# Lab Report 2

## Zachary Foster

## April 18, 2014

Dr. Muller suppiled us with contigs generated from a metagemoic assembly of DNA from a sample of marine water. The contigs are thought to be an undescribed lineage wihtin *verrucomicrobia*. We used BLAST, GLIMMER, and web-based tools to idenitify ORFs within one of the contigs and pridict possible functions. Using an interpolated markov model supplied by the instructors, we used GLIMMER to pridict ORFs. We used a custom `perl` script `parse_Glimmer.pl` to extract the amino acid sequence of the predicted ORFs. We then made a BLAST database, using sequences supplied by the instructors, and used the database to identify the ORFs. 4 ORFs were found, 2 of which return significant BLAST hits using an -e threshold of 0.01. Each of the 2 ORFs with siginificant hits had one identical hit and two very weak hits. We then used a custom `perl` script `grab_seq.pl` to exctract one of the ORFs with and identical hit to the reference database.

The extracted sequence was search aginst the genbank nr database using blastp with default settings 1. All of the hits were hypothetical protiens and only one was full length. Next, a PSI blast search with 3 iterations, useing default parameters, was used in order to discover more distaly realted putative homologs. This revealed many more full length matches, but all were stil hypothetical protiens 2. We used the COG and CCD databases to search for conserved domains using default parameters. The COG database search returned a hit to a chromosome segregation ATPase 3, but the CCD databasse search did not return any hits with known function 4 Finally, we searched pfam in order to detect any protien families that the unknown sequence may be a part of, but this to return no significant result with known function. The 2 hits that were found were of unknown function 5. From all of the analyses overall, it appears that this ORF is ether a novel protien that has not been assigned a faunction yet, or a conserved sequence suseptible to missaignment by automatic *ab inito* gene annotaing program, since there are many hits to multiple databases, but little functional information. The one signifiacnt hit with functional information was a chromosome segregation ATPase. but this alone is not suffiecnt evidence to predict the function of this unknown sequence.
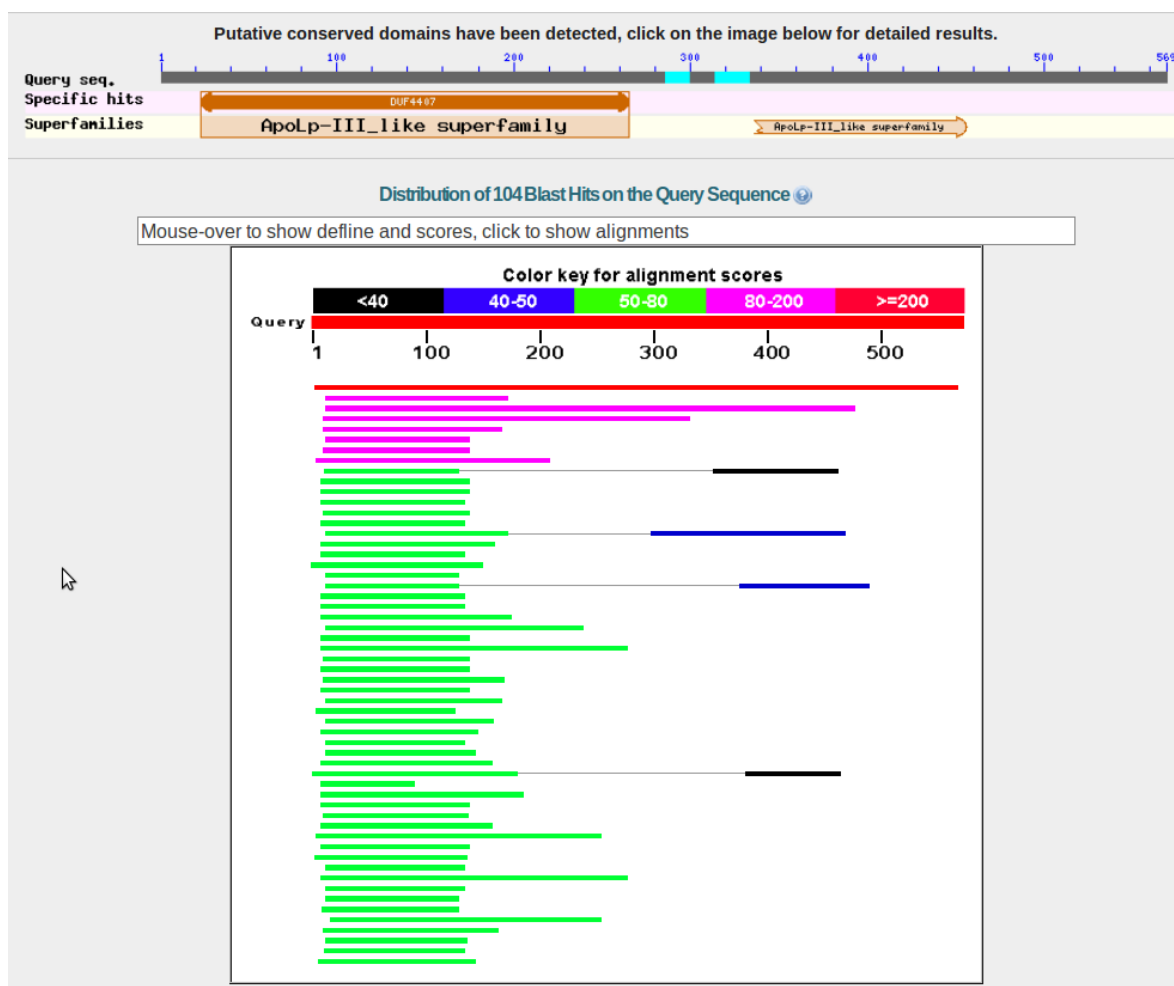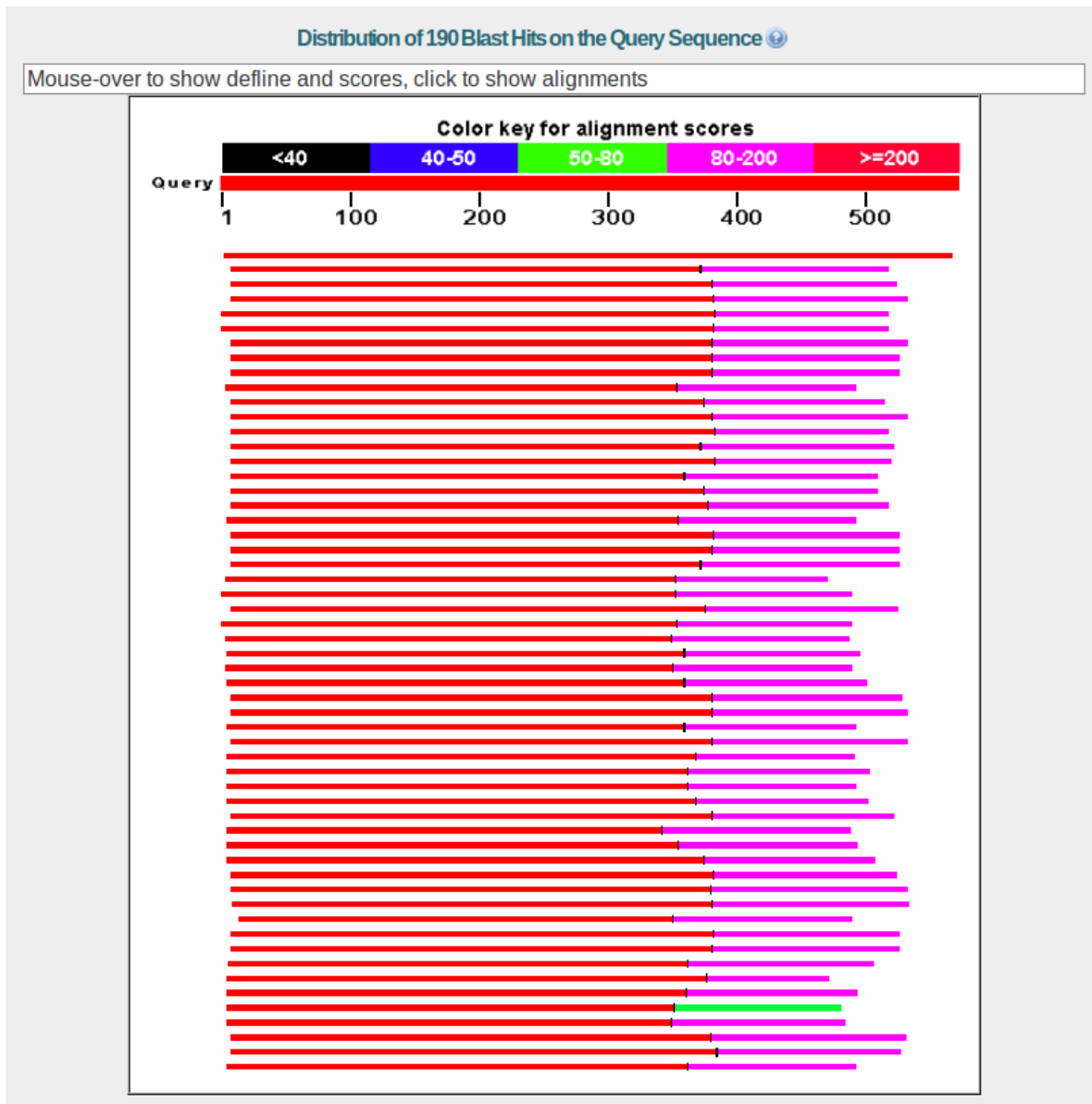
Figure 1: Screenshot of blastp results
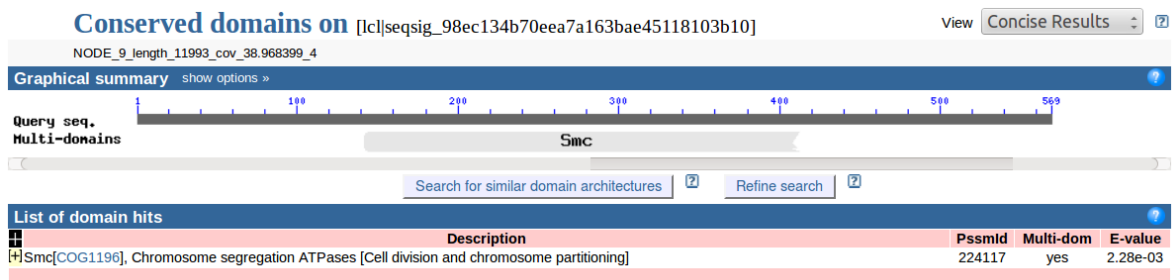
2

Figure 2: Screenshot of PSI-BLAST results
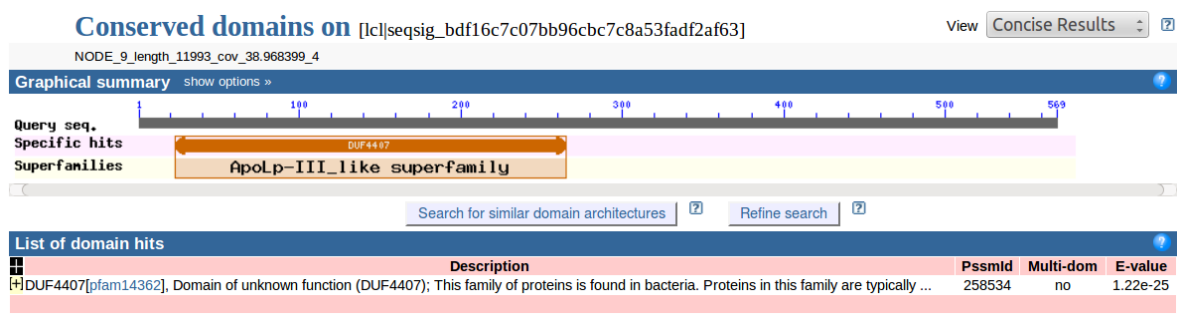


Figure 3: Screenshot of COG search results.

3

Figure 4: Screenshot of CCD search results



Figure 5: Screenshot of pfam results