

Subs Paper Things to Do:

- ~~# of coding and non-coding sites~~
- ~~# of subs in each of  $\uparrow$~~
- ~~Look into *Streptomyces* non-coding issue~~
- ~~Look into *E. coli* coding issue~~
- ~~Look into pSymB coding/non-coding trend weirdness~~
- ~~Figure out why *Streptomyces* appears to have tons of coding data missing~~
- ~~Figure out what is going on with cod/non-cod code and why it is still not working!~~
- ~~write up methods for coding/non-coding~~
- ~~write methods and results for clustering~~
- ~~start code to split alignment into multiple alignments of each gene~~
- ~~figure out how to deal with overlapping genes~~
- ~~figure out how to deal with gaps in gene of ref taxa~~
- ~~split up the alignment into multiple alignments of each gene~~
- ~~check if each gene alignment is a multiple of 3 (proper codon alignment)~~
- ~~get dN/dS for coding/non-coding stuff per gene~~
- ~~Or get 1st, 2nd, 3rd codon pos log regs~~
- ~~write up coding/non-coding results~~
- ~~take out gene expression from this paper~~
- ~~write better intro/methods for distribution of subs graphs~~
- ~~write discussion for coding/non-coding~~
- ~~write coding/non-coding into conclusion~~
- ~~figured out pipeline for CODEML to calculate dN/dS for each gene~~
- ~~make a list of what should be in supplementary files for subs paper~~
- ~~put everything in list into supplementary file for subs paper~~
- ~~write dN/dS methods~~
- ~~write dN/dS results~~
- ~~write dN/dS discussion~~

- write dN/dS into conclusion
- ~~new bar graph with coding and non-coding sites separated~~
- mol clock for my analysis?
- GC content? COG? where do these fit?

#### Gene Expression Paper Things to Do:

- ~~look for more GEO expression data for *S. meliloti*~~
- ~~look for more GEO expression data for *Streptomyces*~~
- ~~look for more GEO expression data for *B. subtilis*~~
- ~~format paper and put in stuff that is already written~~
- ~~look for more GEO expression data for *E. coli*~~
- ~~Get numbers for how many different strains and multiples of each strain I have for gene expression~~
- ~~re-do gene expression analysis for *B. subtilis*~~
- ~~re-do gene expression analysis for *E. coli*~~
- ~~find papers about what has been done with gene expression~~
- ~~read papers ↑~~
- ~~put notes from ↑ papers into word doc~~
- write abstract
- ~~write intro~~
- add stuff from outline to Data section
- create graphs for expression distribution (no sub data)
- add # of genes to expression graphs (top)
- average gene expression
- write discussion
- write conclusion
- add into methods: filters for Hiseq, RT PCR and growth phases for data collection
- update supplementary figures/file

#### Inversions and Gene Expression Letter Things to Do:

- ~~get as much GEO data as possible~~
- ~~find papers about inversions and expression~~
- ~~see how many inversions I can identify in these strains of *Escherichia coli* with gene expression data~~
- ~~read papers about inversions~~
- check if opposite strand in progressiveMauve means an inversions (check visual matches the xmfa)
- check if PARSNP and progressiveMauve both identify the same inversions (check xmfa file)
- create latex template for paper
- ~~put notes from papers into doc~~
- use large PARSNP alignment to identify inversions
- confirm inversions with dot plot
- write outline for letter
- write Abstract
- write intro
- write methods
- compile tables (supplementary)
- write results
- write discussion
- write conclusion
- do same ancestral/phylogenetic analysis that I did in the subs paper

## Last Week

- ✓ re-run coding and non-coding substitution analysis for all bacteria

I worked on re-running the substitution analysis for the coding and non-coding sections for all the bacteria. The results from this are summarized in the table 1 below. The results are pretty consistent for the coding sections with a negative substitution trend (pSymB was the only exception). But for the non-coding sections the results are a bit all over the place, with some non-significant results and the rest split up half and half as to if the substitution trend is positive or negative. I am not really sure what to make of these results or how to explain them. The only thing I can think of is that non-coding sections are more variable than the coding sections in their

content so that is why we are seeing inconsistent results? I would really appreciate your thoughts on this.

I also began working on re-creating the substitution density graphs for both coding and non-coding sections for each replicon but was hitting some snags with how my data is formatted. I am sure that I can have this working and complete by the end of today.

## This Week

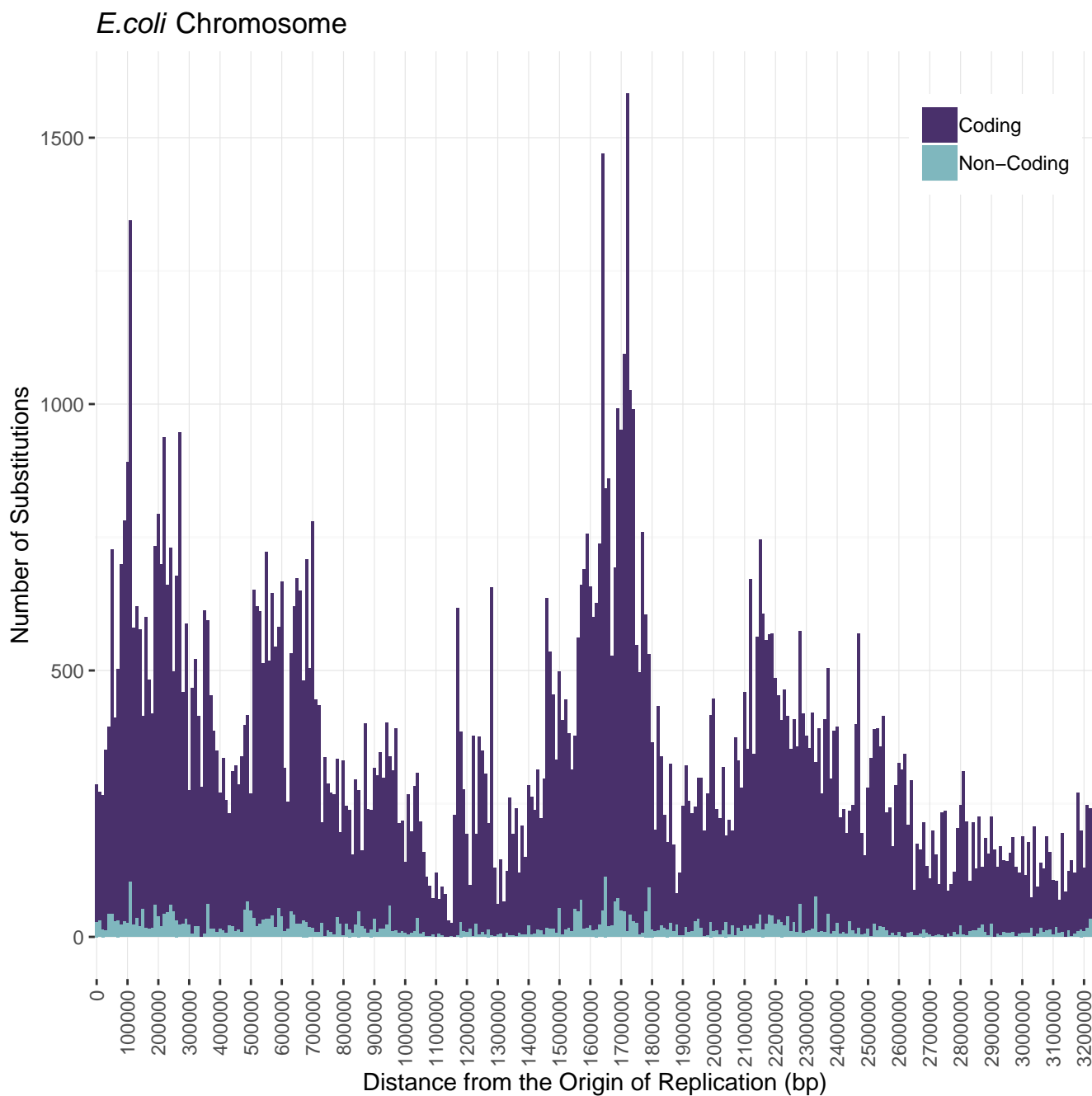
I plan on finishing the graphs for the coding/non-coding analysis, and I would like to again keep working on the pipeline for the dN/dS analysis. This weekend I would like to work on writing more for the gene expression paper and get a solid draft ready for that.

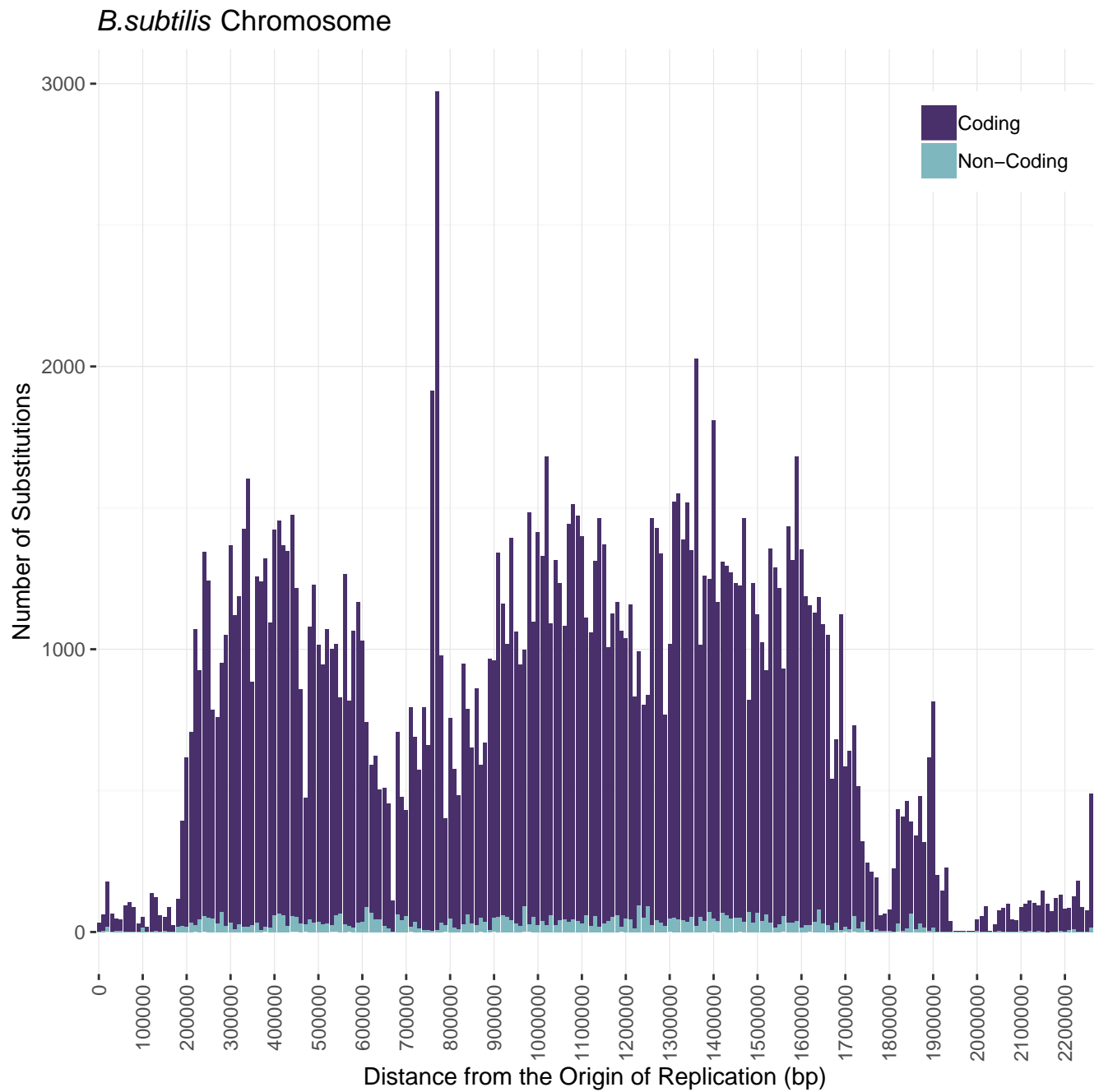
## Next Week

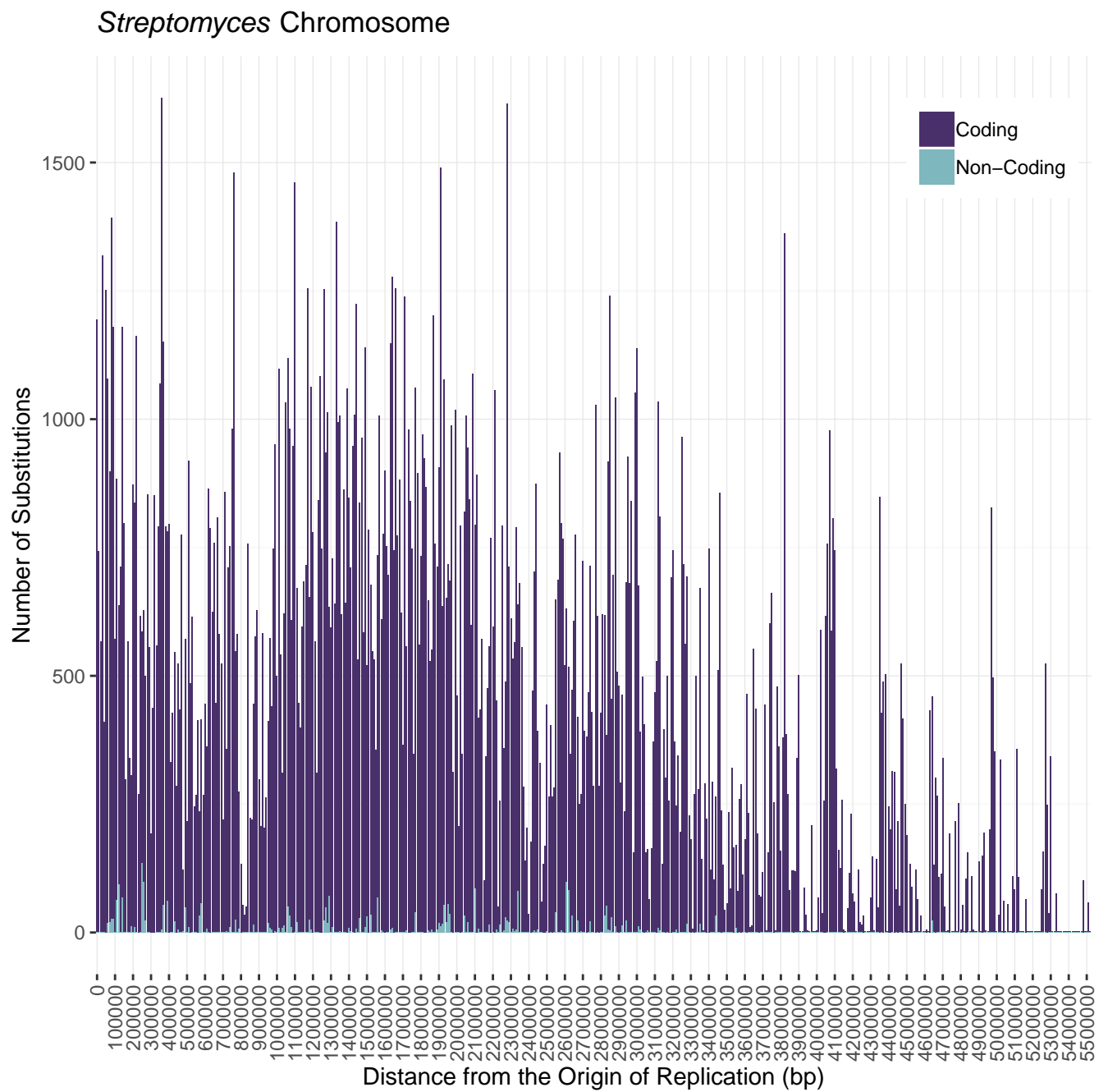
I would like to implement the above pipeline for calculating the dN/dS for each gene and put all these results in to a coherent table.

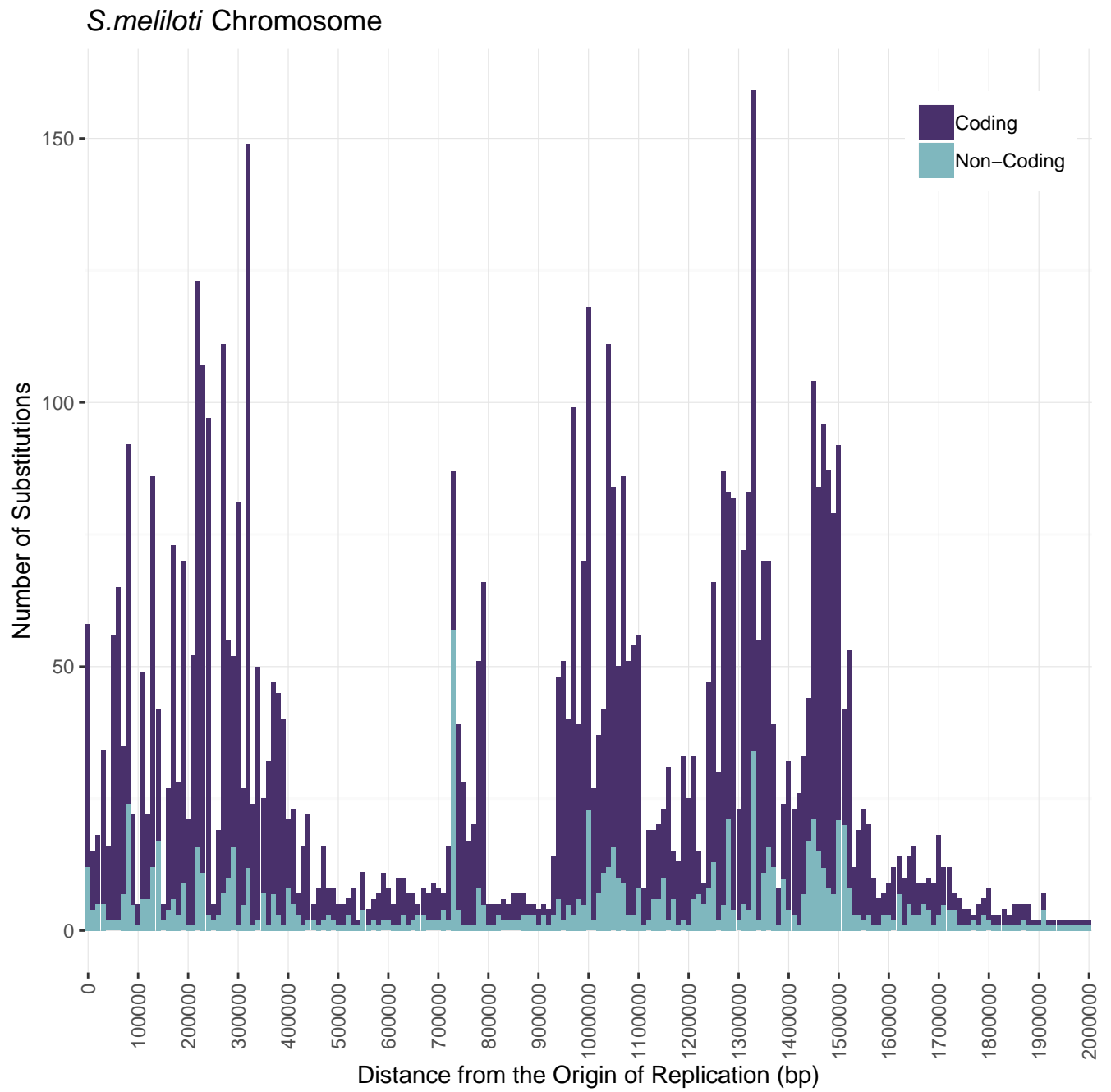
I want to begin figuring out how to obtain all inversions from the Mauve or PARSNP alignment for the inversions analysis.

Sub density graphs with coding and non-coding information

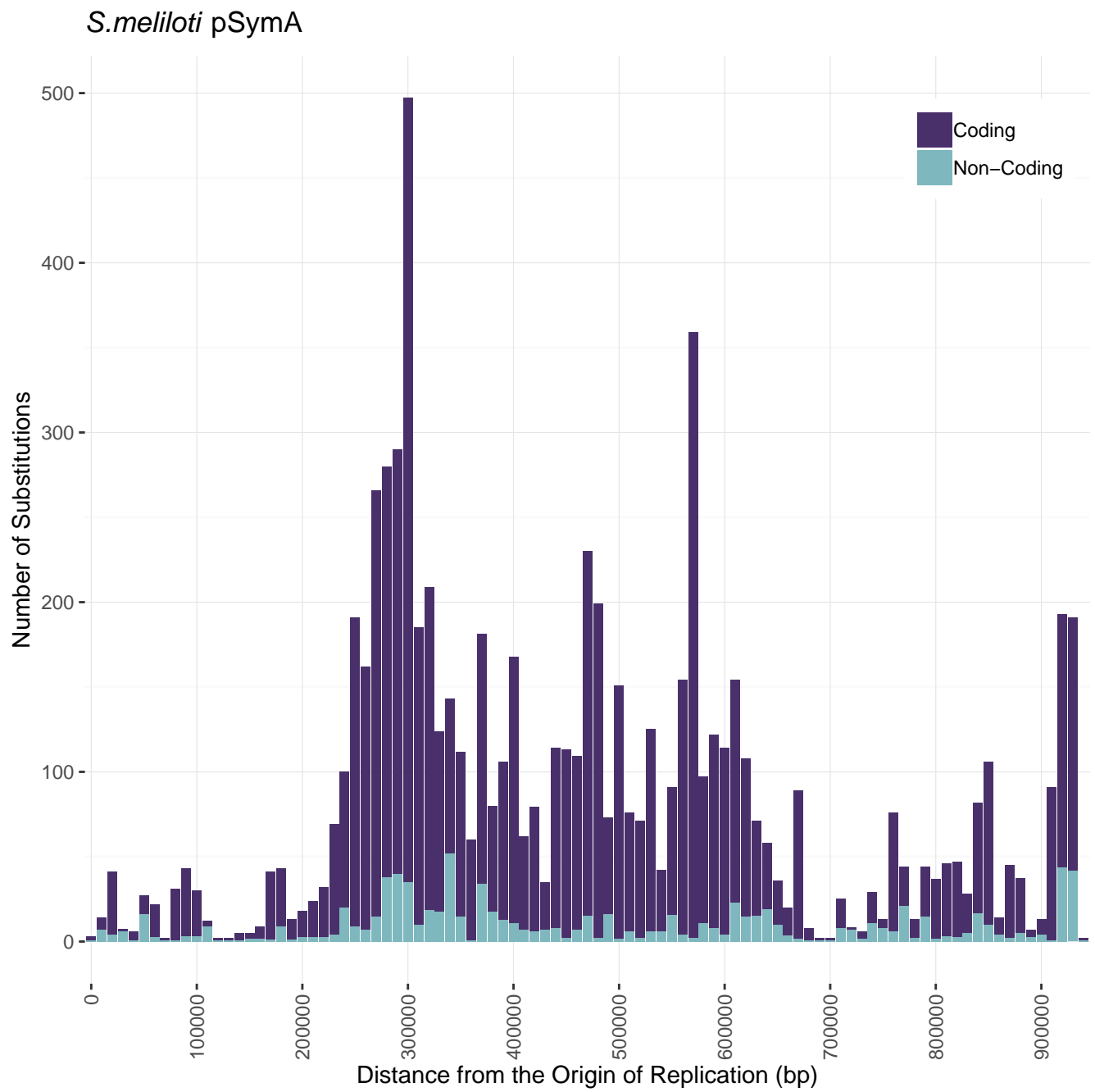


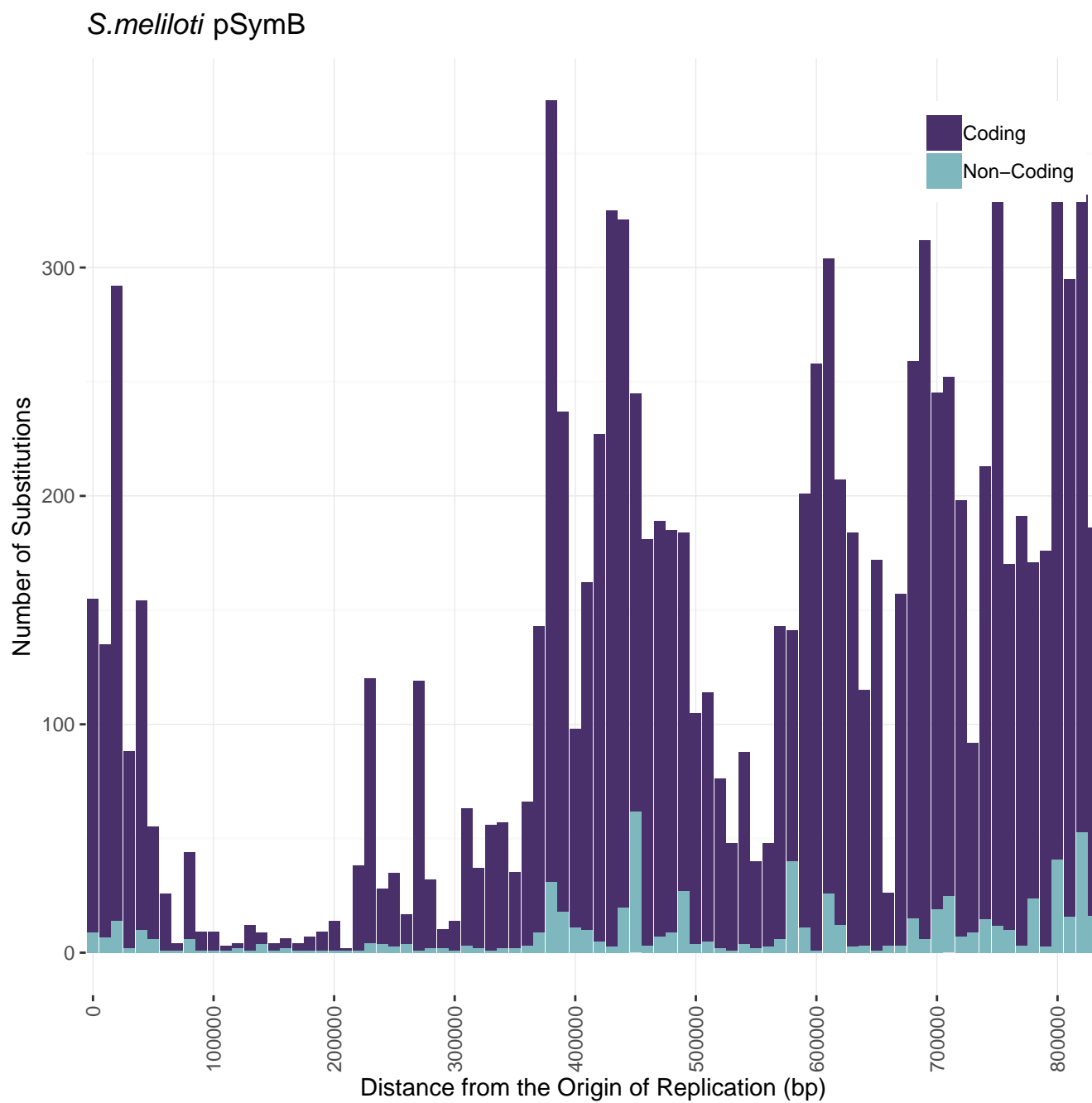


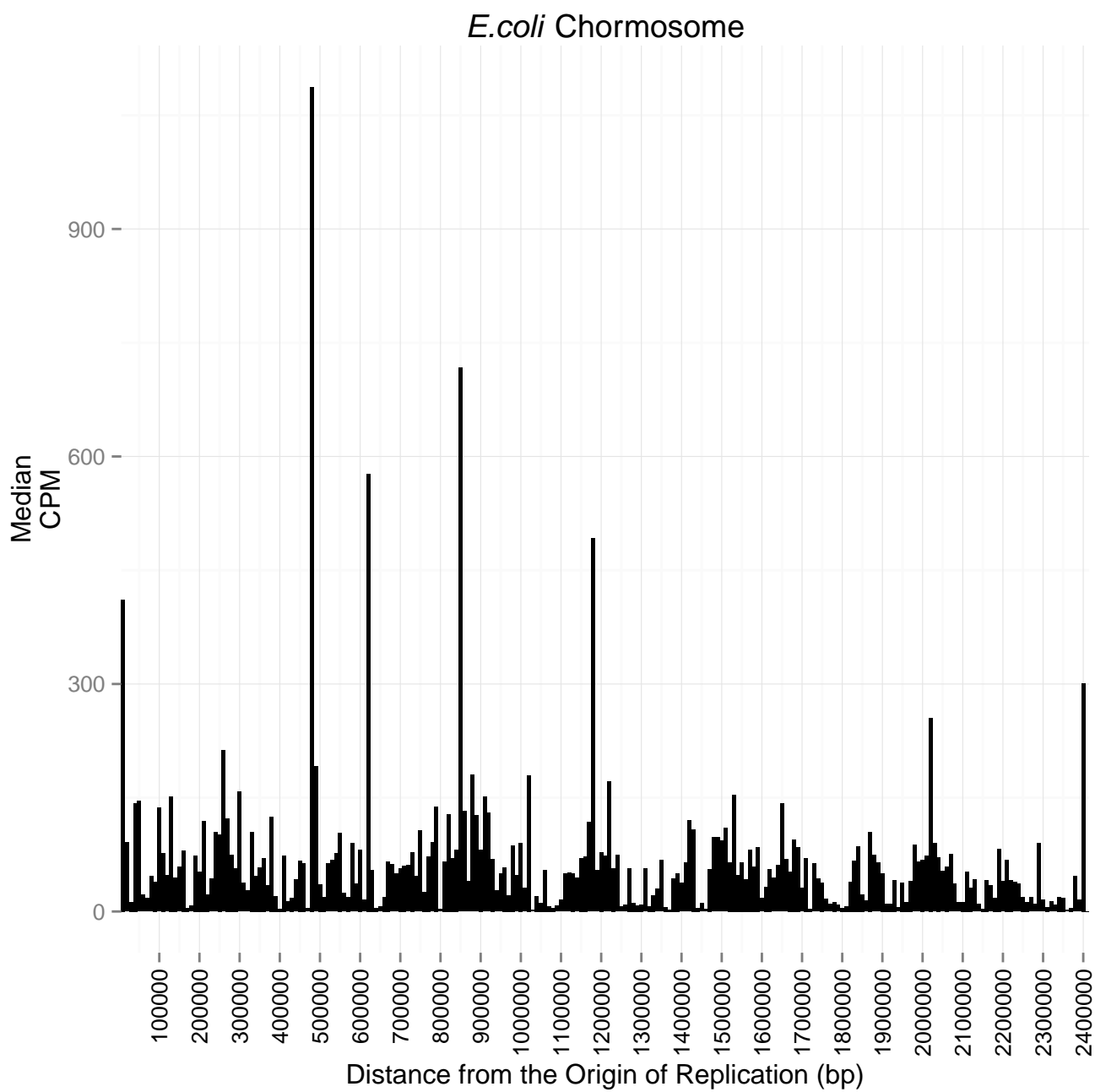


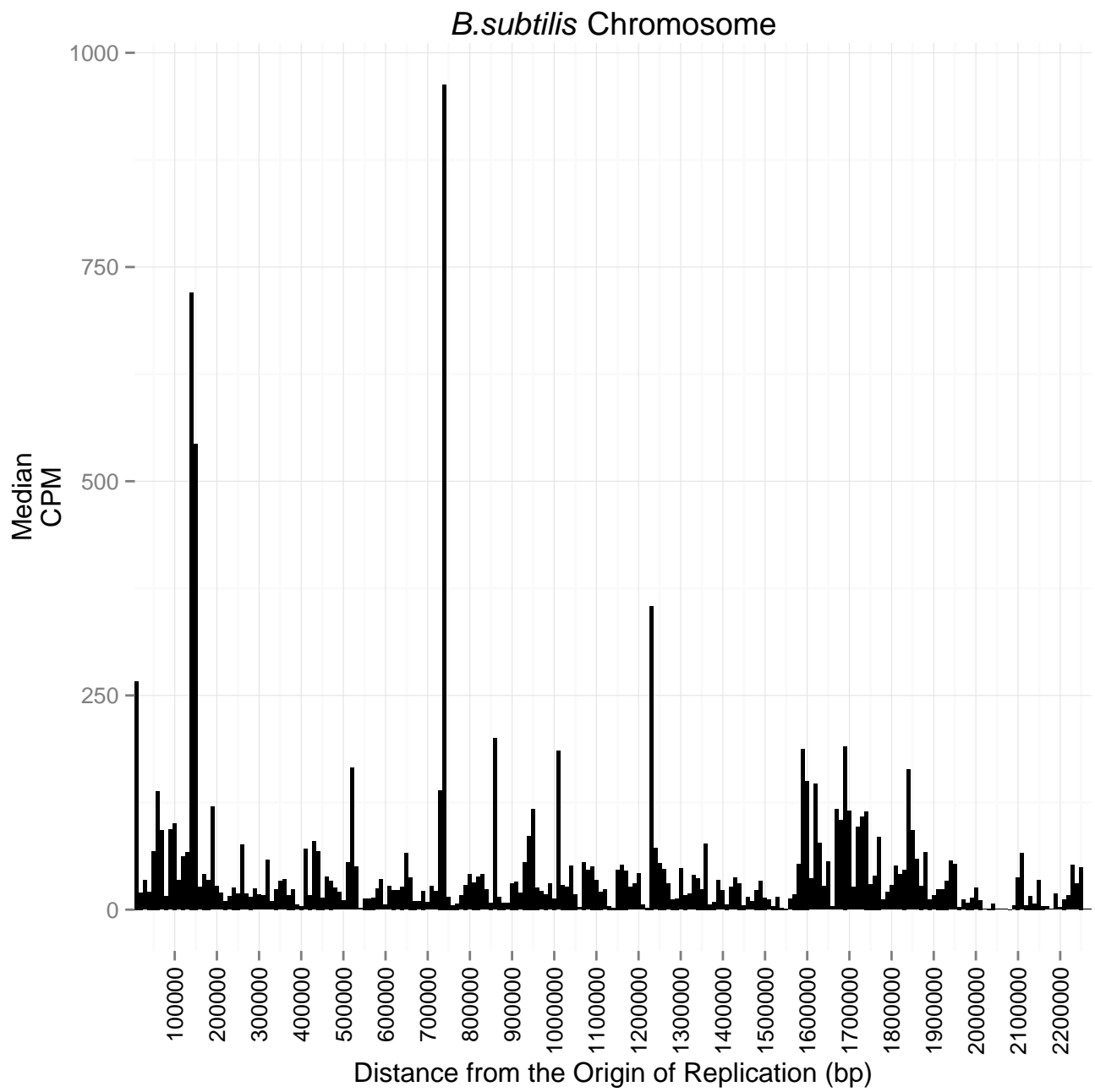


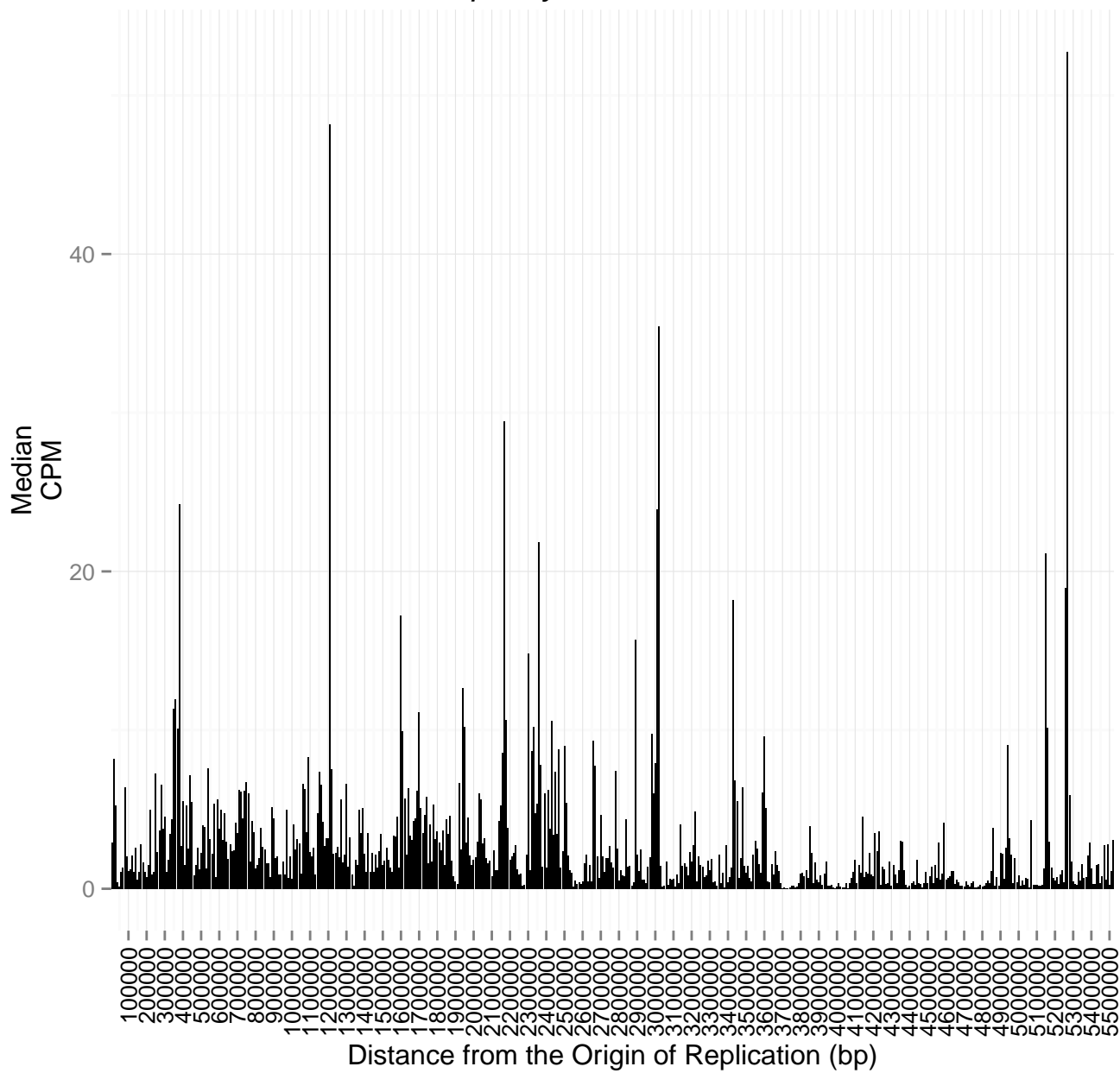


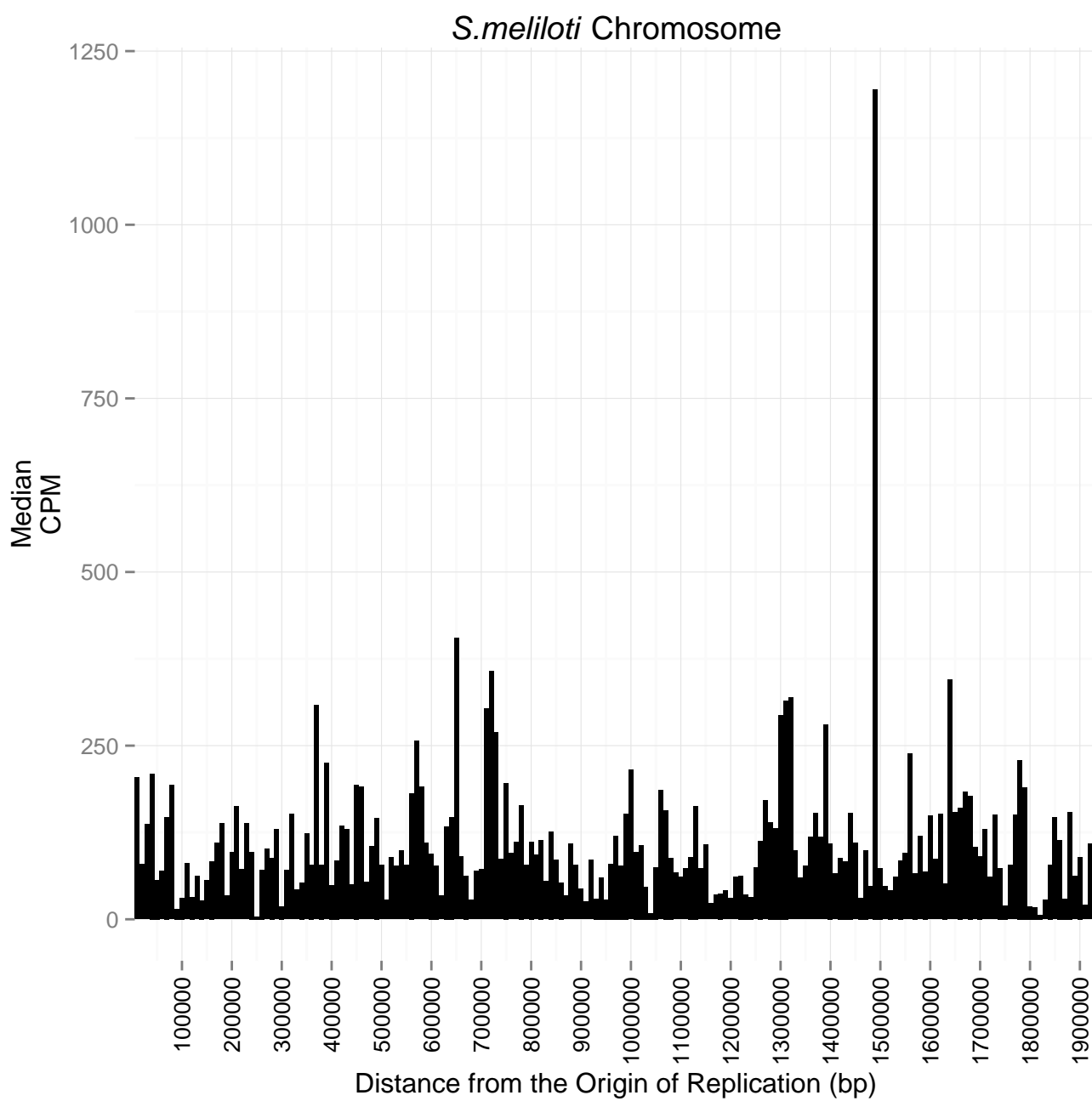


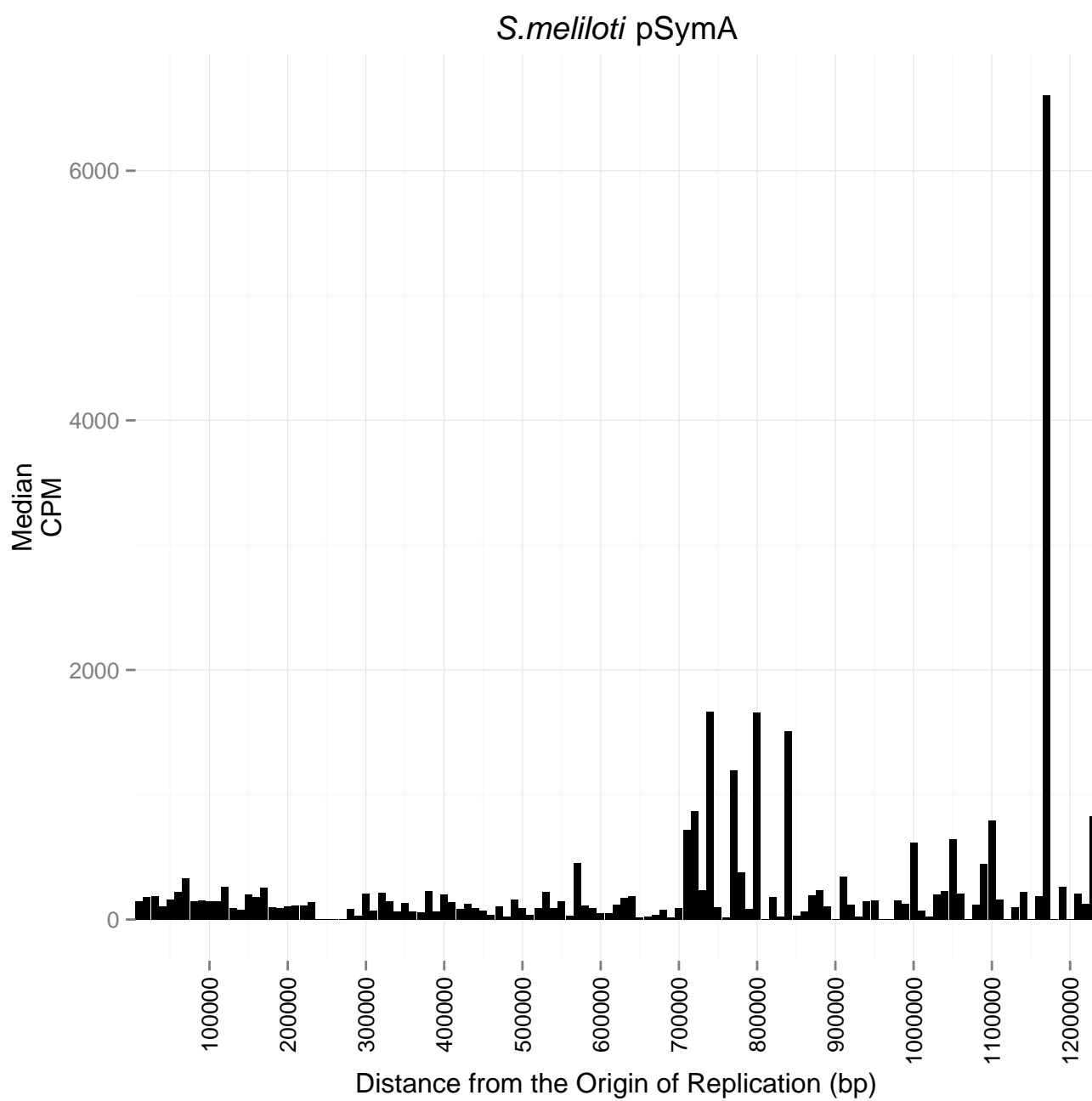


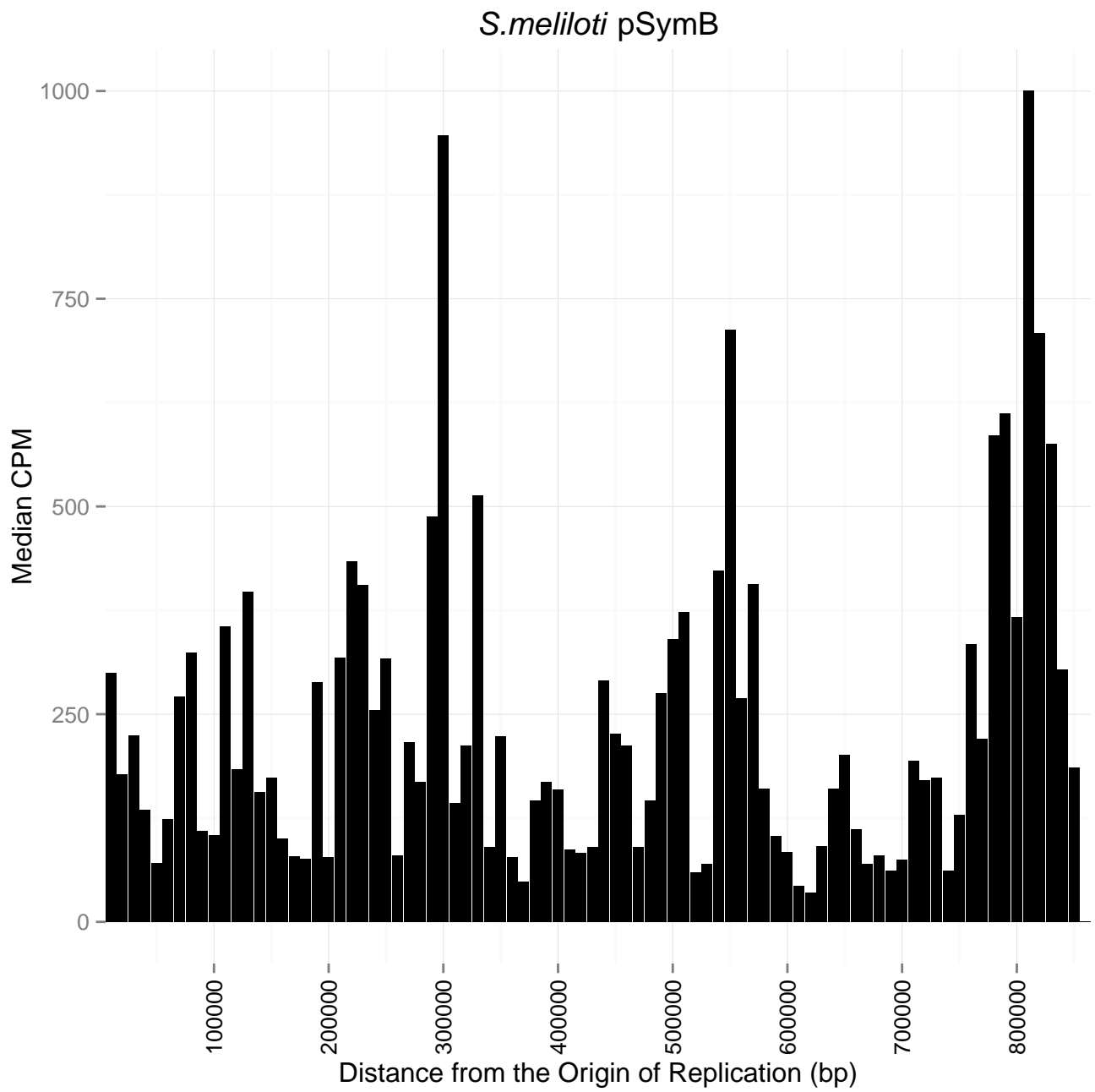




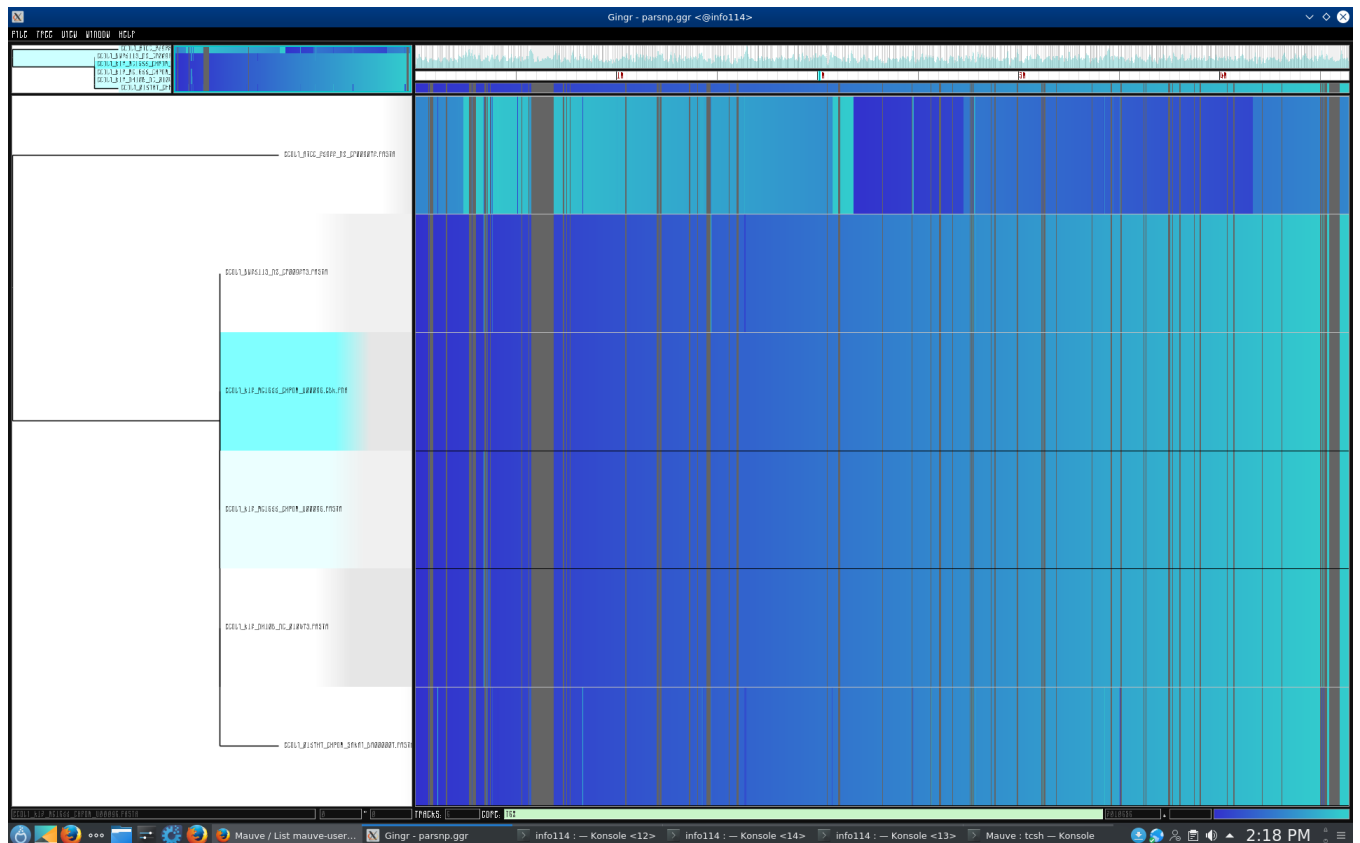
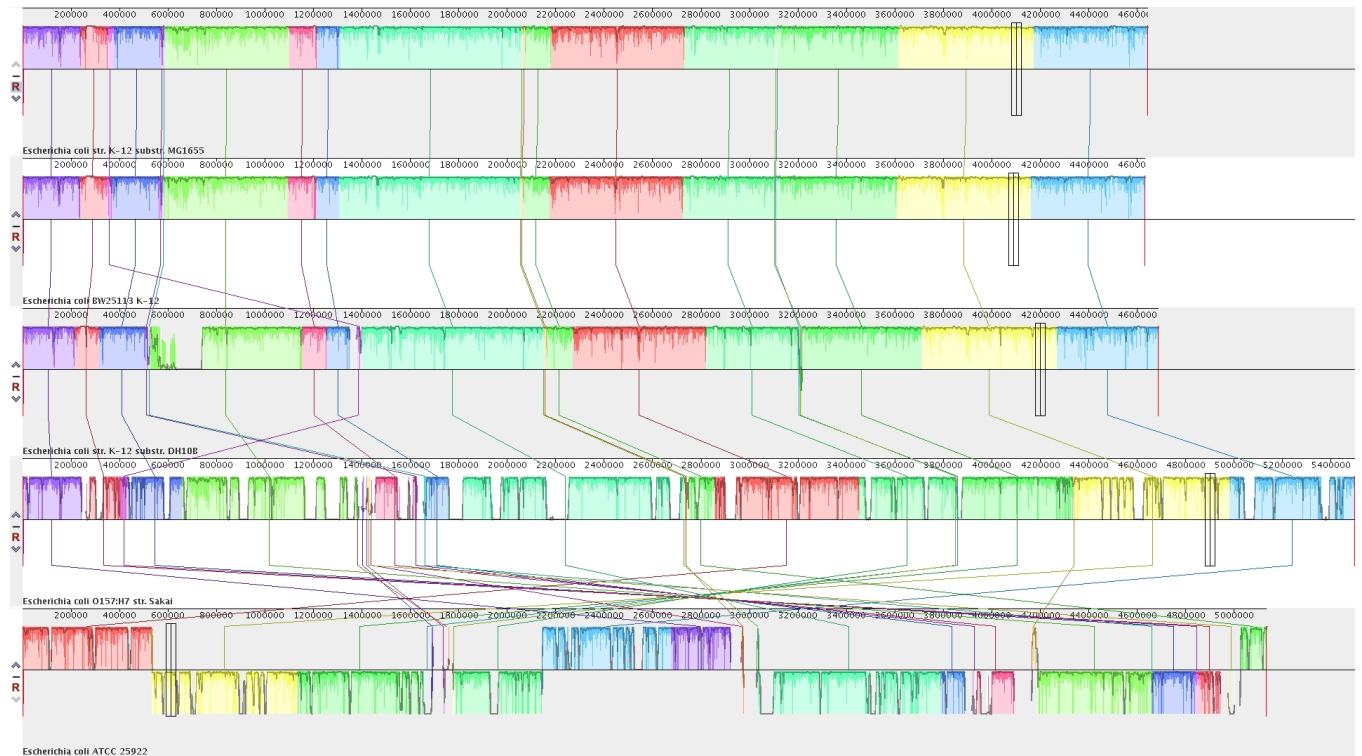
*Streptomyces* Chromosome











Bacteria and Replicon	Coding Sequences	Non-Coding Sequences
<i>E. coli</i> Chromosome	$-6.454 \times 10^{-8}***$	NS
<i>B. subtilis</i> Chromosome	$-1.159 \times 10^{-7}***$	$-9.861 \times 10^{-8}***$
<i>Streptomyces</i> Chromosome	$-8.464 \times 10^{-9}***$	$3.572 \times 10^{-7}***$
<i>S. meliloti</i> Chromosome	$-1.269 \times 10^{-7}***$	$-1.51 \times 10^{-7}*$
<i>S. meliloti</i> pSymA	$-2.02 \times 10^{-7}***$	NS
<i>S. meliloti</i> pSymB	$2.618 \times 10^{-7}***$	$8.591 \times 10^{-7}***$

Table 1: Logistic regression analysis of the number of substitutions along all positions of the genome of the respective bacteria replicons. These genomic positions were split up into the coding and non-coding regions of the genome. 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. All results are marked with significance codes as followed:  $< 0.001 = '***'$ ,  $0.001 < 0.01 = '**'$ ,  $0.01 < 0.05 = '*'$ ,  $> 0.05 = 'NS'$ .

Bacteria Strain/Species	GEO Accession Number	Date Accessed
<i>E. coli</i> K12 MG1655	GSE60522	December 20, 2017
<i>E. coli</i> K12 MG1655	GSE73673	December 19, 2017
<i>E. coli</i> K12 MG1655	GSE85914	December 19, 2017
<i>E. coli</i> K12 MG1655	GSE40313	November 21, 2018
<i>E. coli</i> K12 MG1655	GSE114917	November 22, 2018
<i>E. coli</i> K12 MG1655	GSE54199	November 26, 2018
<i>E. coli</i> K12 DH10B	GSE98890	December 19, 2017
<i>E. coli</i> BW25113	GSE73673	December 19, 2017
<i>E. coli</i> BW25113	GSE85914	December 19, 2017
<i>E. coli</i> O157:H7	GSE46120	August 28, 2018
<i>E. coli</i> ATCC 25922	GSE94978	November 23, 2018
<i>B. subtilis</i> 168	GSE104816	December 14, 2017
<i>B. subtilis</i> 168	GSE67058	December 16, 2017
<i>B. subtilis</i> 168	GSE93894	December 15, 2017
<i>B. subtilis</i> 168	GSE80786	November 16, 2018
<i>S. coelicolor</i> A3	GSE57268	March 16, 2018
<i>S. natalensis</i> HW-2	GSE112559	November 15, 2018
<i>S. meliloti</i> 1021 Chromosome	GSE69880	December 12, 2017
<i>S. meliloti</i> 2011 pSymA	NC_020527 (Dr. Finan)	April 4, 2018
<i>S. meliloti</i> 1021 pSymA	GSE69880	November 15, 18
<i>S. meliloti</i> 2011 pSymB	NC_020560 (Dr. Finan)	April 4, 2018
<i>S. meliloti</i> 1021 pSymB	GSE69880	November 15, 18

Table 2: Summary of strains and species found for each gene expression analysis. Gene Expression Omnibus accession numbers and date accessed are provided.

Bacteria and Replicon	Coefficient Estimate	Standard Error	P-value
<i>E. coli</i> Chromosome	$-6.03 \times 10^{-5}$	$1.28 \times 10^{-5}$	$2.8 \times 10^{-6}$
<i>B. subtilis</i> Chromosome	$-9.7 \times 10^{-5}$	$2.0 \times 10^{-5}$	$1.2 \times 10^{-6}$
<i>Streptomyces</i> Chromosome	$-1.17 \times 10^{-6}$	$1.04 \times 10^{-7}$	$< 2 \times 10^{-16}$
<i>S. meliloti</i> Chromosome	$3.97 \times 10^{-5}$	$4.25 \times 10^{-5}$	NS ( $3.5 \times 10^{-1}$ )
<i>S. meliloti</i> pSymA	$1.39 \times 10^{-3}$	$2.53 \times 10^{-4}$	$4.9 \times 10^{-8}$
<i>S. meliloti</i> pSymB	$1.46 \times 10^{-4}$	$2.03 \times 10^{-4}$	NS ( $5.34.7 \times 10^{-1}$ )

Table 3: 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.