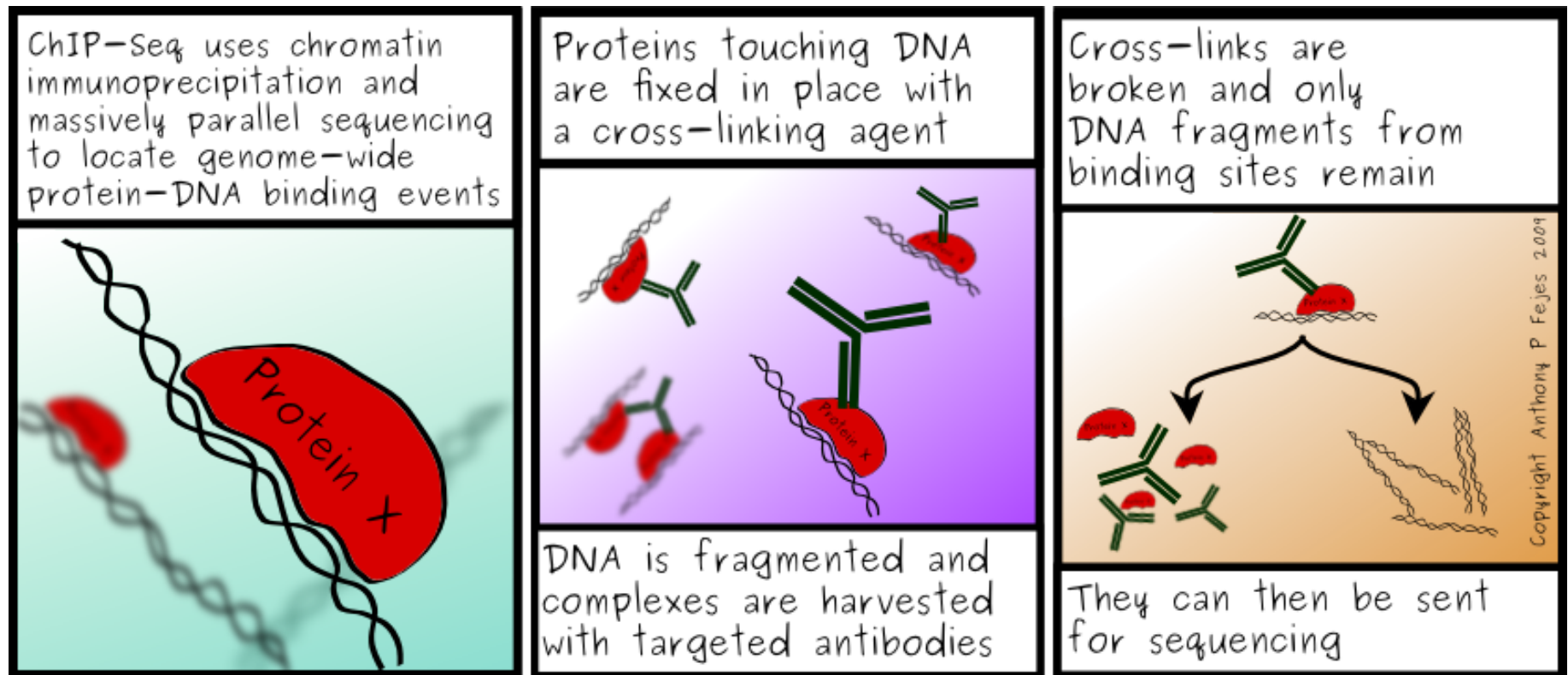


# ChIP-seq

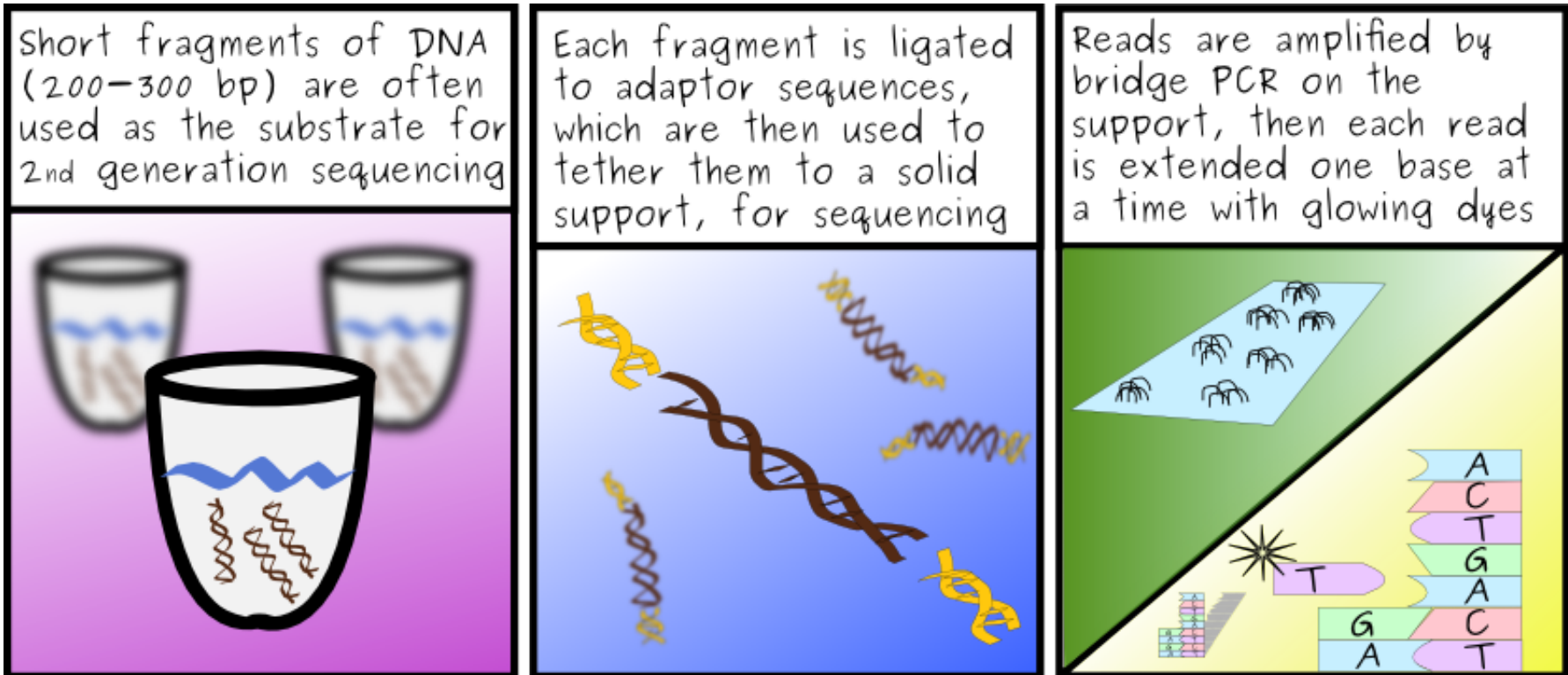
Mark Robinson

Department of Biochemistry and Molecular Biology,  
Michigan State University

# ChIP-seq: A method for determining protein-DNA interactions on a genome-wide scale

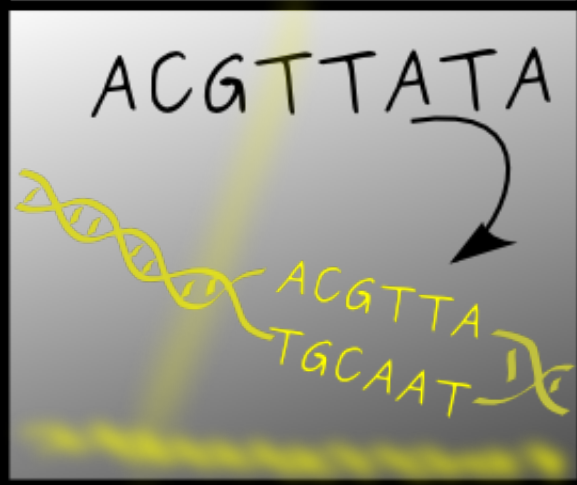


# Next generation sequencing: Short read lengths with massive depth of sequencing

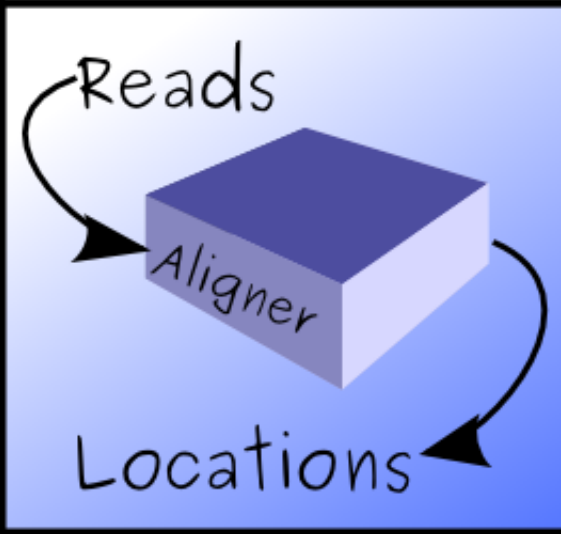


# Alignment of reads to reference genome allows mapping of binding locations

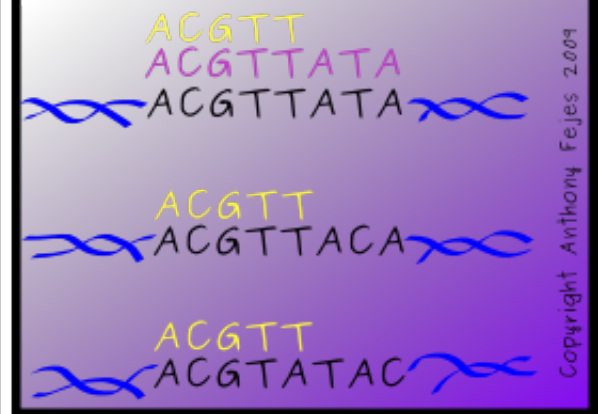
Sequenced reads can be aligned to the reference genome using an aligner, such as MAQ, bwa, Eland, Exonerate or Bowtie



Aligners work as a black box to locate the most likely point of origin of each sequenced read



The longer the reads, the more likely the aligner will find a unique (or best) point of origin  
-Most aligners do not require perfect matches



# Peak Detection software

- QuEST:

<http://mendel.stanford.edu/sidowlab/downloads/quest/>

- Findpeaks:

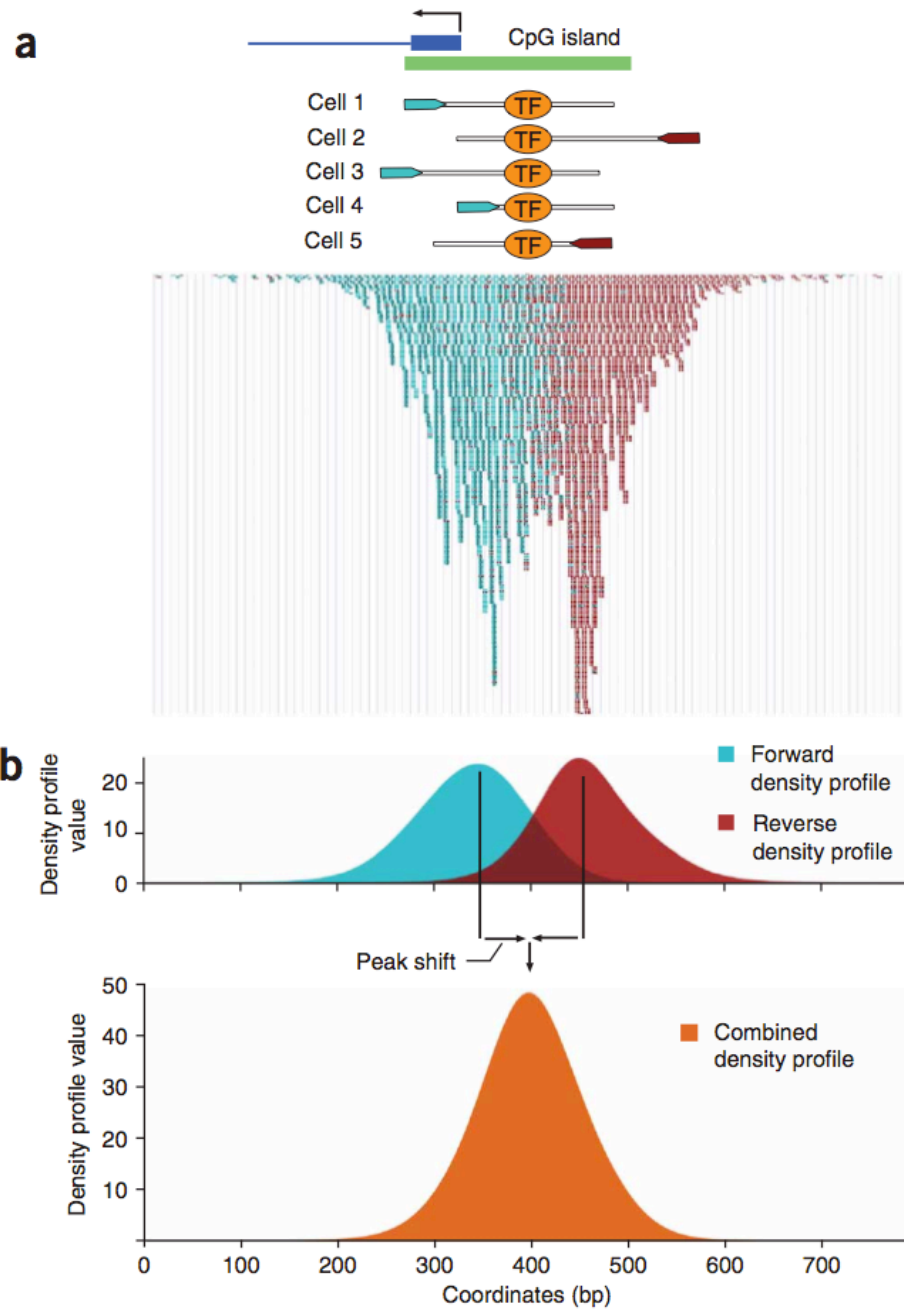
<http://www.bcgsc.ca/platform/bioinfo/software/findpeaks>

- ChIP-seq PeakFinder:

[http://woldlab.caltech.edu/html/chipseq\\_peak\\_finder](http://woldlab.caltech.edu/html/chipseq_peak_finder)

# Correct handling of short reads

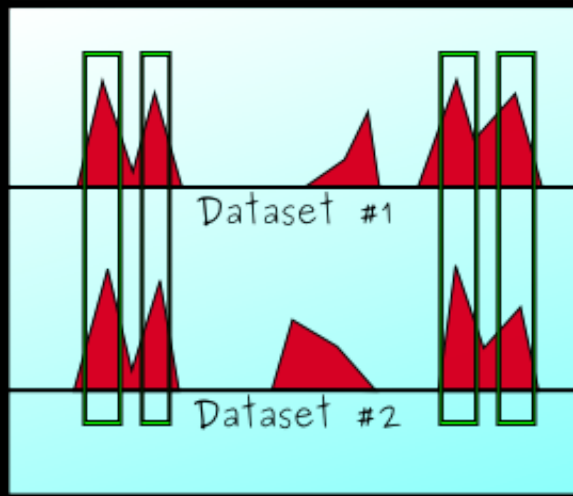
- Transcription factor binding sites are located somewhere on original DNA fragment
  - NOT the sequenced read
- Before attempting to determine binding site location this fact needs to be accounted for.



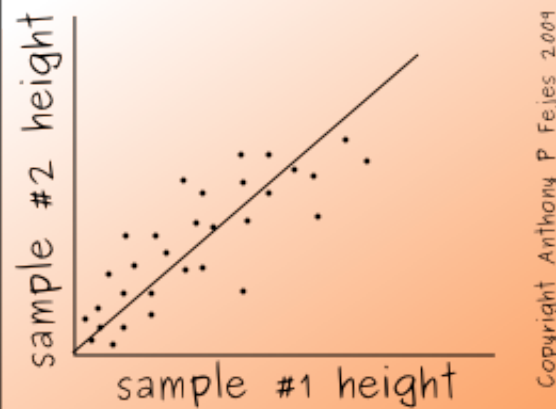
Valouev, A, et al.  
Nature Methods, 2008 Sep; 5(9):  
829-34

# Normalization and Significance detection

The first step of the compare function is to identify locations with similar peaks in both sets of reads

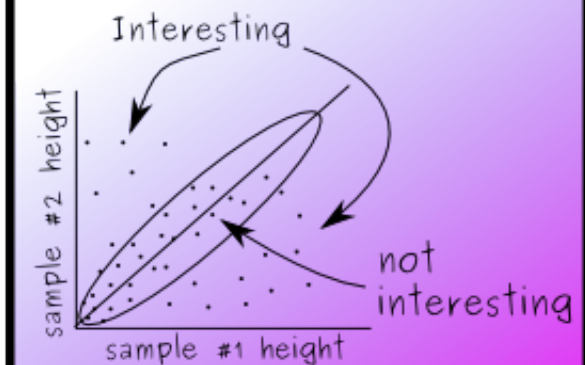


These points are plotted and a symmetrical best-fit regression line is calculated



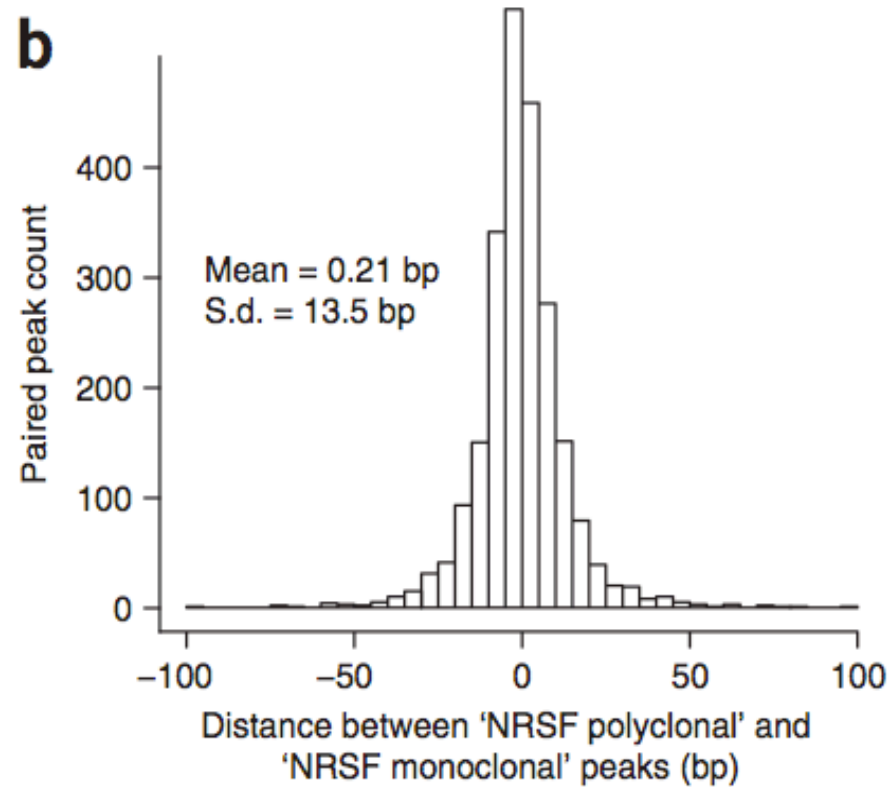
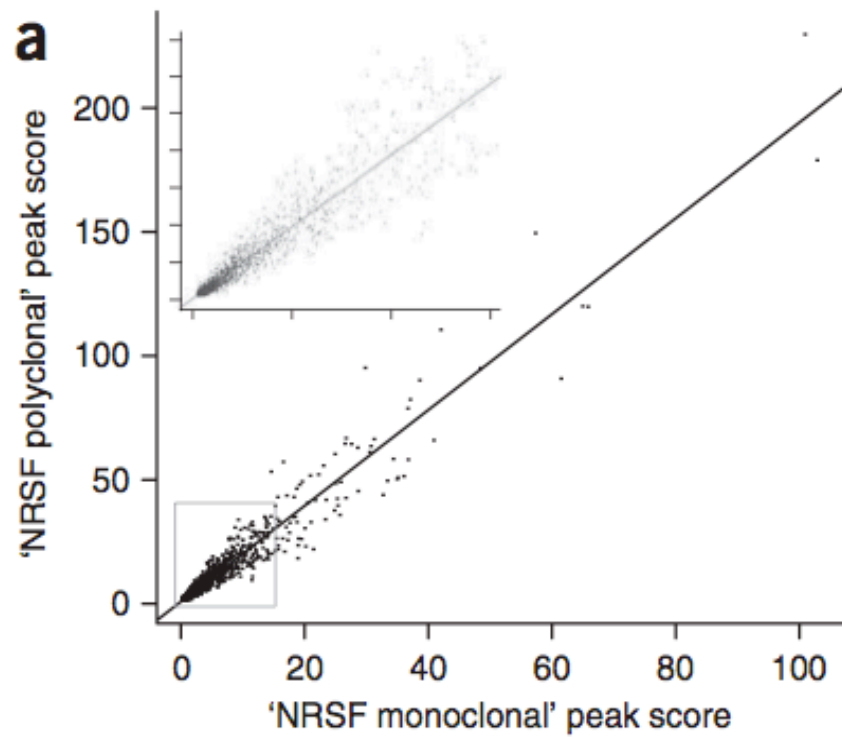
The best-fit line removes the need to normalize

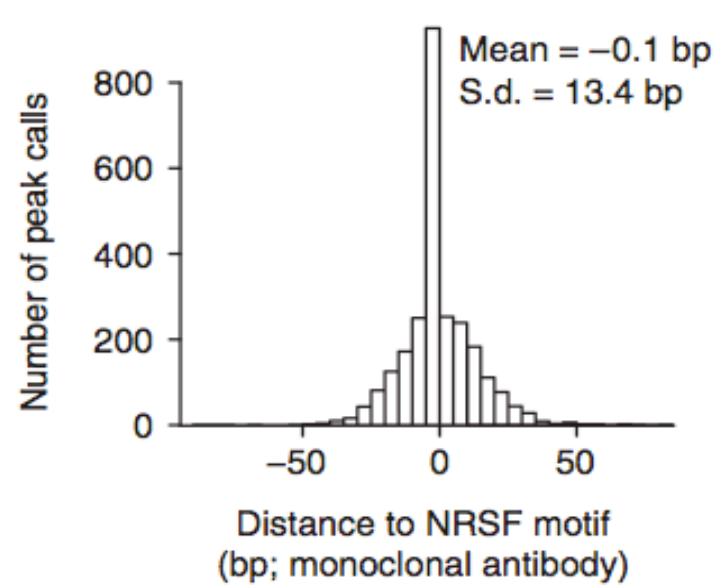
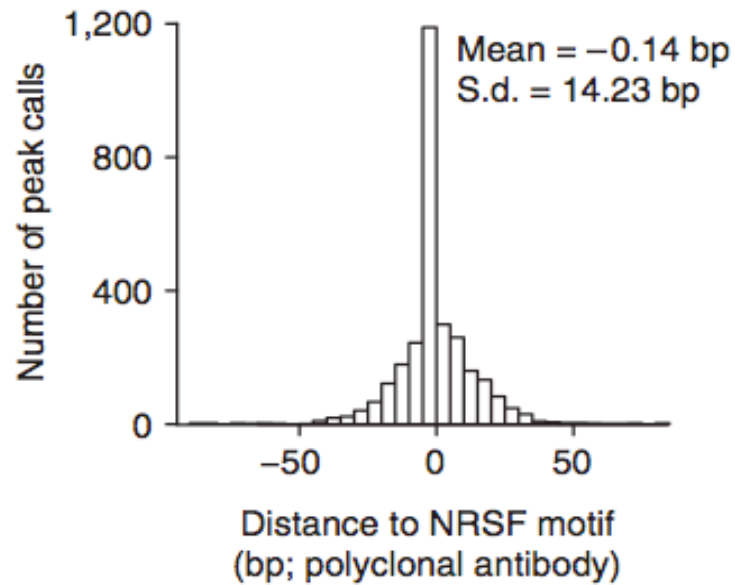
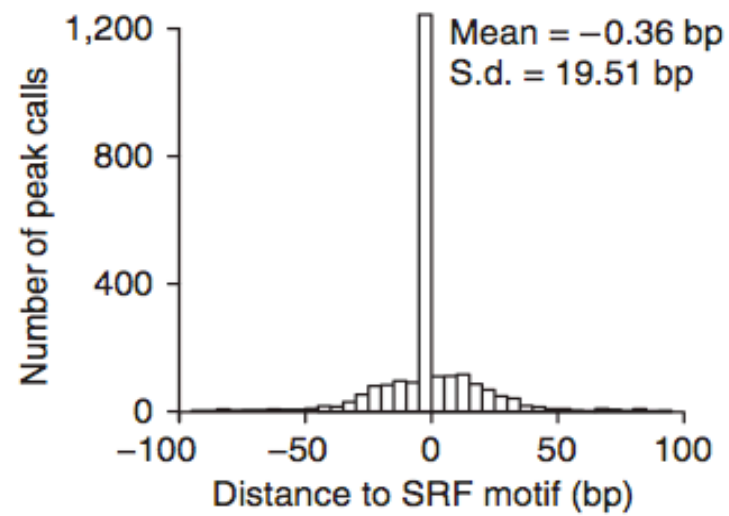
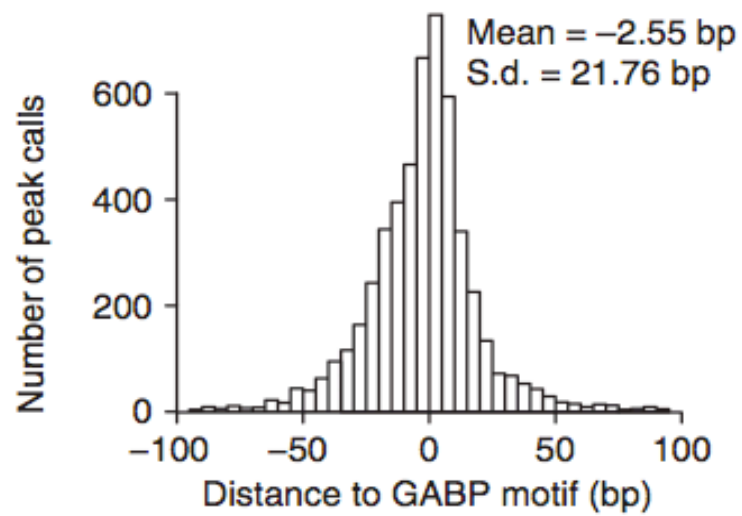
The distribution of all peak pairs around the line are calculated and points close to the line are removed



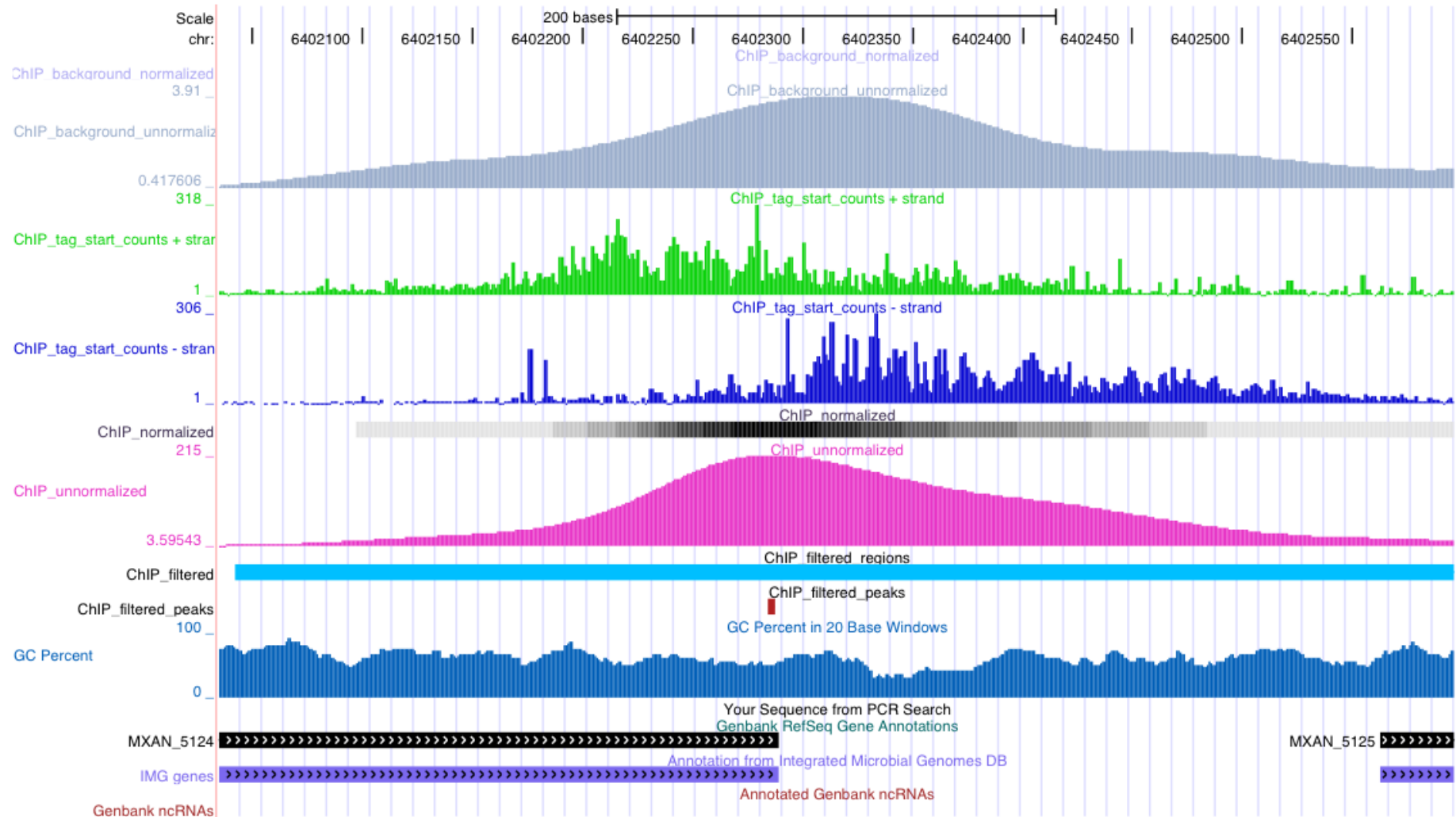
Any points remaining are statistical outliers



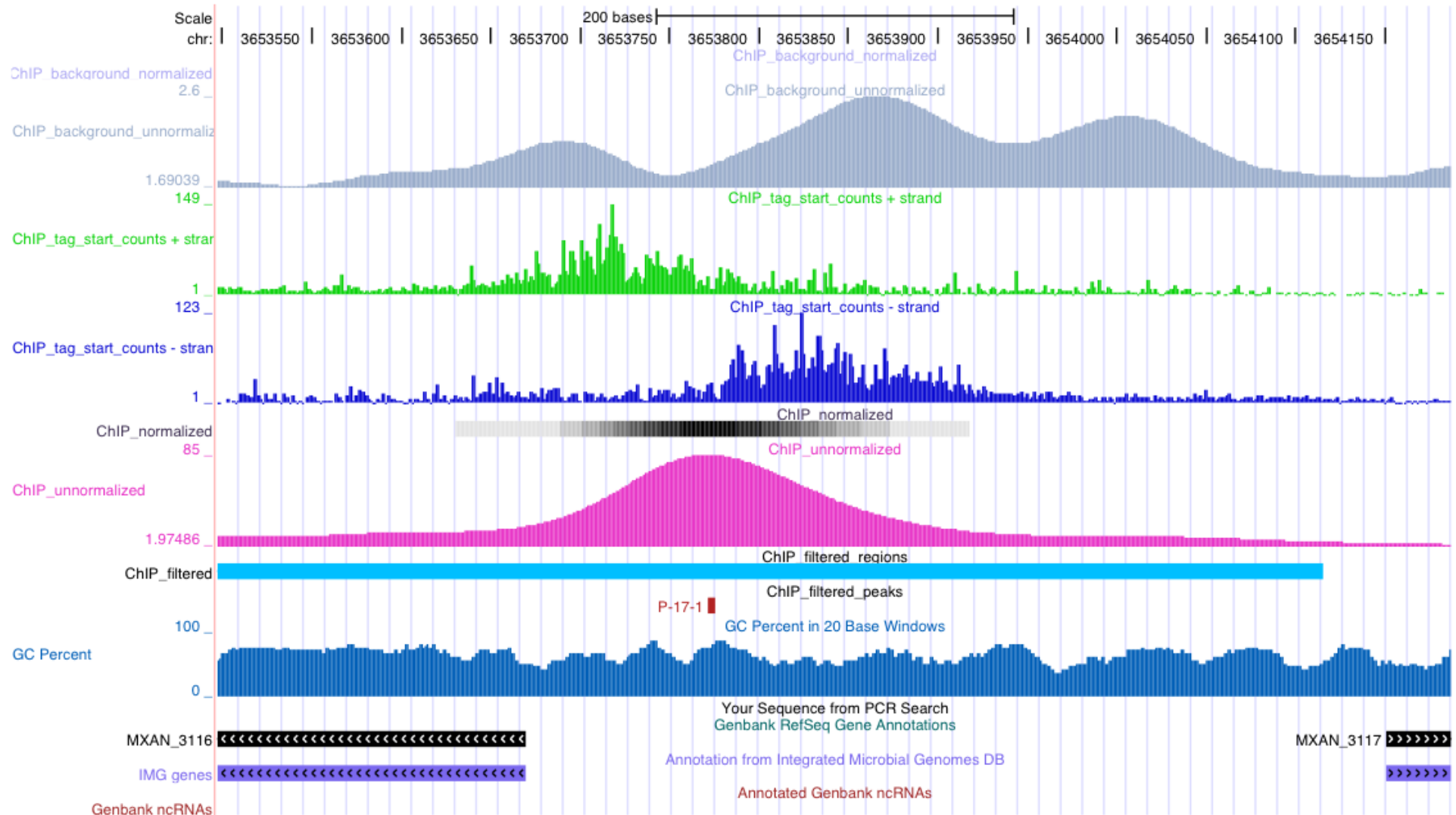




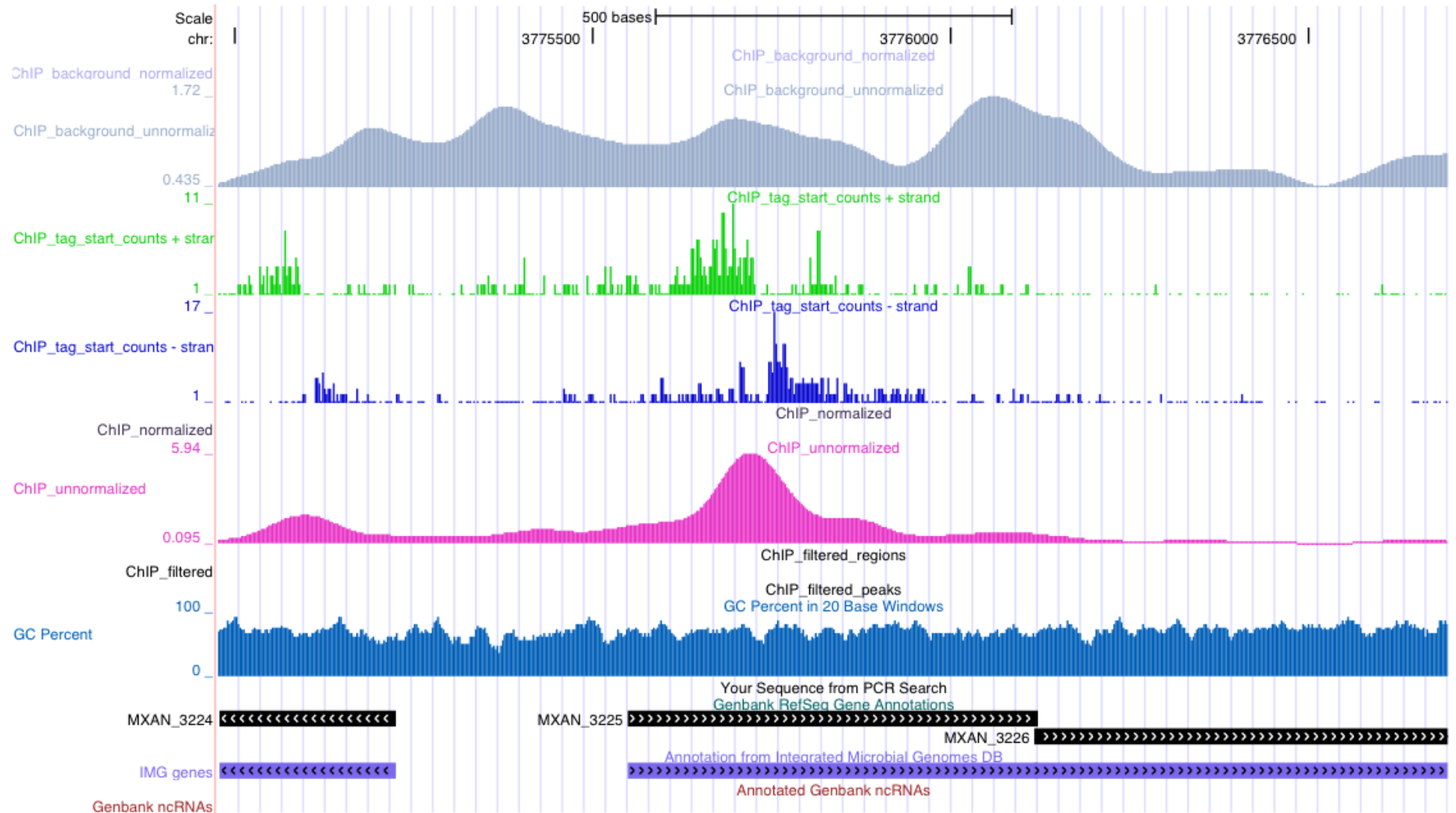
# MrpC promoter (positive control)



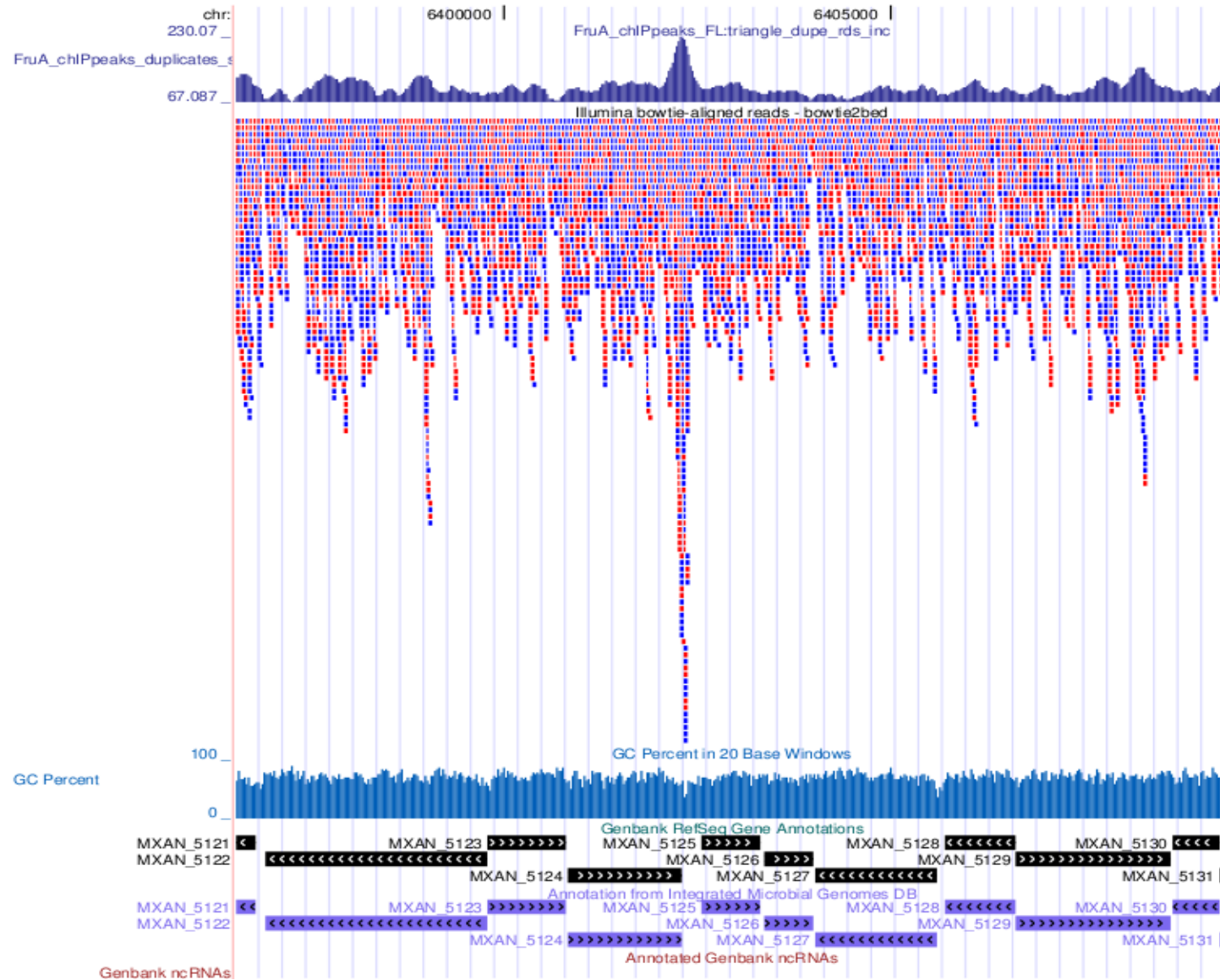
# FruA promoter (positive control)



# FdgA promoter (negative control)



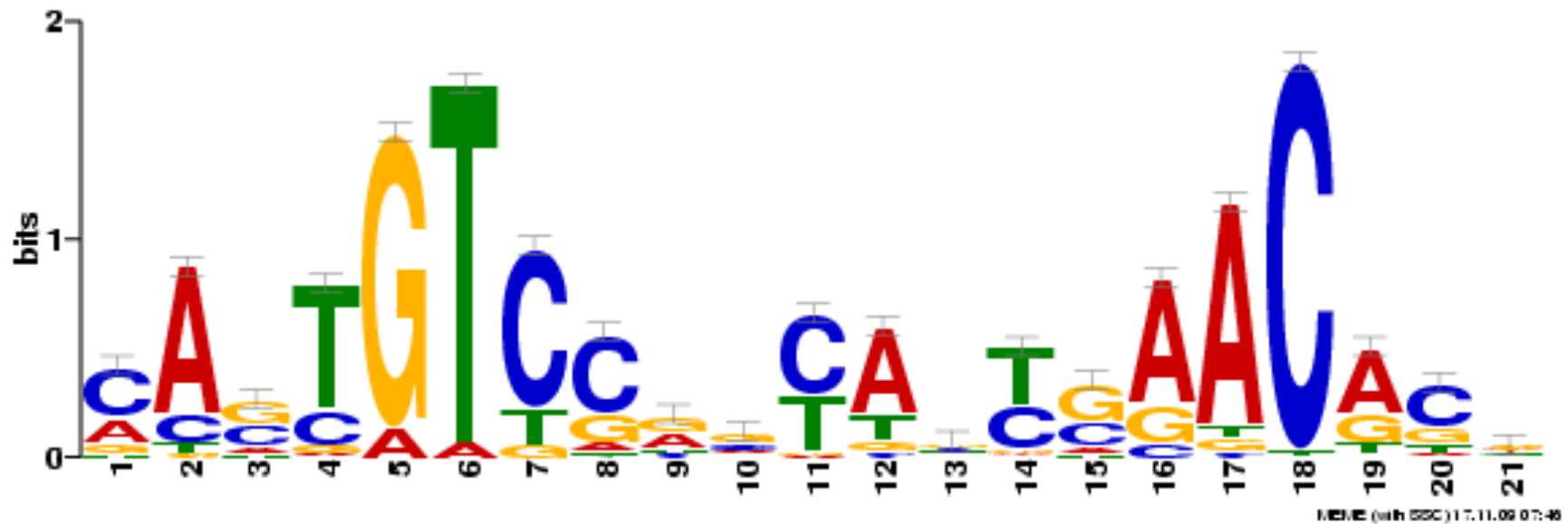
# A cautionary tale!



# Motif searching

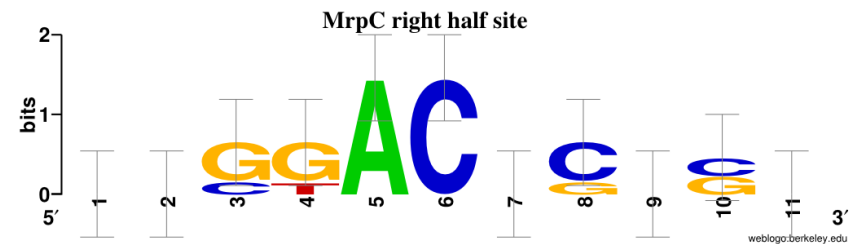
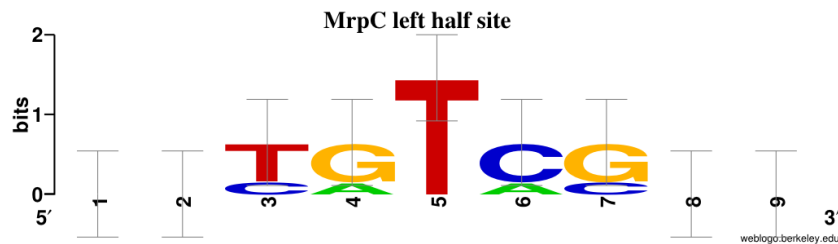
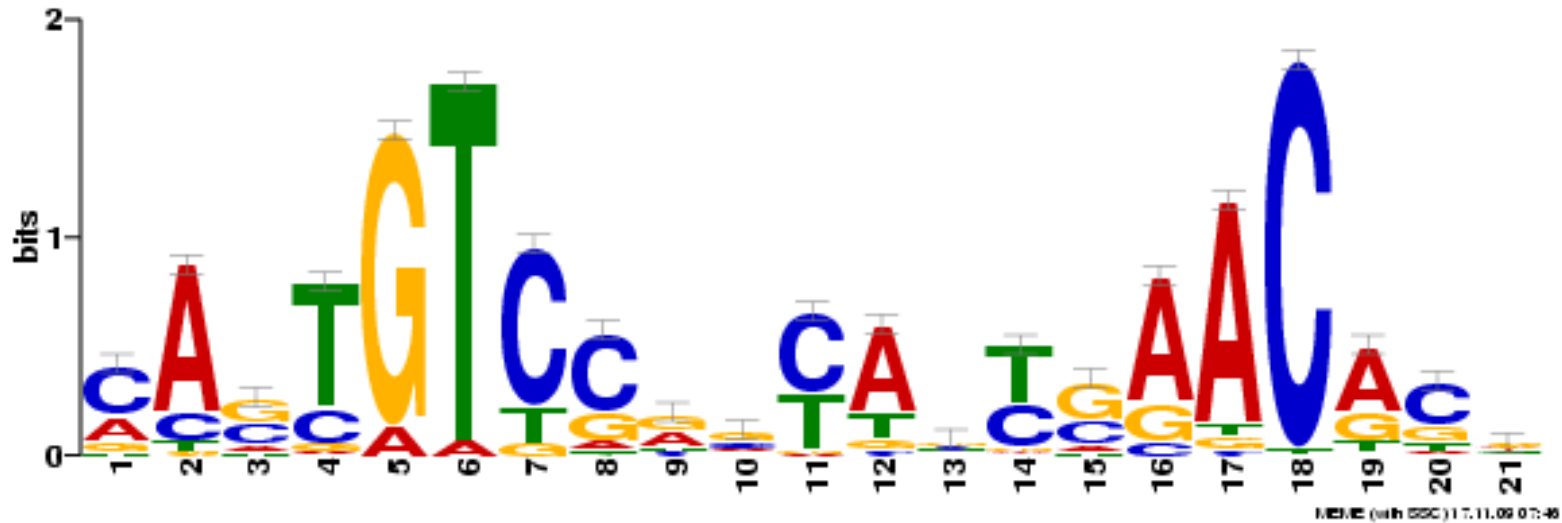
- Can we extract any common motifs from the regions immediately surrounding ChIP-enriched peaks?
  - Do any motifs found resemble MrpC binding sites?
  - Can we use this approach to determine an informative PWM for MrpC binding sites
  - Is such a PWM informative enough for genome wide searching?

Most informative motif found in all 170  
stringently called peaks

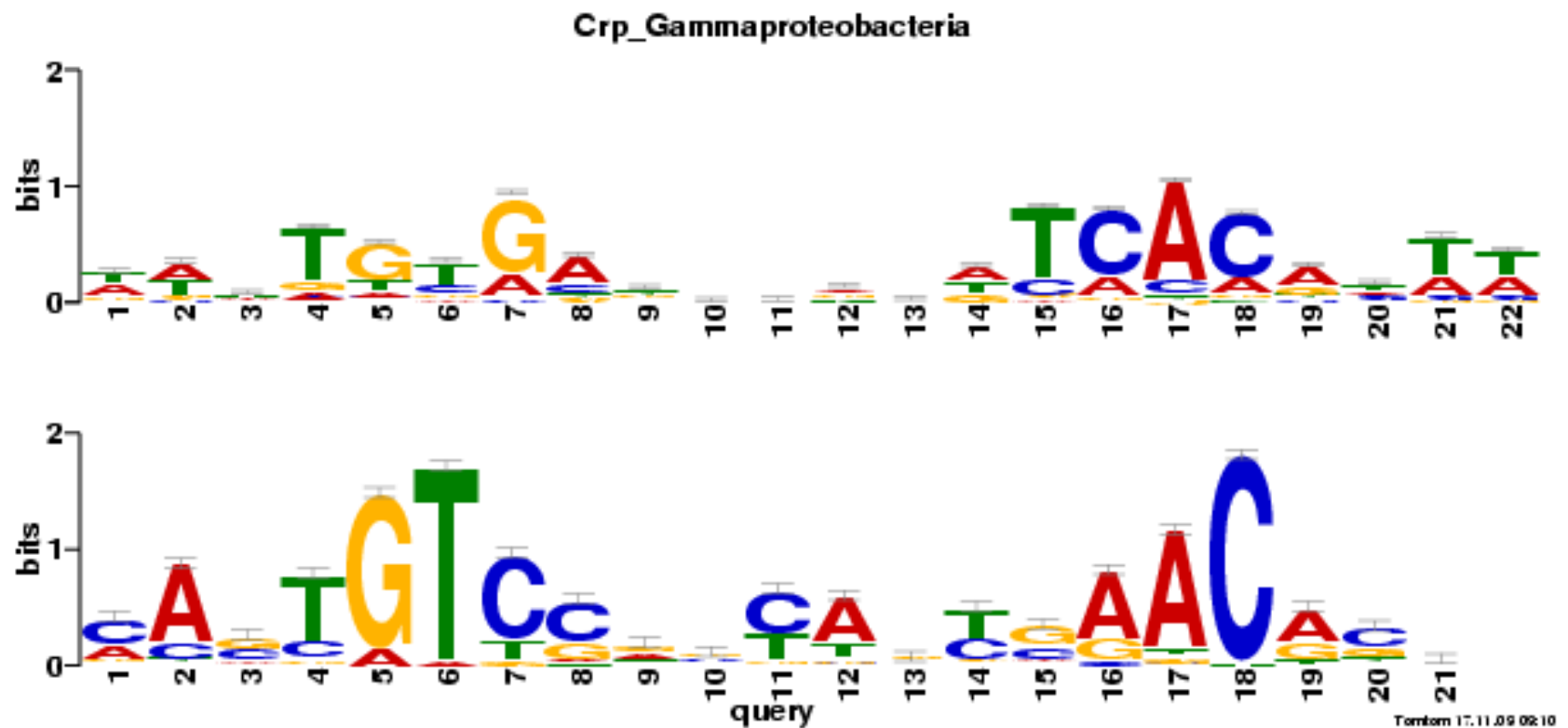




Motif closely resembles PWM  
constructed from known sites



Motif shows homology with Crp motif, a known homolog of MrpC



# Genome-wide computational predictions

- FruA: ~ 26,000 hits genome-wide (exp: 22,000)
- MrpC: ~ 11,000 hits genome-wide (exp: 12,000)
- Co-located motifs in intergenic regions:
  - Expection genome wide: 120
  - Expection restricted to intergenic regions: 12
  - 38 motifs found genome-wide
  - MrpC promoter predicted to be bound by both FruA and MrpC2



# FruA

