



Verifying the Role of RosR Transcription Factor in the Oxidative Stress Response of *Halobacterium Salinarum* NRC-1

Yitaek Hwang, Tosin Omofoye, Steven Pierre, Zhihao Zhu



Abstract

Previous work by Sharma et al. and Tonner et al. have shown that VNG0258H or RosR is a transcription factor that regulates gene expression in *Halobacterium Salinarum* NRC-1 under extreme oxidative stress. The goal of this project was to verify the results of the aforementioned papers and delve deeper into identifying specific genes that RosR either activates or represses in response to oxidative stress. This allows to map the genes whose expression are directly affected by RosR. In comparing the gene expression ratio in the *Dura3* parent and Δ VNG0258H mutant strains before and after exposing each subject to H₂O₂, we validated six expression profiled noted in Sharma et al. Then, we used k-means clustering to find a list of genes regulated by RosR using the same gene expression profile. This list, however, differed significantly from the list provided in Sharma et al., leading us to question our stringent criteria for ChIP-chip analysis. Next, we created a transcription factor network to find the relationship between the duration of H₂O₂ exposure to the genes that RosR regulates. The specific functions of these genes were then analyzed using the arCOG gene ontology analysis. Lastly, we compared our results with those of Tonner et al. to validate our analysis. In conclusion, the analysis showed that there are very few genes that have very high correlation to RosR binding at every stage of response to extreme peroxide levels, but many genes do have high correlation to RosR binding immediately after exposure to extreme levels of oxidation.

Clustering Analysis

We used k-means clustering to start exploring the data as it is a simple and robust clustering algorithm for big datasets. To determine the appropriate number of k-clusters, we looked for an “elbow” in the sum of squared error (SSE) scree plot. The location of the “elbow” suggested 10 as a suitable number of clusters for the k-means analysis of parent strain dataset. A dendrogram was created based on these k-means clusters.

Heatmap based on peak expression levels of each gene was created for the parent strain as well as RosR deletion strands. The gene expression data values are in log₁₀ values. We exponentiated this dataset to get the nominal values and then found the differences between the maximum gene expression in the parent and RosR deletion strains. Genes with differences outside 2 standard deviations of the mean difference have significant change in gene expression due to RosR.

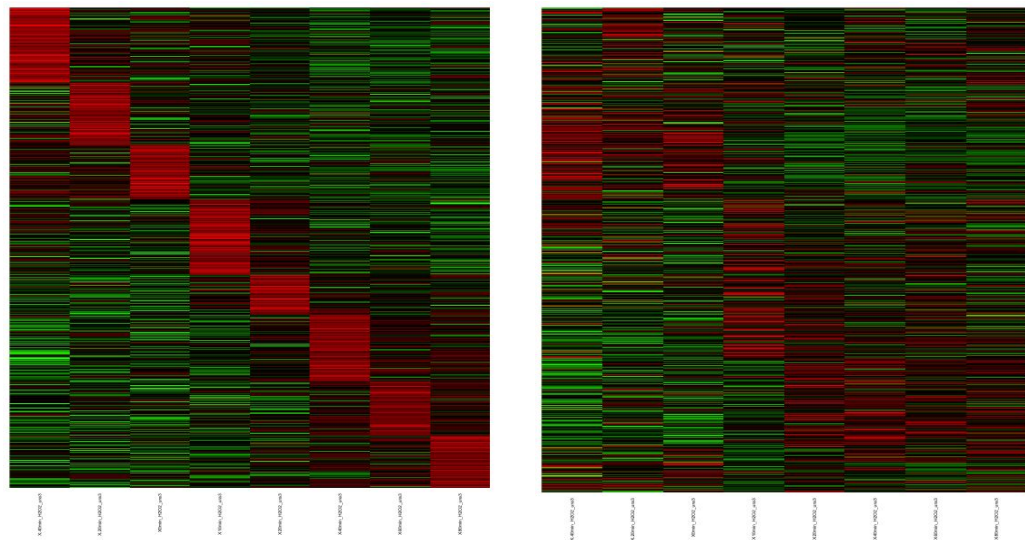


Figure 1: Maximum Gene Expression for Each Gene over the 8 Time Points

We computed the maximum gene expression over time for each gene in the data set and reordered based on peak expression profiles. The left heatmap in Figure 1 shows the peak gene expression (y-axis: genes, x-axis: time points) for the parent strain, whereas the right heatmap in Figure 1 shows the peak gene expression for the RosR deletion strain experiments.

Using all the differences between maximum expression of genes in the parent strain and RosR deletion strain experiments, we found 142 genes that exhibit significant differences in gene expression. The identified genes were then compared against the list in Sharma et al. paper first and then analyzed further using transcription factor network analysis and gene ontology.

VNG0218G	VNG0339H	VNG0349G	VNG0357H	VNG0462C	VNG0506H	VNG0509H	VNG0878G
176	261	273	352	388	391	678	
VNG0931G	VNG0932C	VNG0978H	VNG1087C	VNG1347C	VNG1366H	VNG1497C	VNG1547C
710	711	746	833	1029	1043	1143	1179
VNG1583C	VNG1674H	VNG1748C	VNG1872C	VNG1956H	VNG1997G	VNG2094G	VNG2164C
1206	1270	1323	1418	1483	1514	1585	1646
VNG2338G	VNG2440H	VNG2471G	VNG2472G	VNG2593H	VNG2594C	VNG5039G	VNG5059C
1788	1866	1890	1891	1979	1980	2074	2086
VNG5083H	VNG5090C	VNG5091C	VNG5146H	VNG5157H	VNG5160H	VNG5165H	VNG5167H
2104	2109	2110	2143	2151	2152	2155	2157
VNG6223C	VNG6284H	VNG6312G	VNG6323H	VNG6329H	VNG6330H	VNG6335H	VNG6339H
2250	2291	2312	2322	2326	2327	2330	2332
VNG6353H	VNG6393H	VNG6411H	VNG6416H	VNG6418H	VNG6424H	VNG6432H	
2341	2369	2381	2384	2385	2388	2393	

Replicating Sharma et al. Results

Gene expression in response to H₂O₂ exposure in the *Dura3* parent and Δ VNG0258H mutant strains are shown below. Like Figure 5 in Sharma et al., each line in the graph detail the log₂ gene expression ratio for time points before and after the H₂O₂ exposure noted by the dotted line:

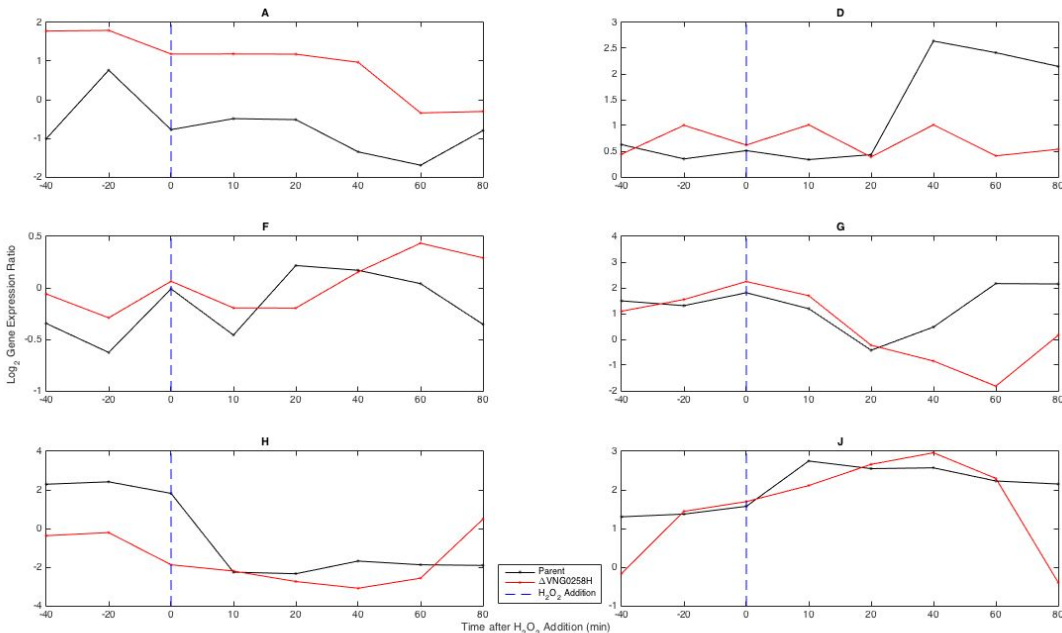


Figure 2: Gene Expression Profiles from Sharma et al.

(A) genes requiring VNG0258H for repression regardless of condition, (D) genes requiring VNG0258H for activation in the presence of H₂O₂, (F) genes requiring VNG0258H for impulse-like dynamic induction, (G) genes requiring VNG0258H for impulse-like dynamic repression, (H) genes induced in response to H₂O₂ but independent of VNG0258H, and (J) genes repressed in response to H₂O₂ but independent of VNG0258H. We wanted to see if the genes we found using the clustering analysis also were found in the expression profiles noted in Sharma et al. In this process, we were trying to validate whether or not our clustering results were due to the fact that we were missing data points.

TF Network Analysis

The ChIP-chip data provided detailed the interactions of a list of genes in response to the addition of H₂O₂ at various time points. We used a threshold of 0.001 on the p-value to find the genes with high probability of significant interaction with RosR. We chose this very small cutoff value as we wanted to focus on smaller group of genes that have very high relation to RosR and study their functions and interactions to produce the response to reactive oxygen stimulus. Next an incidence matrix was made using the data at the time points following addition of H₂O₂. The transcription factor network generated below shows the significant interaction between RosR and individual genes at specific time points after the exposure to oxidative stress.

