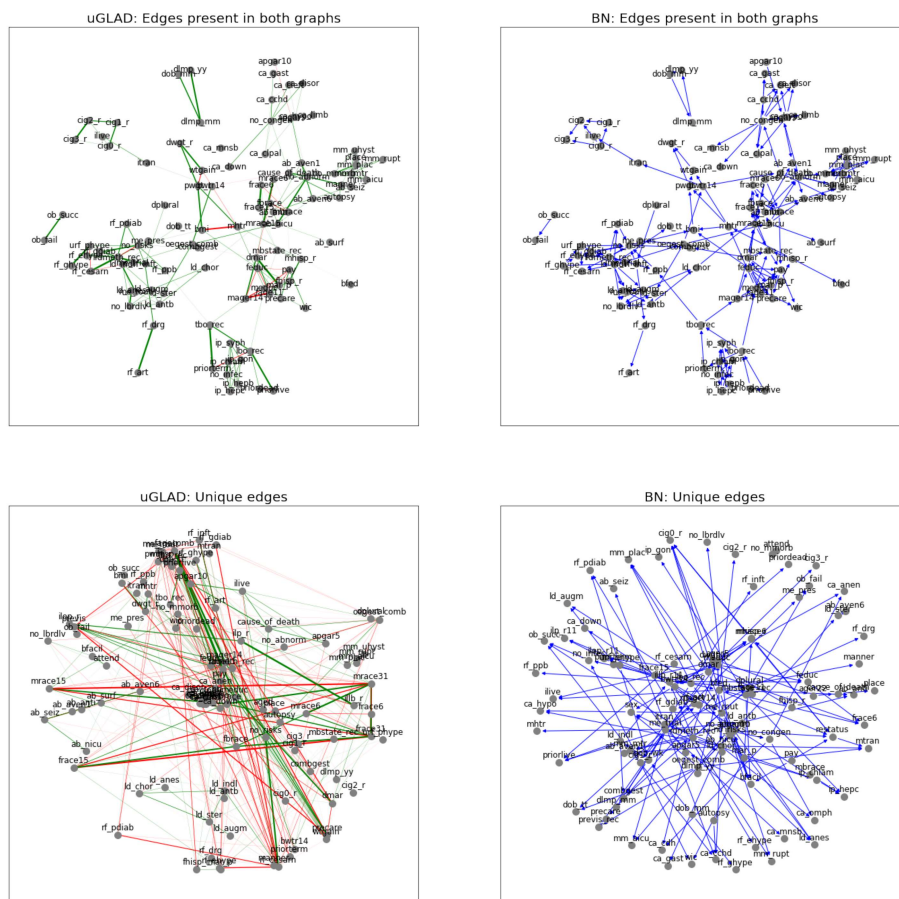


Figure 11: Comparing the graphs recovered by uGLAD and Bayesian Network recovery package (Scutari, 2010).