✓ Dec 23: Obtain gene expression data for each bacteria

X Jan 6: Write up methods for COG and sub paper

✓ Jan 6: Read papers on gene expression

✓ Jan 6: Apply for McMaster Bursaries and Grants

✓ Jan 6: Conference Grants Completed

✓ Feb 16: Have pipeline in R for normalizing raw counts

✓ Mar 2: Have code for plotting gene expression and substitution graphs

✓ Mar 17: Have all *S. meliloti* chrom, *E. coli* and *B. subtilis* data sets combined and into one graph

✓ Mar 31: Have all *Streptomyces* data sets combined and into one graph (will take more time because it is not all the same strain)

✓ April 5: ISMB Chicago Conference Abstract Due

✓ April 7: Have something figured out for the pSymB and pSymA gene expression datasets

✓ April 27: Create regression lines for gene expression

April 27: Write up gene expression stuff

April 27: Manuscript for Substitution and Gene expression paper finished

~~April 27: Make 2nd, 3rd, 4th, order regression lines for substitution data~~

May 31: Have data for other molecular trends (GC content, number of genes, essential gene lists..etc.) combined with graphs (or in supplement) for sub analysis

May 31: Complete COG analysis

Jun 30: Gene Expression analysis write up

Jun 30: COG analysis Paper draft completed

Jul 31: Updated Sub Paper methods and results

Jul 31: Add other mol trends to Sub Paper

# Last Week

I realized that for pSymB I miscalculated the bidirectionality transformation, so I had to fix this and re-run everything. It did not change the logistic regression results (seen below). However, when I re-did the origin shuffling to see if the placement of the origin changed anything, moving the origin 100kb, 90kb and 80kb to the left made the logistic regression negative. I have been trying to figure out why this is happening but I am having no luck. I thought maybe it was because these shufflings are now 700kb away from the terminus, but the actual origin is about the same distance. I am still trying to figure this out but I am not sure what to do or what it means about the robustness of the origin shuffling. I also created simple linear regressions for the gene expression data and the results are summarized in the table below.

I looked at the gene expression data for *S. meliloti* in detail and it appears to look ok. When I graph the raw data and plot the regression line it looks like there is no trend for *S. meliloti*, the points are evenly distributed throughout the genome and there appears to be no increasing or decreasing trend. When looking at the other bacteria's raw data, there is clearly a decreasing trend when moving away from the origin. Additionally the number of genes and number of replicates are all comparable between *S. meliloti* and the other bacteria. So I think that the reason *S. meliloti* does not have significant gene expression regressions is because there is simply no trend.

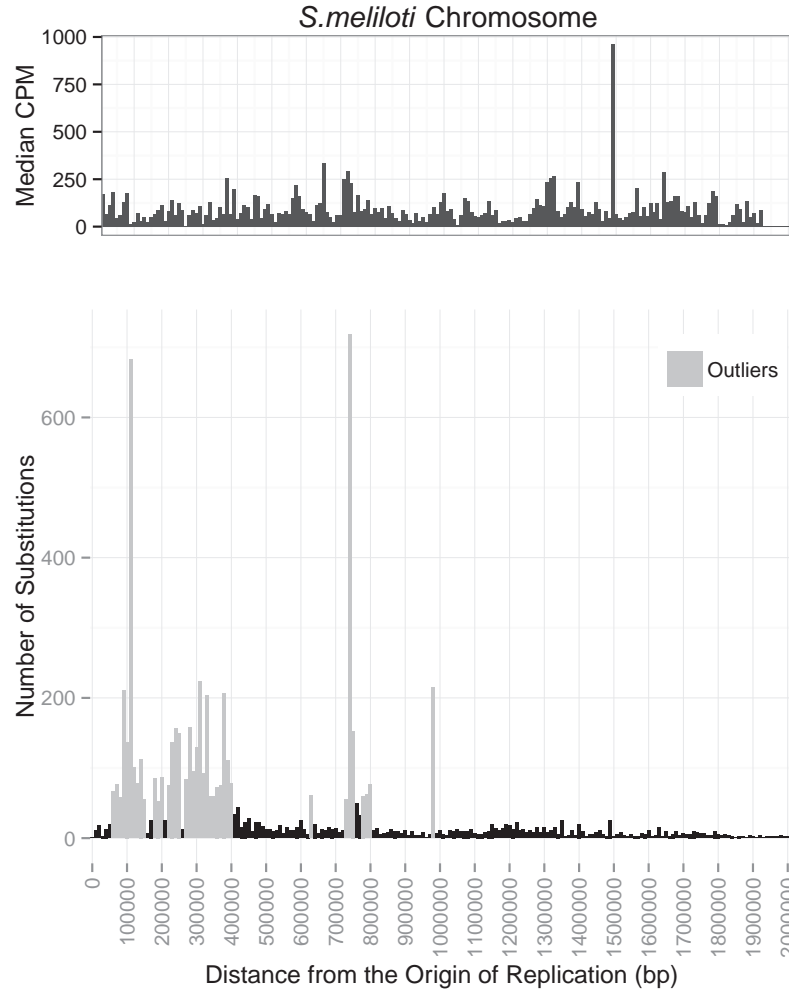I spent a few days last week writing up the manuscript for the gene

expression and substitution stuff.

# This Week

I need to figure out the weird thing happening with the pSymB robust origin shuffling test. I also plan on having my manuscript finished before I leave for camping on Friday.

| Bacteria and Replicon | Coefficient Estimate | Standard Error | P-value |
|---|---|---|---|
| *E. coli* Chromosome | $-6.41 \times 10^{-5}$ | $1.65 \times 10^{-5}$ | $1.1 \times 10^{-4}$ |
| *B. subtilis* Chromosome | $-9.9 \times 10^{-5}$ | $2.18 \times 10^{-5}$ | $6 \times 10^{-6}$ |
| *Streptomyces* Chromosome | $-1.5 \times 10^{-6}$ | $1.4 \times 10^{-7}$ | $<2 \times 10^{-16}$ |
| *S. meliloti* Chromosome | $3.19 \times 10^{-5}$ | $3.57 \times 10^{-5}$ | $3.7 \times 10^{-1}$ |
| *S. meliloti* pSymA | $-5.36 \times 10^{-5}$ | $6.34 \times 10^{-4}$ | $9.33 \times 10^{-1}$ |
| *S. meliloti* pSymB | $5.05 \times 10^{-4}$ | $2.6 \times 10^{-4}$ | $5.3 \times 10^{-2}$ |

Table 1: Linear regression analysis of the median counts per million expression data along the genome of the respective bacteria replicons. Grey coloured boxes indicate statistically significant results at the 0.5 significance level. Linear regression was calculated after the origin of replication was moved to the beginning of the genome and all subsequent positions were scaled around the origin accounting for bidirectionality of replication.

| Bacteria and Replicon | Coefficient Estimate | Standard Error | P-value |
|---|---|---|---|
| *E. coli* Chromosome | -1.394×10$^{-7}$ | 2.425×10$^{-9}$ | <2×10$^{-16}$ |
| *B. subtilis* Chromosome | -2.538×10$^{-8}$ | 1.58×10$^{-9}$ | <2×10$^{-16}$ |
| *Streptomyces* Chromosome | 1.736×10$^{-8}$ | 7.231×10$^{-10}$ | <2×10$^{-16}$ |
| *S. meliloti* Chromosome | -1.541×10$^{-6}$ | 3.042×10$^{-8}$ | <2×10$^{-16}$ |
| *S. meliloti* pSymA | -9.130×10$^{-7}$ | 1.975×10$^{-8}$ | <2×10$^{-16}$ |
| *S. meliloti* pSymB | 2.488×10$^{-7}$ | 1.964×10$^{-8}$ | <2×10$^{-16}$ |

Table 2: Logistic regression analysis of the number of substitutions along the genome of the respective bacteria replicons. Grey coloured boxes indicate a negative logistic regression coefficient estimate. All results are statistically significant. Logistic regression was calculated after the origin of replication was moved to the beginning of the genome and all subsequent positions were scaled around the origin accounting for bidirectionality of replication.

4

*S.meliloti* pSymA

*S.meliloti* pSymB

*B.subtilis* Chromosome

*Streptomyces* Chromosome