



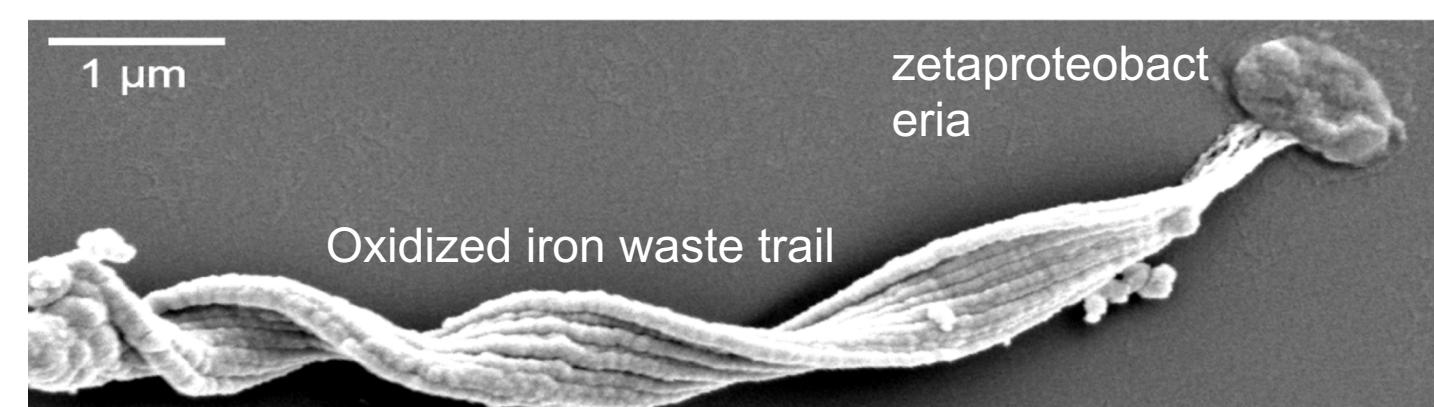
Clustering Temporal Gene Expressions of Iron-Oxidizing Zetaproteobacteria

Students: Stephen Herbein and Sean McAllister

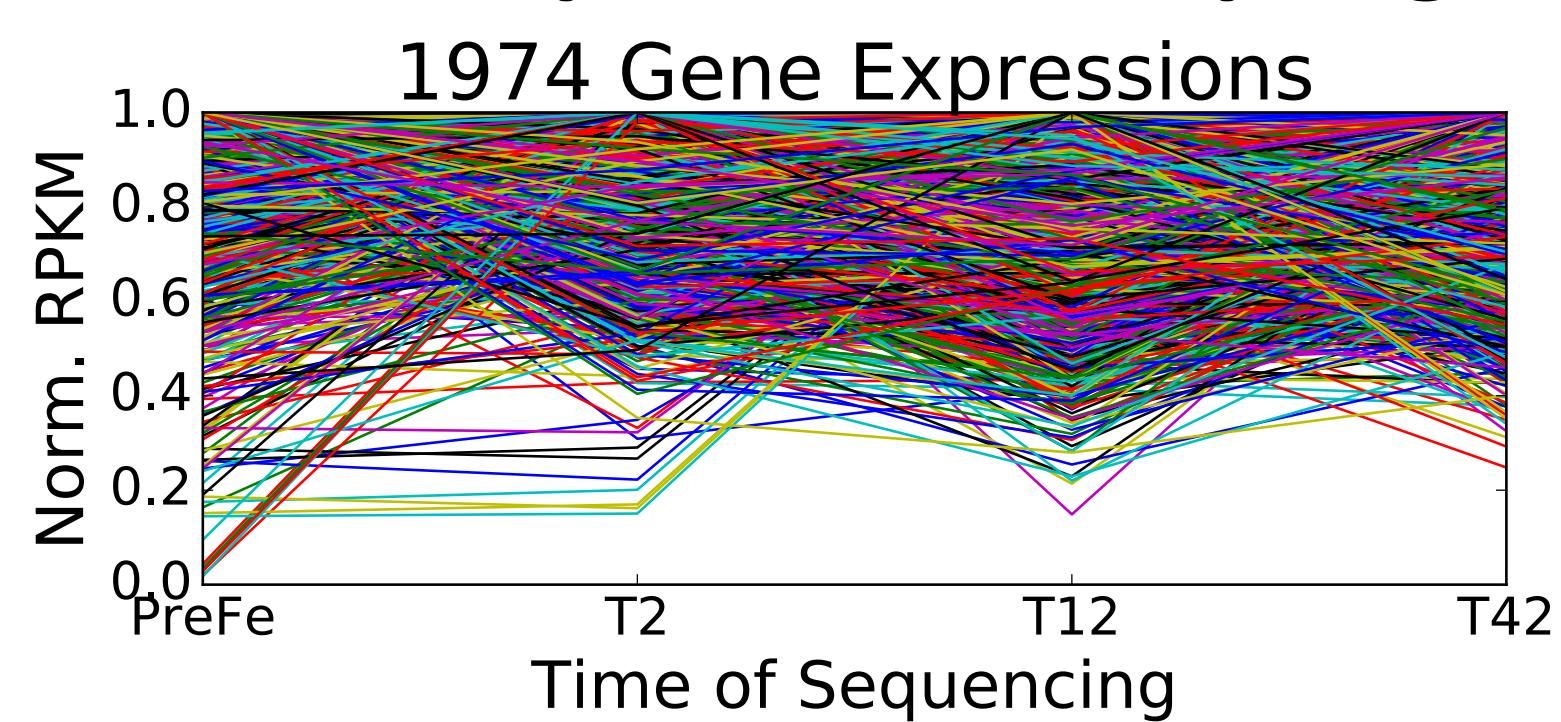
Advisors: Michela Taufer and Clara Chan

Zetaproteobacteria in a Nutshell

- Zetaproteobacteria are an iron-oxidizing bacteria commonly found at deep-sea hydrothermal vents
- The bacteria play a big role in the rusting of ship hulls, metal pilings, and pipelines
- Scientists want to better understand the genes involved in the oxidation of iron



Gene Expression Sampling



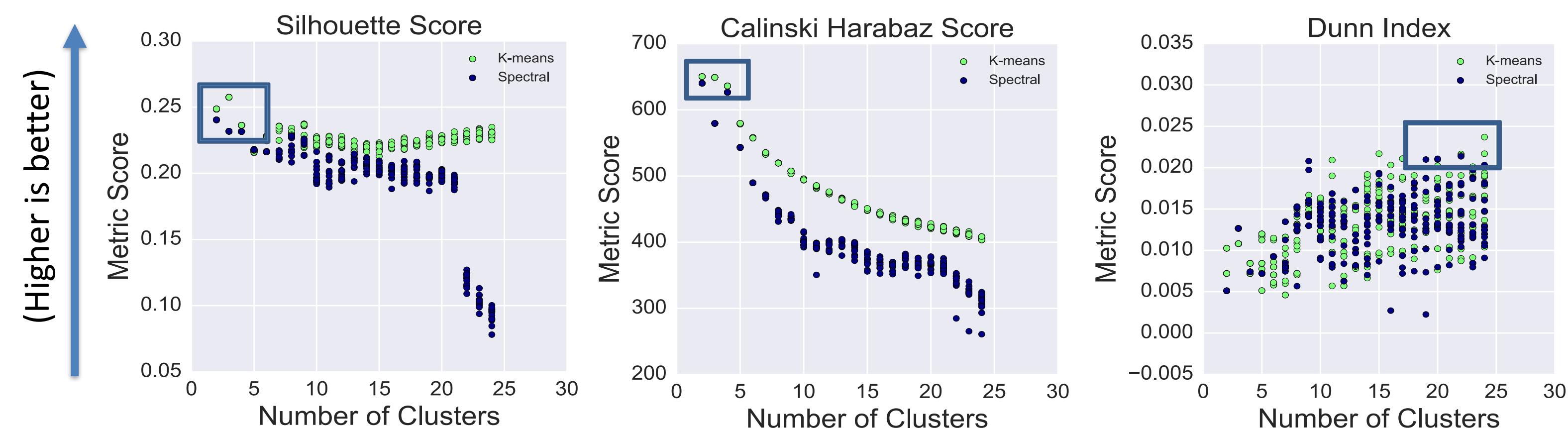
Can we find patterns in the gene expressions to determine which genes are involved in iron oxidation?

Clustering Methods

	K-means	Spectral
Use-cases	General-purpose, even cluster size, flat geometry	Few clusters, even cluster size, non-flat geometry
Scalability	Very large # of samples, medium # clusters	Medium # samples, small # clusters
Distance Metric	Euclidian Distance	Graph Distance
Parameters	Number of clusters	Number of clusters

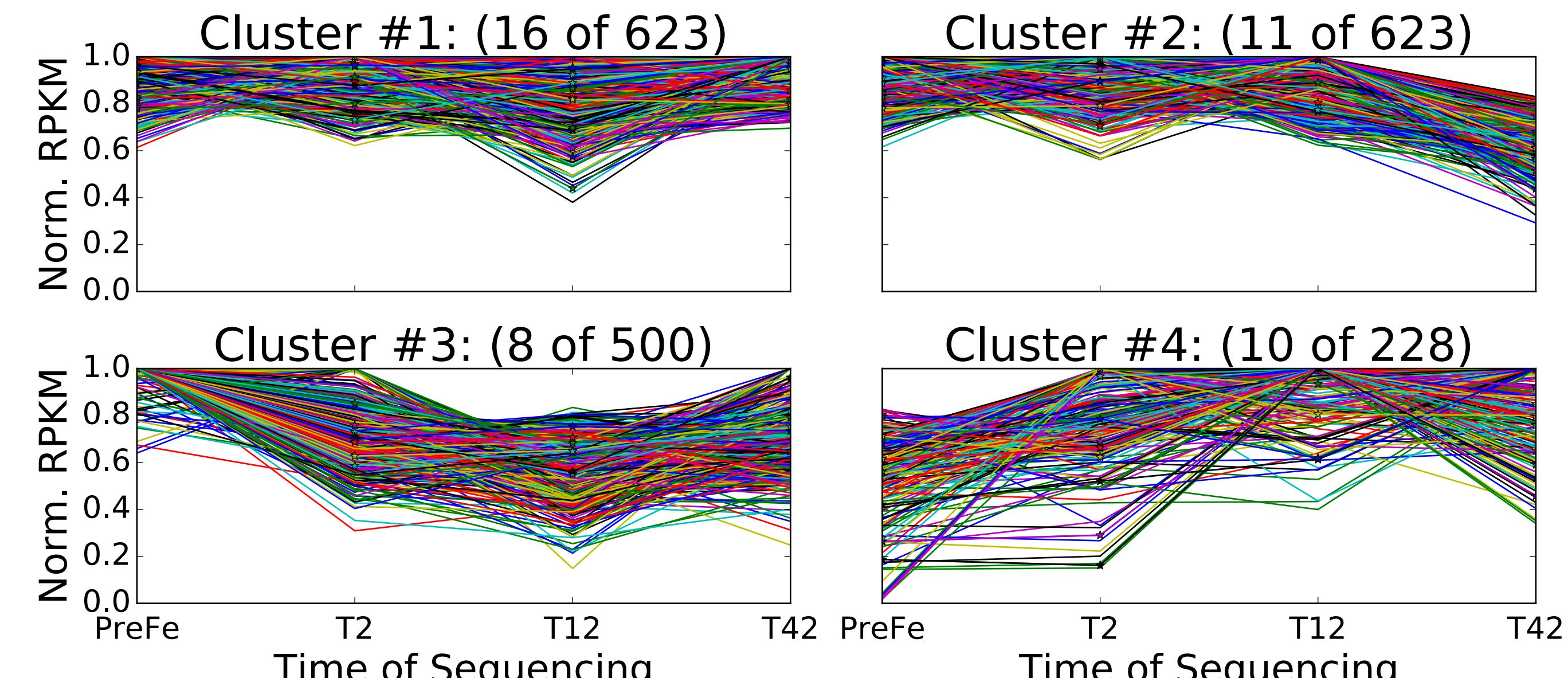
How do we determine the optimal number of clusters?

Determining the Optimal Number of Clusters



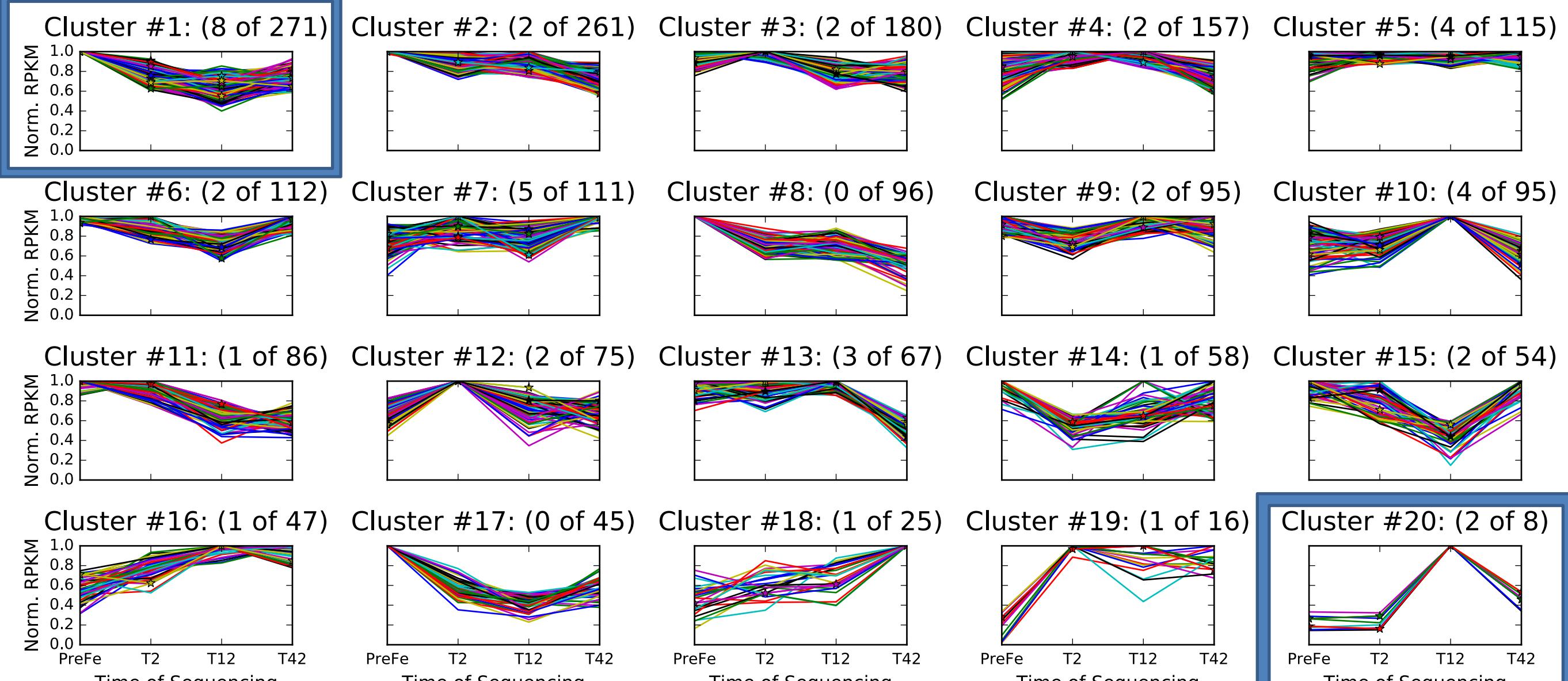
The metrics suggest that using either 4 or 20 clusters is optimal

4 Clusters Under Spectral Clustering



- Results of 4 clusters with k-means are very similar to Spectral

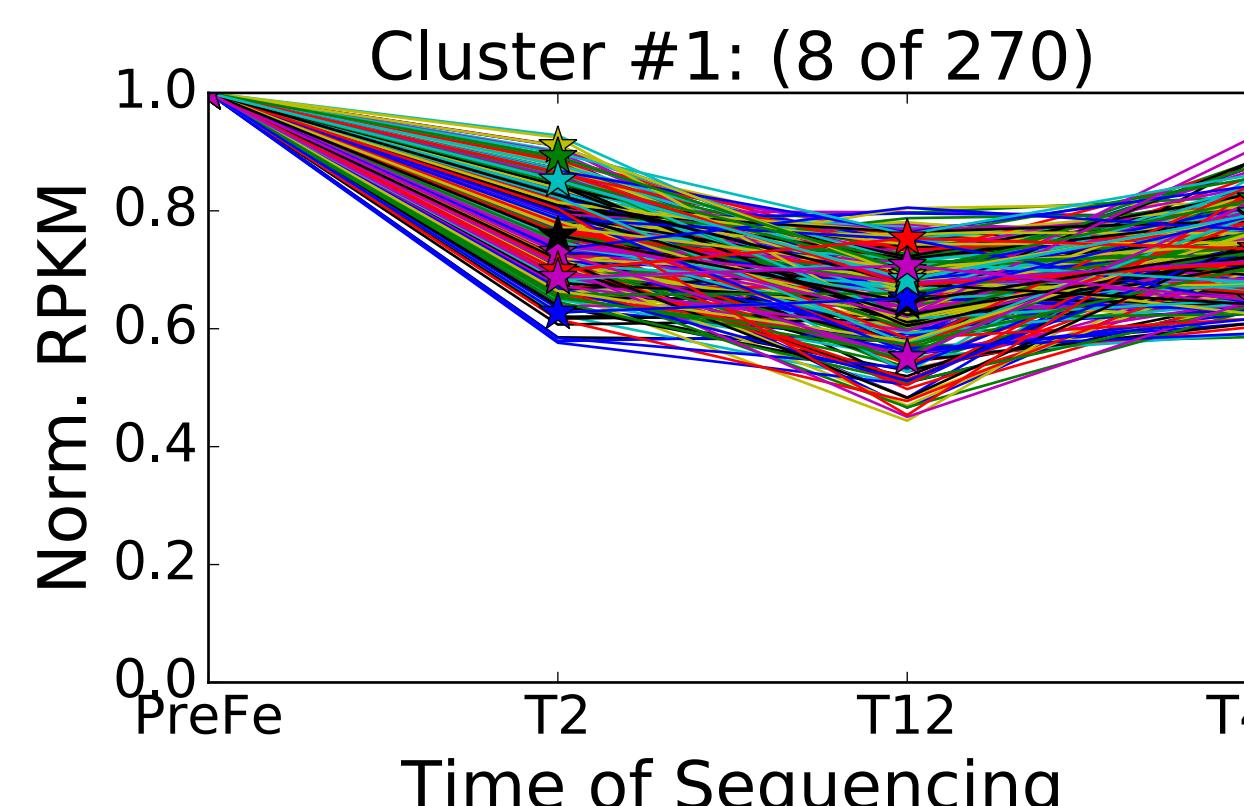
20 Clusters Under Spectral Clustering



- Results of 20 clusters with k-means are visually noisier than Spectral

Conclusions

- 4 clusters result in high metric scores, but the clusters are too large to be useful for the next steps in scientific inquiry
- 20 clusters provide a finer level of detail and help narrow the search down to a smaller number of genes
- Cluster #1 has a large number of genes that were pre-selected by the scientists as interesting
- Cluster #20 has a high percentage of interesting genes



The optimal number of clusters depends on the scientific objective

