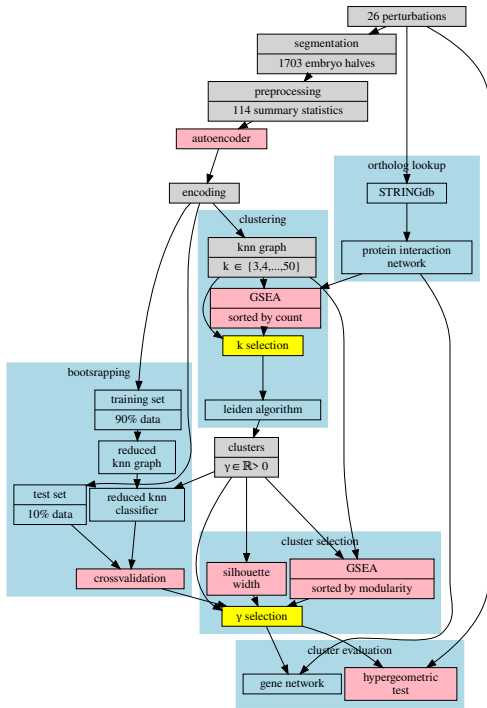


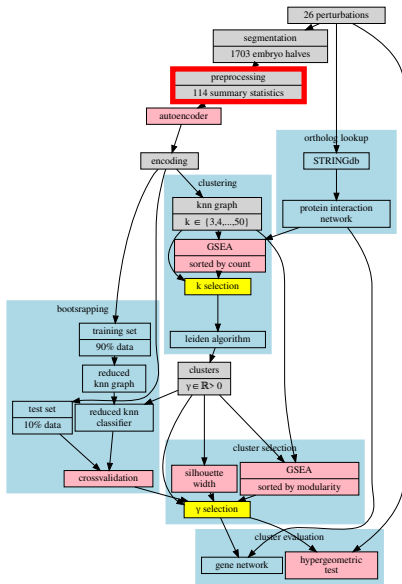
CRISPR Screen

Unsupervised Clustering with Automated Hyperparameter Selection

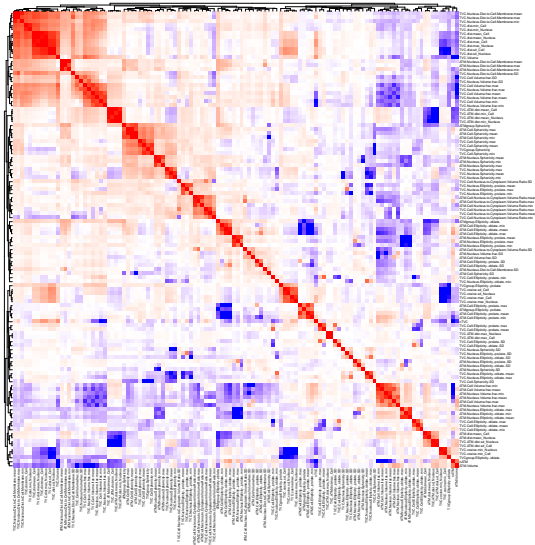
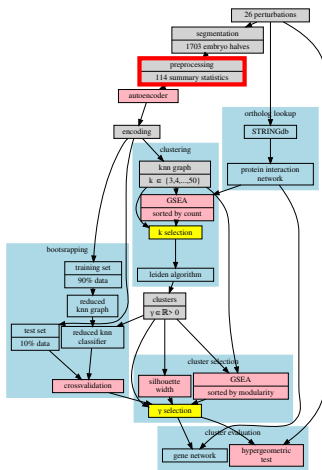
Keira Wiechecki

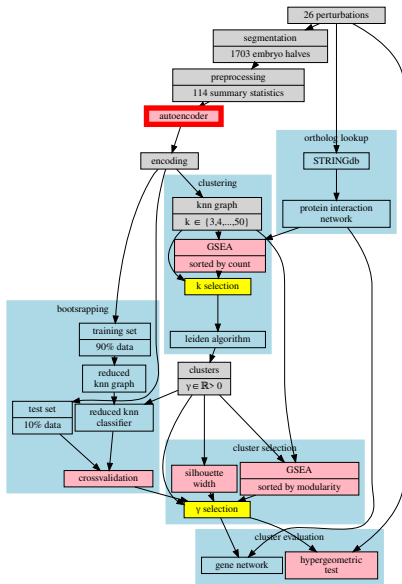
October 23, 2022



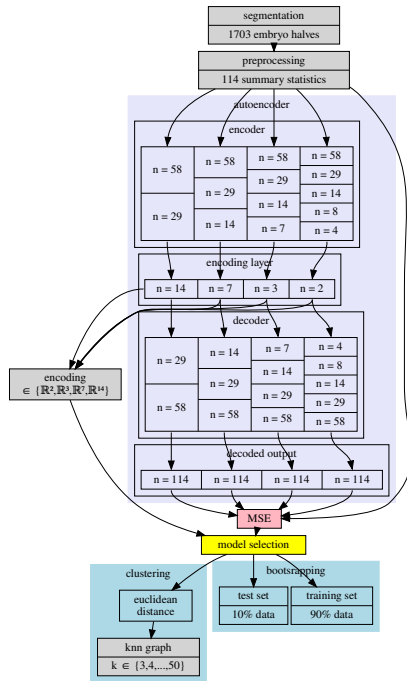
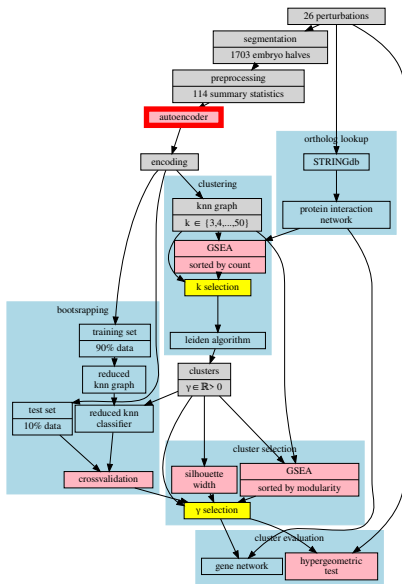


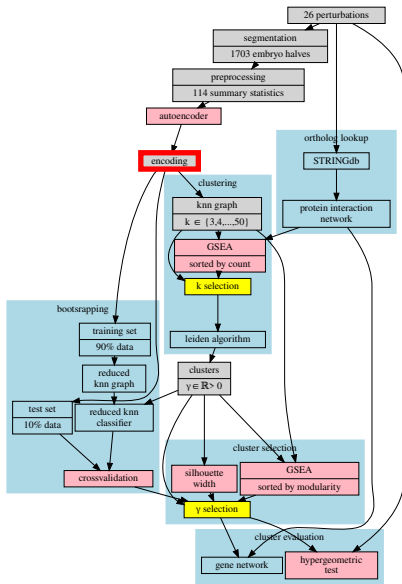
Many parameters are strongly correlated. This is undesirable because each parameter additively contributes to distance used for clustering, resulting in disproportionate weight being given to phenotypes captured by multiple parameters.



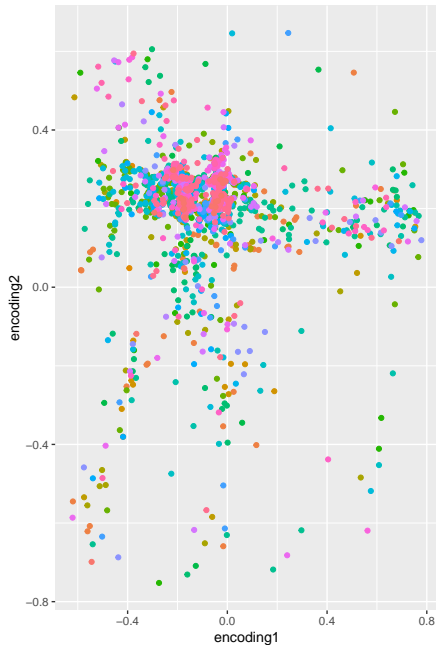
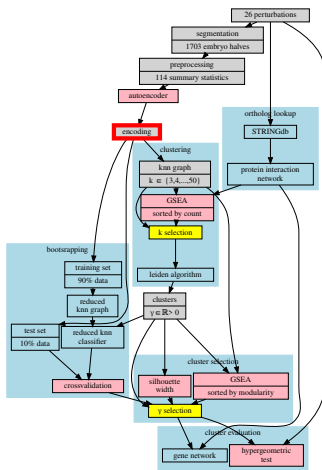


An autoencoder is a method of dimension reduction that uses a neural network to find a lower dimensional encoding which can be decoded to recover the input. This reduces exaggeration of distance due to the number of parameters measured.

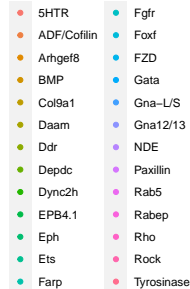


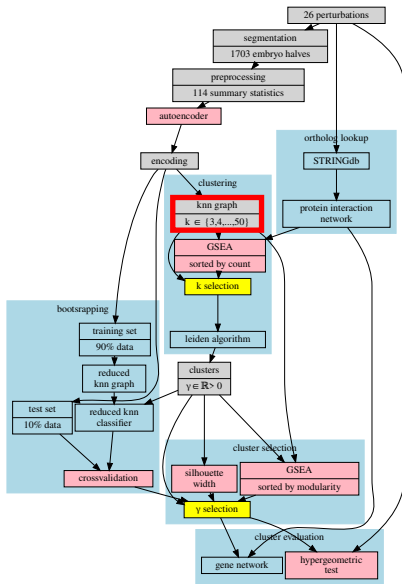


More than 2D produced only marginal improvement.

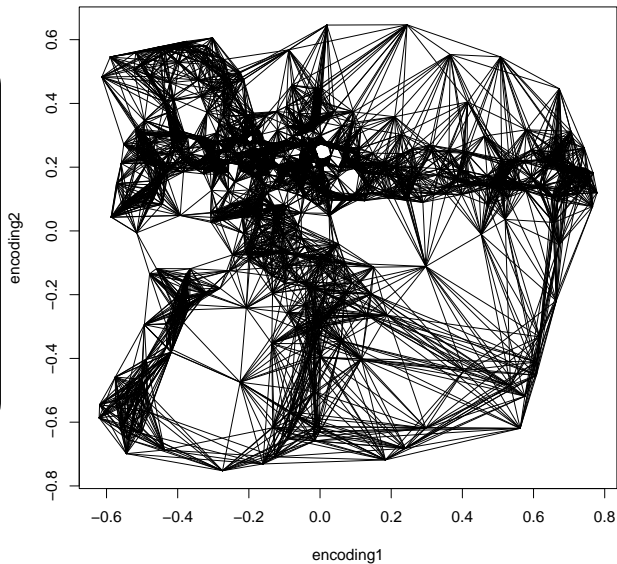
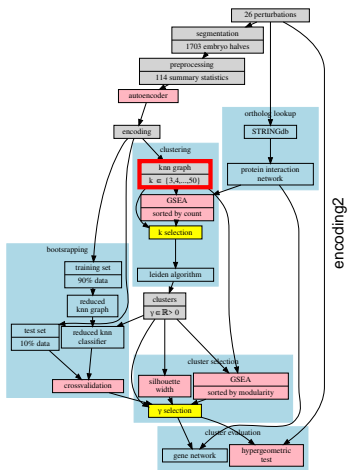


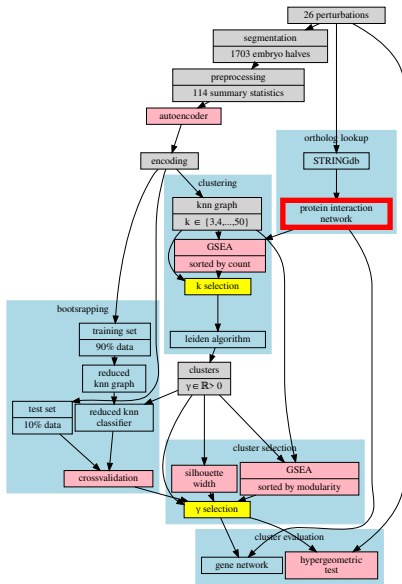
condition



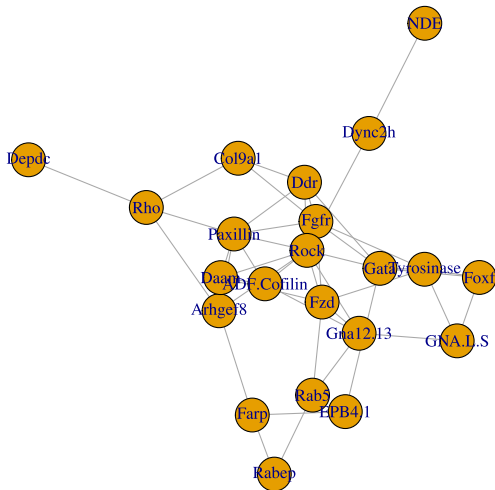
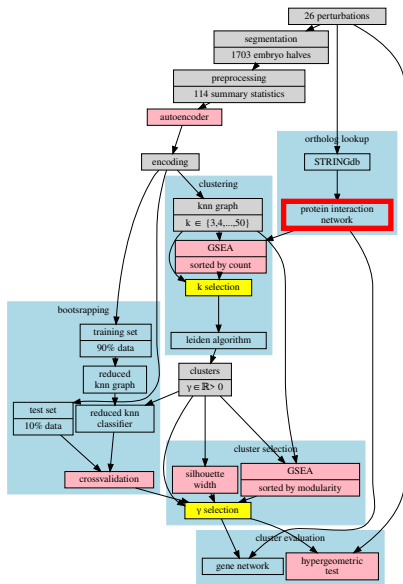


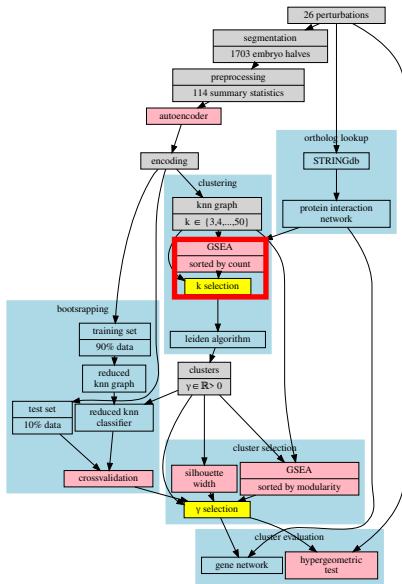
A kNN graph is constructed from euclidean distance calculated from the encoding layer.



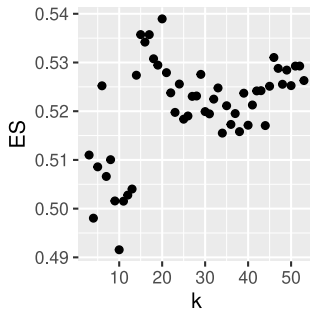
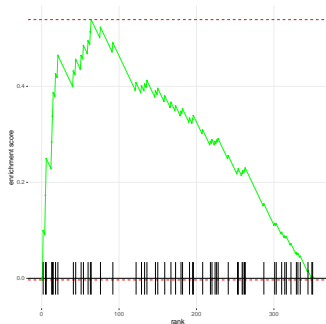
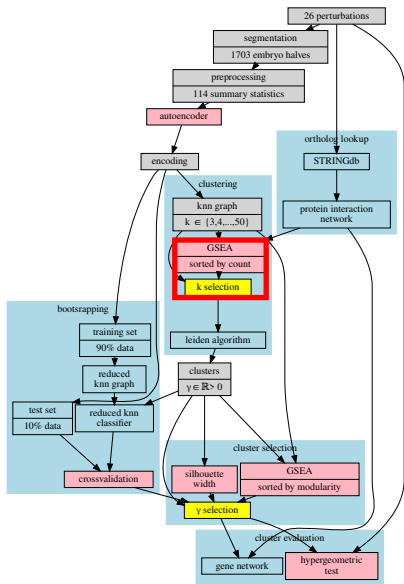


I obtained orthologs from *M. musculus* and *H. sapiens* from ENSEMBL. I used STRINGdb to obtain known protein interactions.





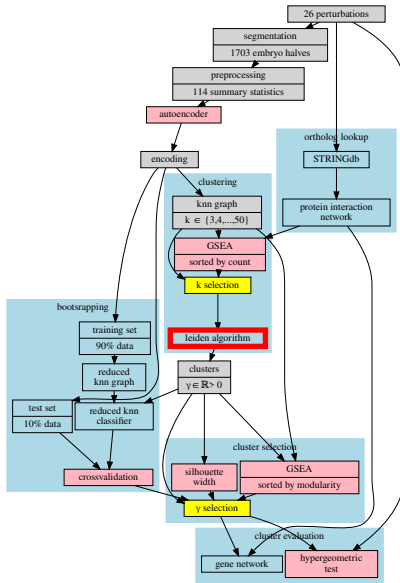
The known protein interactions can be treated as a gene set for GSEA. Interactions can be ranked by edge count between embryos in two conditions. An enrichment score is calculated based on occurrence of known interactions near the top of the ranked list. An optimal k can be selected by maximizing enrichment score.

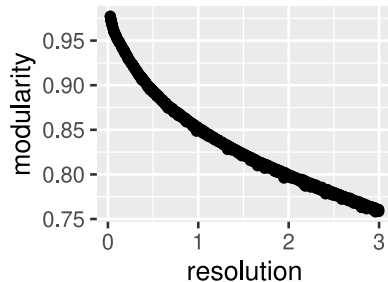
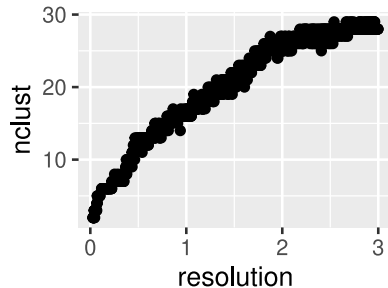
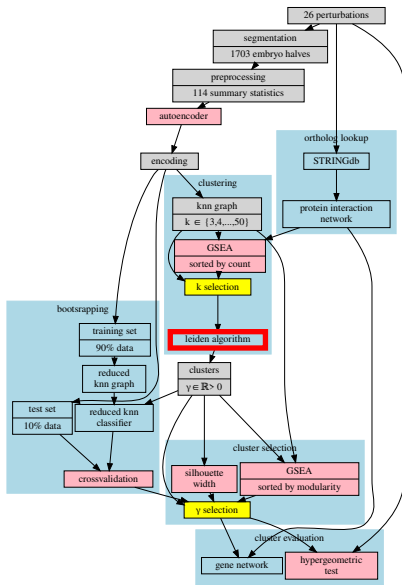


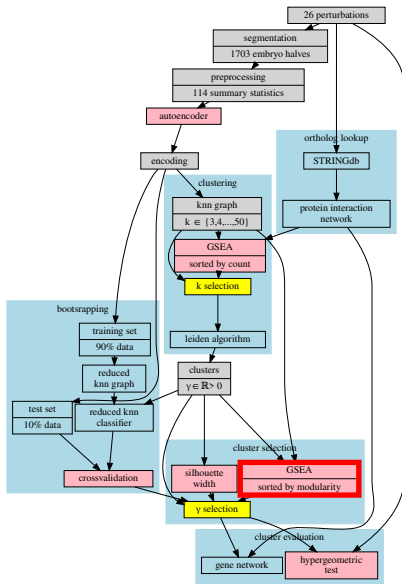
The leiden algorithm attempts to find a clustering that maximizes modularity H for a given graph and resolution γ . Modularity is defined as

$$H = \frac{1}{2m} \sum_c (e_c - \gamma \frac{K_c^2}{2m})$$

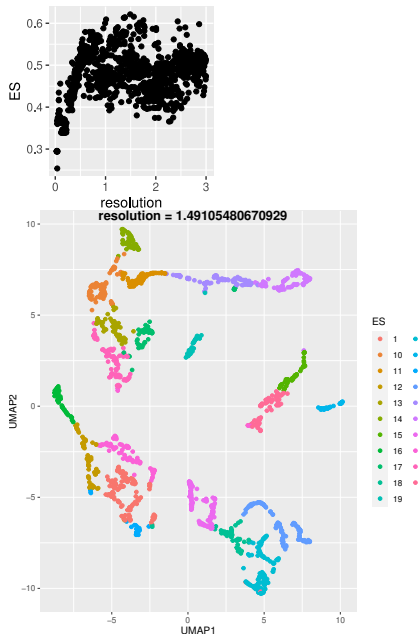
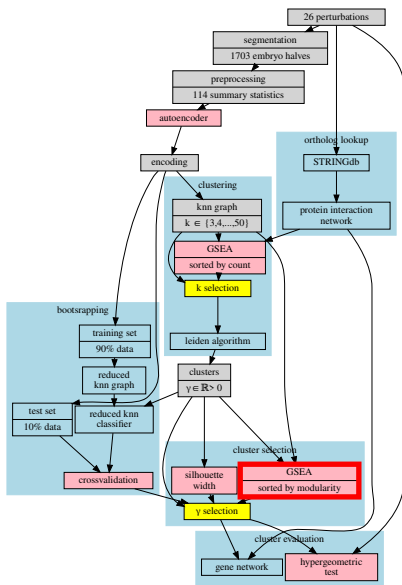
where m is the average degree of the graph, e_c is the number of edges in cluster c , and K_c is the number of nodes in cluster c . This gives a measure of how well-connected clusters compared to expectation based on average degree of the graph and number of nodes in a cluster. A higher γ results in more clusters.

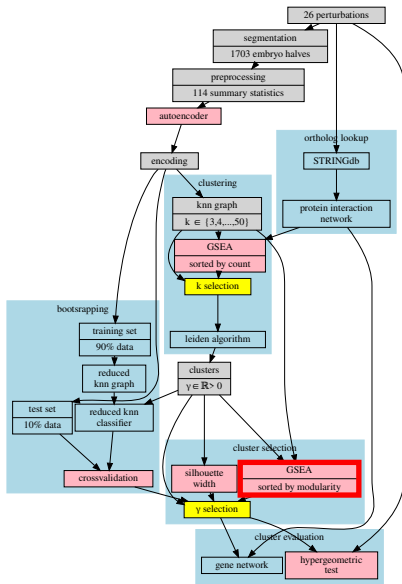




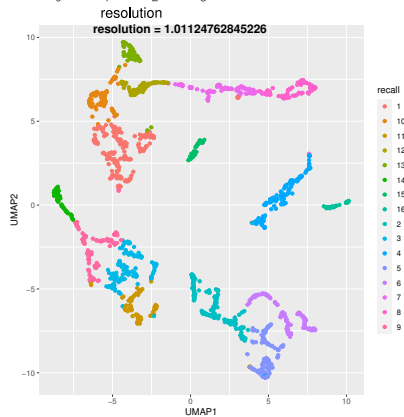
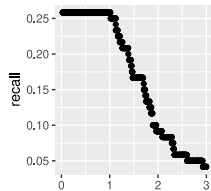
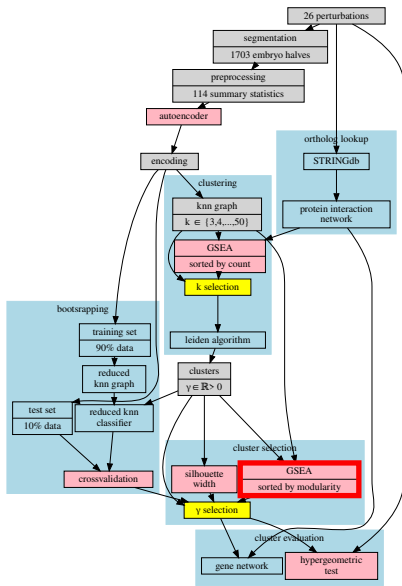


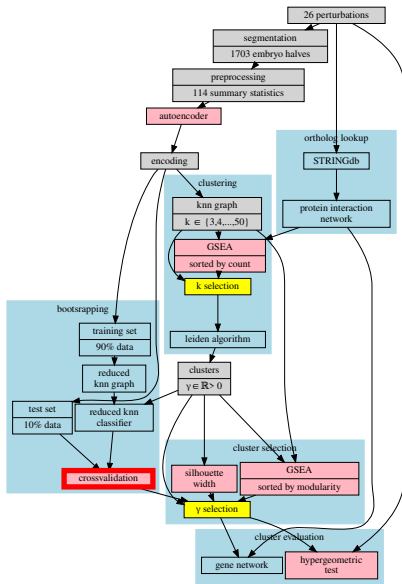
I obtained 1000 clusterings on random γ values between 0.01 and 3. I assessed clustering by several metrics. I obtained an enrichment score from a second GSEA with interactions ranked by modularity of a subgraph consisting of only edges between two given conditions.





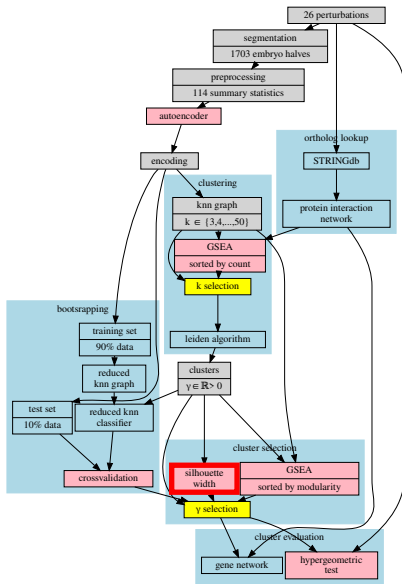
I used this same metric to construct a gene network. If an interaction has positive modularity the genes are considered connected. A recall score can be calculated by comparing this network to the protein interaction network.



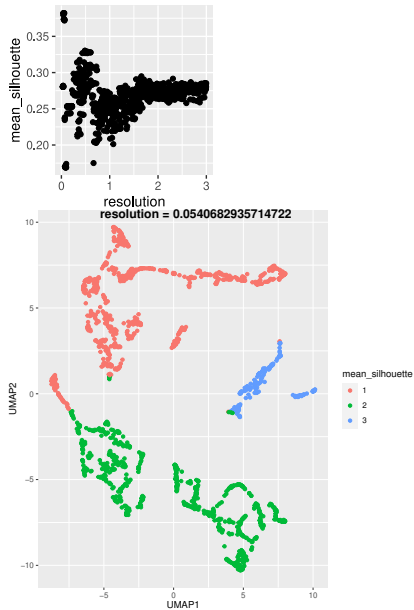
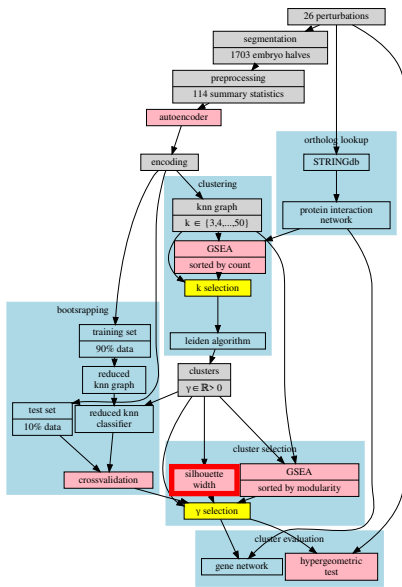


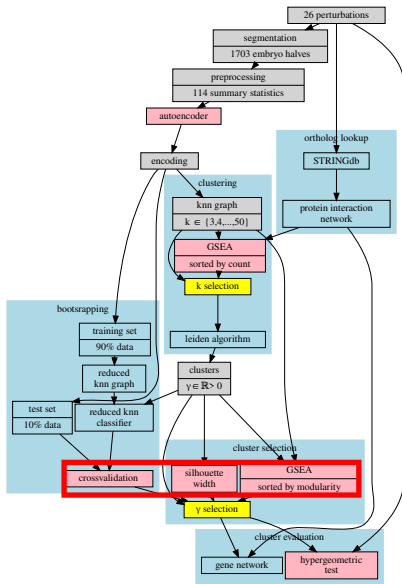
I used this same metric to construct a gene network. If an interaction has positive modularity the genes are considered connected. A recall score can be calculated by comparing this network to the protein interaction network.

— $\log_2 \text{error}$ is calculated by building a kNN classifier for the clusters from a subset of the data and repeating for 1000 permutations.

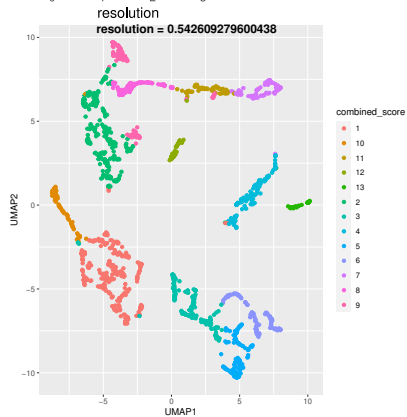
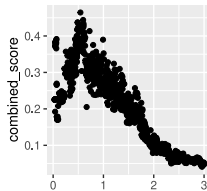
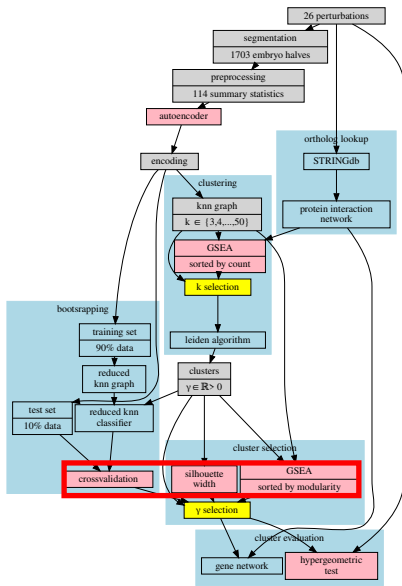


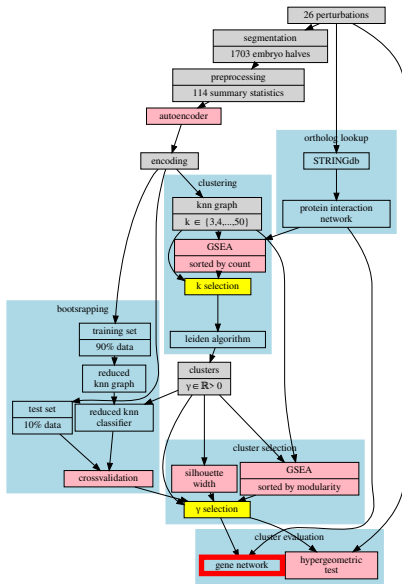
Silhouette score gives the relative distance between points within a cluster compared to distance between points in different clusters.



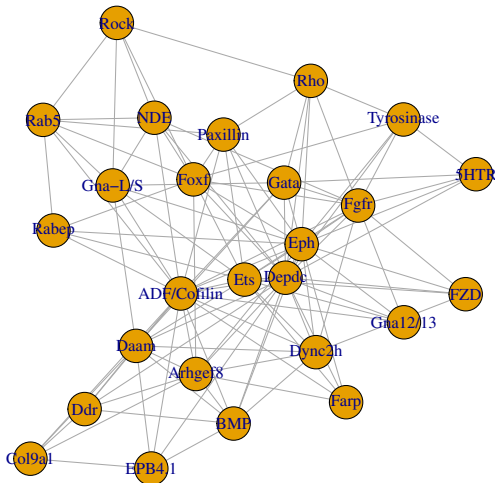
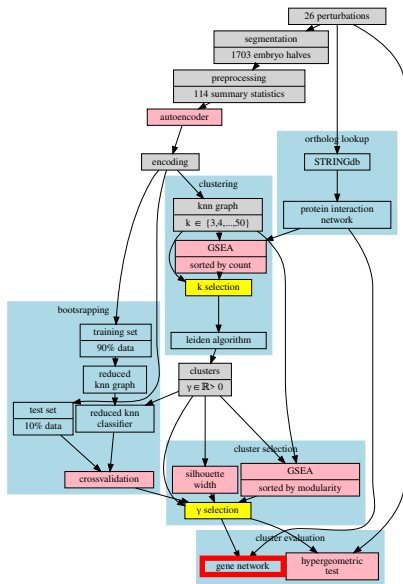


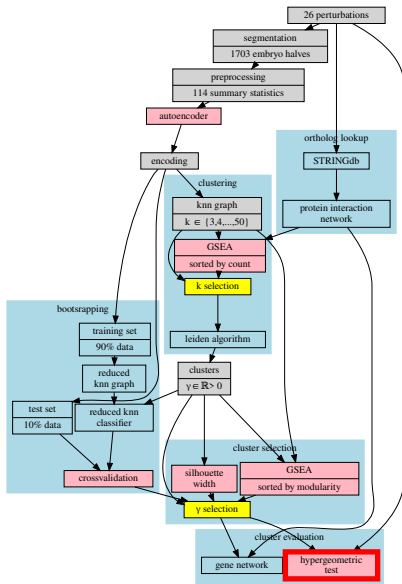
Because there was no clear best metric, I calculated a combined score from $ES \times recall \times silhouettewidth \times -\log_2 error$.





I constructed a network using the top 5 modularity values for each gene.





For each cluster I performed hypergeometric tests for enrichment of embryos in each treatment and enrichment of experimenter-identified phenotypes.

