

Chip seq

Tools required

- i. Mapping tools- Bowtie2 or BWA ii. Sam to bam conversion – samtools iii. MACS 2 – Peak caller
- iv. Homer v. Bedtools vi. Meme- motif analysis tools

Workflow

Mapping reads to the reference genome

#used bwa

#Indexing reference genome

bwa index E.coli_BW25113.fasta

#Alignment to the reference genome

bwa aln -q 20 -t 4 E.coli_BW25113.fasta C_CHR1_R1.fastq.gz > C_CHR1_R1.sai

bwa samse E.coli_BW25113.fasta C_CHR1_R1.sai C_CHR1_R1.fastq.gz > C_CHR1_R1.sam

#Converting sam to bam

samtools view -bSq 20 C_CHR1_R1.sam -o C_CHR1_R1.bam -@4

#Sorting the bam file

samtools sort C_CHR1_R1.bam -o C_CHR1_R1_sorted.bam -@4

#Indexing the sorted bam file

samtools index C_CHR1_R1_sorted.bam [-@4](#)

```
ibab@IBAB-MSc14-Comp003:~/data_analysis/chip_Seq$ bwa aln -q 20 -t 4 E.coli_BW25113.fasta SO_4933_C_CHR1_R1.fastq.gz > C_CHR1_R1.sai
[bwa_aln] 17bp reads: max_diff = 2
[bwa_aln] 38bp reads: max_diff = 3
[bwa_aln] 64bp reads: max_diff = 4
[bwa_aln] 93bp reads: max_diff = 5
[bwa_aln] 124bp reads: max_diff = 6
[bwa_aln] 157bp reads: max_diff = 7
[bwa_aln] 190bp reads: max_diff = 8
[bwa_aln] 225bp reads: max_diff = 9
[bwa_read_seq] 0.1% bases are trimmed.
[bwa_aln_core] calculate SA coordinate... 5.84 sec
[bwa_aln_core] write to the disk... 0.02 sec
[bwa_aln_core] 262144 sequences have been processed.
[bwa_read_seq] 0.1% bases are trimmed.
[bwa_aln_core] calculate SA coordinate... 5.81 sec
[bwa_aln_core] write to the disk... 0.02 sec
[bwa_aln_core] 524288 sequences have been processed.
[bwa_read_seq] 0.1% bases are trimmed.
[bwa_aln_core] calculate SA coordinate... 5.89 sec
```

```

ibab@IBAB-MSc14-Comp003:~/data_analysis/chip_Seq$ bwa samse E.coli_BW25113.fasta
C_CHR1_R1.sai SO_4933_C_CHR1_R1.fastq.gz > C_CHR1_R1.sam
[bwa_read_seq] 0.1% bases are trimmed.
[bwa_aln_core] convert to sequence coordinate... 0.48 sec
[bwa_aln_core] refine gapped alignments... 0.09 sec
[bwa_aln_core] print alignments... 0.30 sec
[bwa_aln_core] 262144 sequences have been processed.
[bwa_read_seq] 0.1% bases are trimmed.
[bwa_aln_core] convert to sequence coordinate... 0.53 sec
[bwa_aln_core] refine gapped alignments... 0.09 sec
[bwa_aln_core] print alignments... 0.30 sec
[bwa_aln_core] 524288 sequences have been processed.
[bwa_read_seq] 0.1% bases are trimmed.
[bwa_aln_core] convert to sequence coordinate... 0.50 sec
[bwa_aln_core] refine gapped alignments... 0.09 sec
[bwa_aln_core] print alignments... 0.30 sec

```

```

#Peakcalling using MACS2
macs2 callpeak -t C_CHR1_R1_sorted.bam -c C_INR1_R1_sorted.bam -f BAM -g 4.6e6 -n
C_R1_sorted -B -p 0.001 --nomodel

```

macs2 callpeak

This is the **peak-calling function** in MACS2 — it identifies genomic regions where the ChIP sample (treatment) has significantly higher read coverage than the input (control), representing potential **protein-DNA binding sites**.

-t C_CHR1_R1_sorted.bam

Treatment file (ChIP sample)

This is your **experimental ChIP-Seq BAM file** (after sorting and indexing).
MACS2 will look for enriched regions (peaks) in this sample.

-c C_INR1_R1_sorted.bam

Control or Input file

Represents background signal — DNA fragments not specifically bound by the protein of interest. MACS2 uses this to estimate noise and correct for biases like open chromatin or GC content.

-f BAM

File format

Specifies that your input files (-t and -c) are in **BAM format** (binary alignment files).

-g 4.6e6

Genome size

Approximate number of **mappable bases** in your reference genome.

You used 4.6e6, which is correct for *E. coli* (~4.6 Mb genome).
For humans, you'd typically use 2.7e9, and for mouse 1.87e9.

-n C_R1_sorted

Output name prefix

All result files will start with this prefix (e.g., C_R1_sorted_peaks.xls, C_R1_sorted_summits.bed, C_R1_sorted_treat_pileup.bdg).

-B

Generate signal tracks

Tells MACS2 to output **BedGraph** files:

- `_treat_pileup.bdg` → ChIP signal track
- `_control_lambda.bdg` → background model

These can be visualized in genome browsers like IGV or UCSC.

-p 0.001

p-value cutoff for peak detection

Only regions with a **p-value** < **0.001** are considered significant peaks.
(You can use `-q 0.05` instead for **FDR-based** cutoff if you prefer.)

--nomodel

Disables automatic model building

Normally MACS2 tries to infer fragment size from the data, but for bacterial or low-coverage ChIP-Seq (like *E. coli*), this can fail.

`--nomodel` prevents model building and uses raw read coverage directly.

```
ibab@IBAB-MSc14-Comp003:~/data_analysis/chip_Seq$ macs3 callpeak -t D_CHR1_R1_sorted.bam -c D_INR1_R1_sorted.bam -f BAM -g 4.6e6 -n D_R1_sorted -B -p 0.001 --nomodel
INFO @ 15 Oct 2025 11:04:55: [39 MB]
# Command line: callpeak -t D_CHR1_R1_sorted.bam -c D_INR1_R1_sorted.bam -f BAM -g 4.6e6 -n D_R1_sorted -B -p 0.001 --nomodel
# ARGUMENTS LIST:
# name = D_R1_sorted
# format = BAM
# ChIP-seq file = ['D_CHR1_R1_sorted.bam']
# control file = ['D_INR1_R1_sorted.bam']
# effective genome size = 4.60e+06
# band width = 300
# model fold = [5, 50]
# pvalue cutoff = 1.00e-03
# qvalue will not be calculated and reported as -1 in the final output.
```

1. narrowpeak with the peak info
2. bedgraph file for visualization
3. summit for the info of the highest intensity region of the peak

#Peak annotation using Homer

annotatePeaks.pl C_R1_sorted_peaks.narrowPeak E.coli_BW25113.fasta -gff

E.coli_BW25113_annotation.gff> C_R1_sorted_annotated_peaks.tsv

```
D_INR1_R1.sai E.coli_BW25113.fasta.pac
D_INR1_R1.sam E.coli_BW25113.fasta.sa
D_INR1_R1_sorted.bam S0_4933_C_CHR1_R1.fastq.gz
D_INR1_R1_sorted.bam.bai S0_4933_C_INR1_R1.fastq.gz
D_R1_sorted_annotated_peaks.tsv S0_4933_D_CHR1_R1.fastq.gz
D_R1_sorted_control_lambda.bdg S0_4933_D_INR1_R1.fastq.gz
D_R1_sorted_peaks.narrowPeak
ibab@IBAB-MSc14-Comp003:~/data_analysis/chip_Seq$ cat D_R1_sorted_peaks.narrowPeak | wc -l
912
ibab@IBAB-MSc14-Comp003:~/data_analysis/chip_Seq$ samtools faidx E.coli_BW25113.fasta
cut -f1,2 E.coli_BW25113.fasta.fai > genome.chrom.sizes
E.coli_BW25113_annotation.gff> C_R1_sorted_annotated_peaks.
```

```
ibab@IBAB-MSc14-Comp003:~/data_analysis/chip_Seq$ annotatePeaks.pl D_R1_sorted_peaks.narrowPeak E.coli_BW25113.fasta -gff E.coli_BW25113_annotation.gff> D_R1_sorted_annotated_peaks.tsv
```

```
Using Custom Genome
Peak file = D_R1_sorted_peaks.narrowPeak
Genome = E.coli_BW25113.fasta
Organism = unknown
Custom annotation GFF file: E.coli_BW25113_annotation.gff (better to get GTF file)
```

```
Peak/BED file conversion summary:
BED/Header formatted lines: 912
peakfile formatted lines: 0
Duplicated Peak IDs: 0
```

```
Peak File Statistics:
```

2145414+	847 NA	Intergenic
1816167+	846 NA	Intergenic
281199+	846 NA	Intergenic
1541144+	844 NA	Intergenic
1252307+	843 NA	Intergenic
1578691+	842 NA	Intergenic
2093487+	838 NA	Intergenic
2051719+	836 NA	TSS (ID=exon-IVL04_RS10130-1;Parent=rna-IVL04_RS10130;anticodon=(pos:com
2227770+	835 NA	Intergenic
1526428+	833 NA	Intergenic
1847617+	832 NA	Intergenic
1668343+	831 NA	Intergenic
2378236+	828 NA	Intergenic
2198255+	828 NA	Intergenic
1092385+	828 NA	TSS (ID=exon-IVL04_RS05230-1;Parent=rna-IVL04_RS05230;Dbxref=RFAM:RF00
1807762+	826 NA	Intergenic
2853874+	825 NA	Intergenic
560412+	823 NA	promoter-TSS (ID=exon-IVL04_RS02710-1;Parent=rna-IVL04_RS02710;anticodon=
1582396+	823 NA	Intergenic
1337943+	821 NA	Intergenic
675294+	820 NA	Intergenic

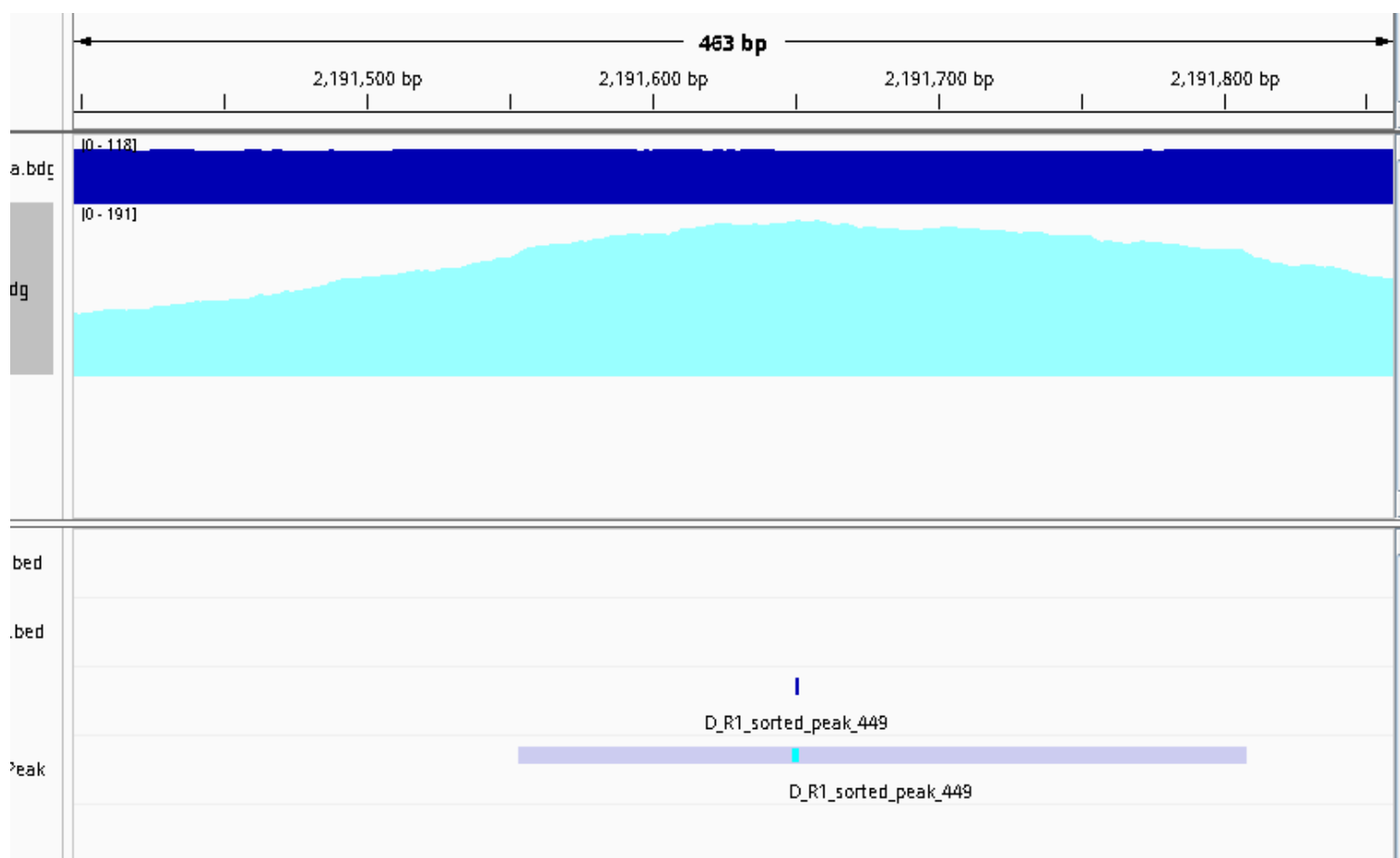
Merging the peaks

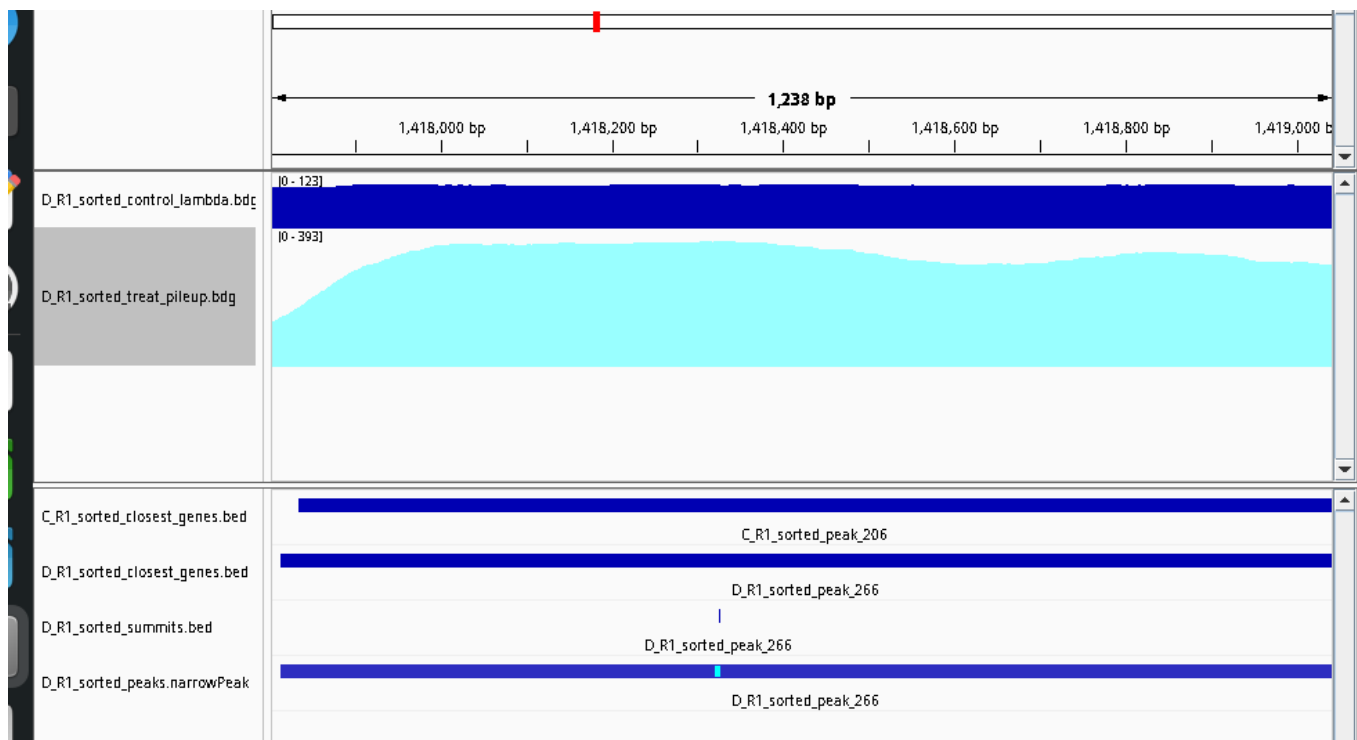
```
mergePeaks D_R1_sorted_peaks.narrowPeak>D_R1_sorted_peaks_merged.txt
```

Finding peaks nearer the features like gene

```
closestBed -a D_R1_sorted_peaks.narrowPeak -b E.coli_BW25113_annotation.bed
```

The results are viewed in IGV





Extract the sequences from the reference file that corresponds to the narrow peak-calling

Using bed tools we can extract the sequences corresponding to peaks

```
bedtools getfasta -fi E.coli_BW25113.fasta -bed D_R1_sorted_peaks.narrowPeak -fo D_R1_sorted_peaks.fa
```

```
ibab@IBAB-MSc14-Comp003:~/data_analysis/chip_Seq$ bedtools getfasta -fi E.coli_BW25113.fasta -
bed D_R1_sorted_peaks.narrowPeak -fo D_R1_sorted_peaks.fa
ibab@IBAB-MSc14-Comp003:~/data_analysis/chip_Seq$ cat D_R1_sorted_peaks.fa
>NZ_CP064677.1:4953-5186
CTTGCTTTACATAATCTGCCGCCGATTTTGCTGCGTTGCGTAAATTGATGATGAATCATCAGTAAAATCTATTCATTATCTCAATCAGGCCG
GGTTTGCTTTATGCAGCCCGGCTTTTATGAAGAAATTATGGAGAAAAATGACAGGGAAAAAGGAGAAATCTCAATAAATGCGGTAACCTA
GAGATTAGGATTGCGGAGAATAACAACCGCGTTCTCATCGAGTA
>NZ_CP064677.1:10395-10612
AGCGAAATAACCCACGTTGTGACAGTTAAGCAGAATGGTGGTCATGCCGAAGCCCATCAGGCCAGCGGTGCCGATTAGCCAACCTAGTGTTG
CCCATAATTCCTCAAAAATCATCATCGAATGAATGGTGAAATAATTTCCCTGAATAACTGTAGTGTTCAGGGCGCGGCATAATAATCAGCCA
GTGGGCGAGTGCTACGATCTTTTGAGGG
>NZ_CP064677.1:11093-11553
AATTTACCGTGTCCGCGCAGTTTGTGGCGATACTATCGCCACCAAAATGCTGTAATTCCTCGGCAATCAGCTGCCAGTTGCGGCGATGTTGCT
CGGGATGCCCTTCATCGATTTAAACAGTTCGTTGCGCATCAGTACGCTGGAGAGGCGAGTTTGCCTTTTTCATTATGGGTGAGCAATCGGGC
GAAATTTGCCAACTGTTCTCACTACAATGCTGAAGAAAATCCAGATCTGAATCATTACAGTAATTAACATTCATTTTTGTGGCTTCTATATT
CTGGCGTTAGTCGTCGCCGATAATTTTCAGCGTGGCCATATCCGATGAGTTCACCGTATGACCCGAAAAGGTGATTTTTGAGACGCAGCGTTTA
TTGTCGTTATCGCTGTTAATGTTGATCCAGTCAGTGGTTTGCCCTTCTTTATTTCTGAAGGAATATTCAGGCTCTGACTGGCG
>NZ_CP064677.1:11871-12160
TTCTTTTATTACAATTTTTGTGTAATGCCTTGGCTGCGATTCAATCTTTATATGAATAAAATTGCTGTCAATTTTACGTCTTGTCTGCCAT
ATCGCGAAATTTCTGCGCAAAAGCACAAAAATTTTTGCATCTCCCCCTTGATGACGTGGTTTACGACCCCATTTAGTAGTCAACCGCAGTGAG
TGAGTCTGCAAAAAATGAAATGGGCAGTTGAAACCAGACGTTTCGCCCTATTACAGACTCACAAACCATGATGACCGAATATATAGTGA
GACGTTT
>NZ_CP064677.1:16961-17609
TAGGAGGCCTCGGGTTGATGGTAAATATCACTCGGGGCTTTTCTCTATCTGCCGTTACAGTAATGCCTGAGACAGACAGCCTCAAGCACCCGC
CGCTATTATATCGCTCTCTTTAACCCATTTTGTGTTATCGATTCTAATCCTGAAGACGCCTCGCATTTTGTGGCGTAATTTTTTAATGATTTA
ATTATTTAATTTAATTTATCTCTTCATCGCAATTATTGACGACAAGCTGGATTATTTTGAATATTGGCCTAACAGCATCGCCGACTGACA
```

```

CATTTTTAAACAACTCAACCGTTAGTACAGTCAGGAAATAGTTTAGCCTTTTTTAAGCTAAGTAAAGGGCTTTTCTGCGACTTACGTTAAGAA
TTTGTAATTTCGCACCGCGTAATAAGTTGACAGTGATCACCCGTTTCGCGTTATTTGATCAAGAAGAGT
ibab@IBAB-MSc14-Comp003:~/data_analysis/chip_Seq$ cat D_R1_sorted_peaks.fa |wc -l
1824
ibab@IBAB-MSc14-Comp003:~/data_analysis/chip_Seq$

```

There are 1824 sequences

Motif analysis using Hommer

```

findMotifs.pl D_R1_sorted_peaks.fa fasta homer_motif_output/ -len 8,10,12 -norevopp
D_R1_sorted_peaks.fa → input sequences

```

fasta → input type

homer_motif_output/ → output directory

```

ibab@IBAB-MSc14-Comp003:~/data_analysis/chip_Seq$ findMotifs.pl D_R1_sorted_peaks.fa fasta hom
er_motif_output/ -len 8,10,12 -norevopp

Selected Options:
  Input file = D_R1_sorted_peaks.fa
  Promoter Set = fasta
  Output Directory = homer_motif_output/
  Motif length set at 8, 10, 12,
  Will not search the reverse strand

!Warning - no background FASTA file specified (Highly recommended)
!Your input sequences will be randomized to serve as a background instead.

Found 912 sequences
Using custom gene IDs for GO analysis
Parsing FASTA format files...
Found 912 sequences
!! 1 of 912 contained bad nucleotide characters [not ACGTN], replaced with N
Found 4560 sequences
!! 5 of 4560 contained bad nucleotide characters [not ACGTN], replaced with N

```

-len 8,10,12 → motif lengths to search

-norevopp → only forward strand (optional)

HOMER outputs:

Known motifs enriched in your peaks

De novo motifs discovered

Motif logos and statistics

Phylogenetic analysis using MAFFT and FastTree

mafft alanine.fasta>align_seq.fasta

```
--dash :          Add structural information (Rozewicki et al, submitted)
ibab@IBAB-MSc14-Comp003:~/data_analysis/phylo$ mafft alanine.fasta>align_seq.fasta
nthread = 0
nthreadpair = 0
nthreadtb = 0
ppenalty_ex = 0
stacksize: 8192 kb
generating a scoring matrix for nucleotide (dist=200) ... done
Gap Penalty = -1.53, +0.00, +0.00

Making a distance matrix ..
    1 / 63
done.
```

```
ibab@IBAB-MSc14-Comp003:~/data_analysis/phylo$ cat align_seq.fasta
>Saccharomyces_cerevisiae
-ggggttatagttaaatttggtagaacgactgcggttgcatttaatatgagttcaag
tctcattaact-----
>Saccharomyces_weihenstephan
-ggggttatagttcaatttggtagaacgactgcggttgcatttaatatgagttcaag
tctcattaactcca---
>Saccharomyces_paradoxus
-ggggttatagttaaatttggtagaacgattgcggttgcatttaatatgagttcaag
tctcattaactccaata
>Saccharomyces_mikatae
-ggggttatagttaaatttggtagaacgactgcggttgcatttaatatgagttcaag
tctcattaactccaata
>Saccharomyces_uvarum
-ggggttatagttcaatttggtagaacgactgcggttgcatttaatatgagttcaag
tctcattaactcca---
>Saccharomyces_eubayanus
aggggttatagttcaatttggtagaacgactgcggttgcatttaatatgagttcaag
tctcattaactccaata
>Saccharomyces_pastorianus
```


Synteny using mummer

nucmer --prefix=alignment alanine.fasta glycine.fasta

```
Try '/usr/bin/nucmer -h' for more information.
ibab@IBAB-MSc14-Comp003:~/data_analysis/phylo$ nucmer --prefix=alignment alanine
.fasta glycine.fasta
1: PREPARING DATA
2,3: RUNNING mummer AND CREATING CLUSTERS
# reading input file "alignment.ntref" of length 4621
# construct suffix tree for sequence of length 4621
# (maximum reference length is 536870908)
# (maximum query length is 4294967295)
# CONSTRUCTIONTIME /usr/bin/mummer alignment.ntref 0.00
# reading input file "/home/ibab/data_analysis/phylo/glycine.fasta" of length 45
25
# matching query-file "/home/ibab/data_analysis/phylo/glycine.fasta"
# against subject-file "alignment.ntref"
# COMPLETETIME /usr/bin/mummer alignment.ntref 0.00
# SPACE /usr/bin/mummer alignment.ntref 0.01
4: FINISHING DATA
```

Synten plot was not able to generate the time.