# WORKSHOP 1: <u>Analysis of genome-wide transcription factor-binding sites using ChIP-seq data</u>

## TASK1: Create an environment using mamba, and name the environment "workshop1"
Address: 10.139.1.132 Username: binf6_03 Password: binf-GiantDingo76
mamba create -n workshop1 python=3.7
mamba activate workshop1
mamba install -c bioconda fastqc seqtk bowtie2 samtools macs2 idr

## TASK2: PIPELINE

```
binf6_03@binf6000-03-32: ~
GNU nano 6.2                                                                                    pipe.
#!/bin/bash

# Define a new output directory for the full analysis
FULL_OUTPUT_DIR="full_output"

# Step 1: Quality Control
echo "Running FastQC for quality control on the full dataset..."
mkdir -p ${FULL_OUTPUT_DIR}/metrics # Create a metrics directory inside the full_output directory
fastqc data/*.fastq.gz -o ${FULL_OUTPUT_DIR}/metrics/ # Run FastQC on all fastq.gz files in the data directory and output the results to the metrics directory
```

```
# Step 2: Alignment with Bowtie2
echo "Aligning reads with Bowtie2 for the full dataset..."
mkdir -p ${FULL_OUTPUT_DIR}/alignment # Creates a directory to store alignment outputs.
# Align the control sample to the reference genome
bowtie2 -x reference/bowtie2_index/genome \ # Specifies the path to the pre-built Bowtie2 index for the reference genome.
        -U data/control.fastq.gz \ # Specifies the path to the FASTQ file for the control sample.
        > ${FULL_OUTPUT_DIR}/alignment/control.sam # Directs the SAM formatted alignment output to a file named control.sam in the alignment directory.
# Align treatment replicate 1 to the reference genome
bowtie2 -x reference/bowtie2_index/genome \ # Specifies the path to the Bowtie2 index.
        -U data/treatment_rep1.fastq.gz \ # Specifies the path to the FASTQ file for treatment replicate 1.
        > ${FULL_OUTPUT_DIR}/alignment/treatment_rep1.sam # Directs the SAM formatted alignment output for treatment replicate 1 to its file.
# Align treatment replicate 2 to the reference genome
bowtie2 -x reference/bowtie2_index/genome \ # Specifies the path to the Bowtie2 index for the reference genome.
        -U data/treatment_rep2.fastq.gz \ # Specifies the path to the FASTQ file for treatment replicate 2.
        > ${FULL_OUTPUT_DIR}/alignment/treatment_rep2.sam # Directs the SAM formatted alignment output for treatment replicate 2 to its file.
```

```
# Step 3: SAM to BAM Conversion and Mapping Quality Check
echo "Converting SAM to BAM and checking mapping quality for the full dataset..."
# Convert SAM file to BAM for the control sample
samtools view -bS ${FULL_OUTPUT_DIR}/alignment/control.sam \
        > ${FULL_OUTPUT_DIR}/alignment/control.bam # Uses samtools to convert the SAM format file of the control sample to BAM format for efficient storage and analysis.
# Convert SAM file to BAM for treatment replicate 1
samtools view -bS ${FULL_OUTPUT_DIR}/alignment/treatment_rep1.sam \
        > ${FULL_OUTPUT_DIR}/alignment/treatment_rep1.bam # Converts the SAM format file of treatment replicate 1 to BAM format.
# Convert SAM file to BAM for treatment replicate 2
samtools view -bS ${FULL_OUTPUT_DIR}/alignment/treatment_rep2.sam \
        > ${FULL_OUTPUT_DIR}/alignment/treatment_rep2.bam # Converts the SAM format file of treatment replicate 2 to BAM format.
# Generate mapping quality report for the control sample
samtools flagstat ${FULL_OUTPUT_DIR}/alignment/control.bam \
        > ${FULL_OUTPUT_DIR}/metrics/control_mapping_quality.txt # Generates a statistical report on the mapping quality of the control sample, saving it to a text file.
# Generate mapping quality report for treatment replicate 1
samtools flagstat ${FULL_OUTPUT_DIR}/alignment/treatment_rep1.bam \
        > ${FULL_OUTPUT_DIR}/metrics/treatment_rep1_mapping_quality.txt # Generates a statistical report on the mapping quality of treatment replicate 1, saving it to a text file.
# Generate mapping quality report for treatment replicate 2
samtools flagstat ${FULL_OUTPUT_DIR}/alignment/treatment_rep2.bam \
        > ${FULL_OUTPUT_DIR}/metrics/treatment_rep2_mapping_quality.txt # Generates a statistical report on the mapping quality of treatment replicate 2, saving it to a text file.
```

```
(workshop1) binf6_03@binf6000-03-32:~$ cd full_output/metrics
(workshop1) binf6_03@binf6000-03-32:~/full_output/metrics$ ls
control_fastqc.html  control_fastqc.zip  control_mapping_quality.txt  treatment_rep1_fastqc.html  treatment_rep1_fastqc.zip  treatment_rep1_mapping_quality.txt  treatment_rep2_fastqc.html  treatment_rep2_fastqc.zip  treatment_rep2_mapping_quality.txt
(workshop1) binf6_03@binf6000-03-32:~/full_output/metrics$ cd
(workshop1) binf6_03@binf6000-03-32:~$ cd full_output/alignment
(workshop1) binf6_03@binf6000-03-32:~/full_output/alignment$ ls
control.bam  control.sam  treatment_rep1.bam  treatment_rep1.sam  treatment_rep2.bam  treatment_rep2.sam
```

```
# Step 4: Peak Calling with MACS2
echo "Calling peaks with MACS2 using fixed extension size for the full dataset..."
mkdir -p ${FULL_OUTPUT_DIR}/peaks # Create a directory for storing the peak calling results.
# Calling peaks for treatment replicate 1 against control
macs2 callpeak -t ${FULL_OUTPUT_DIR}/alignment/treatment_rep1.bam \ # Specify the treatment replicate 1 BAM file as the target.
        -c ${FULL_OUTPUT_DIR}/alignment/control.bam \ # Specify the control BAM file for comparison.
        -f BAM -g dm -n treatment_rep1 \ # Define the format as BAM, specify the genome size (dm for Drosophila melanogaster), and output name.
        --outdir ${FULL_OUTPUT_DIR}/peaks -q 0.05 \ # Set the output directory for peak files and set the q-value (FDR-adjusted p-value) cutoff for peak detection.
                # The q-value represents the minimum FDR at which the test may be called significant. A q-value of 0.05 means that you're willing to accept a 5% of the identified peaks could be false positives.
        --call-summits --nomodel --extsize 147 # Call peak summits, bypass the model building step, and set a fixed extension size of 147 bp to simplify analysis.
# Calling peaks for treatment replicate 2 against control
macs2 callpeak -t ${FULL_OUTPUT_DIR}/alignment/treatment_rep2.bam \ # Specify the treatment replicate 2 BAM file as the target.
        -c ${FULL_OUTPUT_DIR}/alignment/control.bam \ # Specify the control BAM file for comparison.
        -f BAM -g dm -n treatment_rep2 \ # Define the format as BAM, specify the genome size (dm for Drosophila melanogaster), and name the output.
        --outdir ${FULL_OUTPUT_DIR}/peaks -q 0.05 \ # Set the output directory for peak files and the q-value cutoff for peak detection.
        --call-summits --nomodel --extsize 147 # Call peak summits, bypass the model building step, and set a fixed extension size of 147 bp.
```

```
chr2L   887987  887988  treatment_rep1_peak_83  8.40418
chr2L   1006775 1006776 treatment_rep1_peak_84  117.643
chr2L   1063014 1063015 treatment_rep1_peak_85  3.49789
chr2L   1072389 1072390 treatment_rep1_peak_86  31.1469
chr2L   1074724 1074725 treatment_rep1_peak_87  6.24684
treatment_rep1_summits.bed
```

```
chr2L   885175  885176  treatment_rep2_peak_83  13.8
chr2L   887972  887973  treatment_rep2_peak_84  9.70
chr2L   901587  901588  treatment_rep2_peak_85  3.06
chr2L   1006773 1006774 treatment_rep2_peak_86  82.5
chr2L   1072408 1072409 treatment_rep2_peak_87  31.2
treatment_rep2_summits.bed
```

```bash
# Step 5: Reproducibility Assessment with IDR
echo "Assessing reproducibility with IDR for the full dataset..."
mkdir -p ${FULL_OUTPUT_DIR}/idr # Create a directory for storing idr results.
# Run IDR (Irreproducible Discovery Rate) analysis to compare peak calls between the two treatment replicates
idr --samples ${FULL_OUTPUT_DIR}/peaks/treatment_rep1_peaks.narrowPeak \ # Specify first set of peaks for comparison
     ${FULL_OUTPUT_DIR}/peaks/treatment_rep2_peaks.narrowPeak \ # Specify second set of peaks for comparison
    --output-file ${FULL_OUTPUT_DIR}/idr/treatment_idr_output.txt \ # Designate file to write IDR analysis results
    --plot \ # Generate plots summarizing the IDR analysis
    --log-output-file ${FULL_OUTPUT_DIR}/idr/treatment_idr.log # Log file for detailed IDR analysis output
# Convert p-value to -log10(p-value) in column 12
awk 'BEGIN {OFS="\t"} {$12 = -log($12)/log(10); print}' ${FULL_OUTPUT_DIR}/idr/treatment_idr_output.txt > ${FULL_OUTPUT_DIR}/idr/treatment_idr_output_log10p>


# Step 6: Filtering Peaks Based on IDR Output
echo "Filtering significant, reproducible peaks based on IDR output for the full dataset..."
mkdir -p ${FULL_OUTPUT_DIR}/final_peaks # Create a directory to hold the final set of peaks deemed significant and reproducible.
# Filter peaks based on p-value (in column 12)
awk '$12 < 0.05' ${FULL_OUTPUT_DIR}/idr/treatment_idr_output.txt \
    > ${FULL_OUTPUT_DIR}/final_peaks/significant_reproducible_peaks.bed #Final output destination.
```

**TASK3: Identify a gene bound by the transcription factor CTCF**

less significant_reproducible_peaks.bed.



Given these locations a genome annotation file recording gene locations (dm6_tss.bed), and a getclosestgene.py script, I can identify a gene which has a CTCF binding site in its promoter region.

```
$ python binfpy/getclosestgene.py home/binf6_03/significant_reproducible_peaks.bed
dm6_tss.bed
```

This will output a bed file `tss_gene.bed` detailing:

```
<CTCF bind chrom> <CTCF bind start> <CTCF bind end> <Gene name> <distance between CTCF
and gene> <strand>
```

less `tss_gene.bed`

chr2L is the chromosome where the binding site is located. "chr2L" refers to the left arm of the second chromosome in Drosophila melanogaster. 6498467 is the start position of the CTCF binding site on the chromosome. BED file coordinates are 0-based, meaning that the first base of the chromosome is considered position 0. Thus, this binding site starts at the 6,498,467th base of chr2L. 6498838 is the end position of the CTCF binding site on the chromosome. In BED format, the end position is exclusive, meaning the actual binding site extends up to but does not include this position. Therefore, the binding site spans from base 6,498,467 to base 6,498,837, making it 371 bases long. NM_001298748.1_up_1_chr2L_6498646_f is a unique identifier for the binding site or the peak. It includes the gene name with which the site is associated, in this case, "NM_001298748.1", which could be a gene identifier in a specific database. The additional details (up_1_chr2L_6498646_f) provide context about the binding site's location, such as it being upstream of the gene, its chromosome, a specific base position, and the direction ("f" for forward strand). The "0" represents the score of the peak or binding site, which can indicate the strength or confidence in the site's identification. A score of "0" might suggest a default value in this context. The "+" means that the binding site is on the forward strand, which has implications for the directionality of any genes or regulatory elements associated with this site.

Discovering a CTCF binding site located on the left arm of chromosome 2L, at the precise coordinates of 6,498,467 to 6,498,837, unveils a fascinating glimpse into the intricate regulatory networks within Drosophila melanogaster. This sequence lies in proximity to the gene tagged as NM_001298748.1. Given CTCF's renowned role as an architectural protein, shaping the 3D organization of chromatin and dictating the rhythm of gene expression, this association is more than mere coincidence. It suggests a targeted regulatory influence, where CTCF could be modulating the expression of NM_001298748.1, thereby impacting fundamental biological processes from development and cell differentiation to the safeguarding of chromosomal architecture. This discovery not only highlights the complexity and precision of genetic regulation in D. melanogaster but also opens up avenues for exploring how such regulatory mechanisms contribute to the organism's biology and evolution.

```bash
#!/bin/bash

FULL_OUTPUT_DIR="full_output"

# Step 1: Quality Control

mkdir -p ${FULL_OUTPUT_DIR}/metrics

fastqc data/*.fastq.gz -o ${FULL_OUTPUT_DIR}/metrics/

# Step 2: Alignment with Bowtie2

mkdir -p ${FULL_OUTPUT_DIR}/alignment.

bowtie2 -x reference/bowtie2_index/genome \
    -U data/control.fastq.gz \ # Specifies the path to the FASTQ file for the control sample.
    > ${FULL_OUTPUT_DIR}/alignment/control.sam

bowtie2 -x reference/bowtie2_index/genome \ # Specifies the path to the Bowtie2 index.
    -U data/treatment_rep1.fastq.gz \
    > ${FULL_OUTPUT_DIR}/alignment/treatment_rep1.sam

bowtie2 -x reference/bowtie2_index/genome \
    -U data/treatment_rep2.fastq.gz \ # Specifies the path to the FASTQ file for treatment replicate 2.
    > ${FULL_OUTPUT_DIR}/alignment/treatment_rep2.sam

# Step 3: SAM to BAM Conversion and Mapping Quality Check

samtools view -bS ${FULL_OUTPUT_DIR}/alignment/control.sam \
    > ${FULL_OUTPUT_DIR}/alignment/control.bam

samtools view -bS ${FULL_OUTPUT_DIR}/alignment/treatment_rep1.sam \
    > ${FULL_OUTPUT_DIR}/alignment/treatment_rep1.bam

samtools view -bS ${FULL_OUTPUT_DIR}/alignment/treatment_rep2.sam \
    > ${FULL_OUTPUT_DIR}/alignment/treatment_rep2.bam

samtools flagstat ${FULL_OUTPUT_DIR}/alignment/control.bam \
    > ${FULL_OUTPUT_DIR}/metrics/control_mapping_quality.txt

samtools flagstat ${FULL_OUTPUT_DIR}/alignment/treatment_rep1.bam \
    > ${FULL_OUTPUT_DIR}/metrics/treatment_rep1_mapping_quality.txt

samtools flagstat ${FULL_OUTPUT_DIR}/alignment/treatment_rep2.bam \
    > ${FULL_OUTPUT_DIR}/metrics/treatment_rep2_mapping_quality.txt

# Step 4: Peak Calling with MACS2

mkdir -p ${FULL_OUTPUT_DIR}/peaks

macs2 callpeak -t ${FULL_OUTPUT_DIR}/alignment/treatment_rep1.bam \ # Specify the treatment replicate 1 BAM file as the target.
        -c ${FULL_OUTPUT_DIR}/alignment/control.bam \
        -f BAM -g dm -n treatment_rep1 \
        --outdir ${FULL_OUTPUT_DIR}/peaks -q 0.05
        --call-summits --nomodel --extsize 147

macs2 callpeak -t ${FULL_OUTPUT_DIR}/alignment/treatment_rep2.bam \
```

```
        -c ${FULL_OUTPUT_DIR}/alignment/control.bam \

        -f BAM -g dm -n treatment_rep2 \

        --outdir ${FULL_OUTPUT_DIR}/peaks -q 0.05 \

        --call-summits --nomodel --extsize 147
```

# Step 5: Reproducibility Assessment with IDR

```
mkdir -p ${FULL_OUTPUT_DIR}/idr # Create a directory for storing idr results.

idr --samples ${FULL_OUTPUT_DIR}/peaks/treatment_rep1_peaks.narrowPeak \

      ${FULL_OUTPUT_DIR}/peaks/treatment_rep2_peaks.narrowPeak \

   --output-file ${FULL_OUTPUT_DIR}/idr/treatment_idr_output.txt \

   --log-output-file ${FULL_OUTPUT_DIR}/idr/treatment_idr.log

awk 'BEGIN {OFS="\t"} {$12 = -log($12)/log(10); print}' ${FULL_OUTPUT_DIR}/idr/treatment_idr_output.txt >
${FULL_OUTPUT_DIR}/idr/treatment_idr_output_log10p>
```

# Step 6: Filtering Peaks Based on IDR Output

```
mkdir -p ${FULL_OUTPUT_DIR}/final_peaks

awk '$12 < 0.05' ${FULL_OUTPUT_DIR}/idr/treatment_idr_output.txt \

   > ${FULL_OUTPUT_DIR}/final_peaks/significant_reproducible_peaks.bed
```

SUBSET

```bash
#!/bin/bash

# Step 0: Create subsets of the original fastq files for testing purposes to speed up the process.

echo "Subsetting data for testing..."

mkdir -p sub_data # Create a directory for the subsetted data.

# Use seqtk to sample 1000 reads from the original fastq.gz files, then gzip the output for FastQC compatibility.

seqtk sample -s100 data/control.fastq.gz 10000 | gzip > sub_data/sub_control.fastq.gz

seqtk sample -s100 data/treatment_rep1.fastq.gz 10000 | gzip > sub_data/sub_treatment_rep1.fastq.gz

seqtk sample -s100 data/treatment_rep2.fastq.gz 10000 | gzip > sub_data/sub_treatment_rep2.fastq.gz


# Step 1: Run FastQC for quality control on the subsetted fastq.gz files.

echo "Running FastQC for quality control..."

mkdir -p output/metrics # Create a directory for FastQC reports.

fastqc sub_data/*.fastq.gz -o output/metrics/ # Run FastQC on all subsetted fastq.gz files and output the reports to the metrics directory.


# Step 2: Align reads to the reference genome with Bowtie2.

mkdir -p output/alignment # Create a directory for the alignment output.

# Align each subsetted fastq.gz file to the reference genome using Bowtie2, outputting SAM files.

bowtie2 -x reference/bowtie2_index/genome -U sub_data/sub_control.fastq.gz > output/alignment/sub_control.sam

bowtie2 -x reference/bowtie2_index/genome -U sub_data/sub_treatment_rep1.fastq.gz > output/alignment/sub_treatment_rep1.sam

bowtie2 -x reference/bowtie2_index/genome -U sub_data/sub_treatment_rep2.fastq.gz > output/alignment/sub_treatment_rep2.sam


# Step 3: Convert SAM files to BAM format and check mapping quality.

echo "Converting SAM to BAM and checking mapping quality..."

# Convert SAM to BAM using samtools for each alignment file.

samtools view -bS output/alignment/sub_control.sam > output/alignment/sub_control.bam

samtools view -bS output/alignment/sub_treatment_rep1.sam > output/alignment/sub_treatment_rep1.bam

samtools view -bS output/alignment/sub_treatment_rep2.sam > output/alignment/sub_treatment_rep2.bam

# Generate mapping quality reports for each BAM file using samtools flagstat.

samtools flagstat output/alignment/sub_control.bam > output/metrics/sub_control_mapping_quality.txt

samtools flagstat output/alignment/sub_treatment_rep1.bam > output/metrics/sub_treatment_rep1_mapping_quality.txt

samtools flagstat output/alignment/sub_treatment_rep2.bam > output/metrics/sub_treatment_rep2_mapping_quality.txt
```

# Step 4: Call peaks using MACS2 with a fixed extension size.

mkdir -p output/peaks # Create a directory for peak calling output.

# Use MACS2 to call peaks on each treatment BAM file against the control, specifying parameters like genome size, q-value cutoff, etc.

macs2 callpeak -t output/alignment/sub_treatment_rep1.bam -c output/alignment/sub_control.bam -f BAM -g dm -n sub_treatment_rep1 --outdir output/peaks -q 0.05 --call-summits --nomodel --extsize 147

macs2 callpeak -t output/alignment/sub_treatment_rep2.bam -c output/alignment/sub_control.bam -f BAM -g dm -n sub_treatment_rep2 --outdir output/peaks -q 0.05 --call-summits --nomodel --extsize 147


# Step 5: Use IDR (Irreproducible Discovery Rate) to assess the reproducibility between biological replicates' peak sets.

echo "Assessing reproducibility with IDR..."

idr --samples output/peaks/sub_treatment_rep1_peaks.narrowPeak \

   output/peaks/sub_treatment_rep2_peaks.narrowPeak \

   --output-file output/idr/sub_treatment_idr_output.txt \

   --plot \

   --log-output-file output/idr/sub_treatment_idr.log


# Step 6: Filter for significant, reproducible peaks based on IDR output.

echo "Filtering significant, reproducible peaks based on IDR output..."

# Use awk to filter the IDR output for entries with a column 12 value (adjusted p-value) less than 0.05, indicating significant reproducibility.

awk '$12 < 0.05' output/idr/sub_treatment_idr_output.txt > output/final_peaks/significant_reproducible_peaks.bed


echo "Pipeline execution for subsetted data completed."