# Correlation of RNA-seq and ChIP-seq data

Stéphanie Le Gras
(slegras@igbmc.fr)

# Exercise 1

We want to know how many up-regulated genes contain a peak for MITF. Compare **Gene names** of the chIPseq data (annotation step) and the RNAseq data (up-regulated genes ).

- Tool: use Venny (https://bioinfogp.cnb.csic.es/tools/venny/)
- Datasets:
  - Use the file siMitfvssiLuc.up.annot.txt (upregulated genes detected by SARtools annotated with BioMart).
    - Download it from your history « RNA-seq data analysis ».
    - Open it with Excel and copy paste the content of column « Gene name » in Venny. You can name this part of the diagramme « RNA-seq »
  - Use the file mitf_peaks.annot.tsv (all chIPseq peaks detected in the second run of MACS2 and annotated with BEDtools closest)
    - Download it from your history « ChIP-seq data analysis »
    - Open it with Excel and copy paste the content of column « O » (the one with Gene names) in Venny. You can name this part of the diagramme « ChIP-seq »

# Exercise 2

Use seqMINER to visualize at the same time chIP-seq data along with RNA-seq data.

Tool: seqMINER

Datasets:

- Reference coordinated: MITF_peak_summits.bed (from your history "ChIP-seq data analysis")
- Density
- seqMINER accept a tab separated file formatted like:
  - Gene ID <tab> Expression Values
  - Aligned reads (available in chipseq/mapping):
    - Mitf.sort.bam
    - H3K4me3.sort.bam
    - polII.sort.bam
  - RNA-seq data:
    - seqMINER accept a tab separated file formatted like: Gene ID <tab> Expression Values
      - Expression values used in our example are normalized read counts divided by gene length in Kb.
    - **Let's generate this file!**

# Prepare exercise 2

In seqMINER, we are going to visualize and compare expression values of different genes within the same sample. Read count per gene can not be directly compared together as they are correlated to the gene length. In the file siMitfvssiLuc.complete.txt, you can find normalized values. We are going to scale them as if genes were all of the same size.

Steps:

1. Extract transcript lengths [ensembl/BioMart]
2. As we are working on read counts per gene and not per transcript, compute a median of transcript lengths per gene [Galaxy]
3. Add the median of transcript lengths to the dataset siMitfvssiLuc.complete.txt [Galaxy]
4. Compute normalized and divided by median of transcripts length in kb values per gene on normalized data for siLuc (rounded mean of normalized counts) [Excel]

# Prepare exercise 2

1. In Ensembl/BioMart
   1. For all genes, extract the following information:
      - Gene stable IDs
      - Transcript stable IDs
      - Transcript length (including UTRs and CDS)
   2. Rename the file: hg38_ens105_transcriptLength.txt.gz (compressed .gz)

# Prepare exercise 2

3. Create a new history « Prepare RNA-seq data for seqMINER ». Import the file hg38_ens105_transcriptLength.txt.gz to Galaxy (type: tabular (3), Genome: hg38 (4))

# Prepare exercise 2

2.1. Use the tool « **Datamash** (operations on tabular data) » to group gene by Ensembl Gene Ids and compute the median on transcript length:
- **Input tabular dataset:** hg38_ens105_transcriptLength.txt.gz
- **Group by fields:** [column with Gene stable ID] *(1)*
- **Input file has a header line:** Yes
- **Print header line:** Yes
- **Operation to perform on each group:**
  - **Type:** Median
  - **On column:** [column with Transcript length (including UTRs and CDS)] *(3)*

2.2. Use the tool « **Compute** » to round median values:
- **Add expression:** round(c2) *(change accordingly if needed)*
- **Input has a header line with column names?** Yes
  - **The new column name:** rounded median(Transcript length (including UTRs and CDS))

2.3. Use the tool « **Advanced Cut** » to extract only column of interest:
- **File to cut:** [Compute on data *] *(result of step 2.2)*
- **Cut by:** fields
  - **List of Fields:** Column: 1 Column: 3 *(change accordingly if needed)*

Rename resulting dataset: Median_of_transcript_length.tsv

# Prepare exercise 2

3.

1. Import the file siMitfvssiLuc.complete.txt (history "RNA-seq data analysis") to Galaxy (type: tabular, Genome: hg38).

2. Use the tool "**Join two Datasets**" to join the two datasets siMitfvssiLuc.complete.txt  and Median_of_transcript_length.tsv

   - **Join:** siMitfvssiLuc.complete.txt
   - **Using column:** 1
   - **With:** Median_of_transcript_length.tsv
   - **and column**: 1
   - **Keep the header lines**: Yes

3. Rename the file siMitfvssiLuc.complete.wTranscriptLength.tsv

4. Download siMitfvssiLuc.complete.wTranscriptLength.tsv and open it in excel.

# Prepare exercise 2

## 4. In Excel (or an equivalent), add 1 column:

1. siLuc (normalized and divided by median of transcripts length in kb) filled with the following formula (French ; English)

   =ARRONDI(K2/Y2*1000/50;0) (K is the column : siLuc – Y is the column with round median )

   =ROUND(K2/Y2*1000/50;0)

   Hint: here we divide values by 50 to get them in the range of the chIP-seq data that we are going to visualize along with them.

# Prepare exercise 2

In Excel (or an equivalent), create a new file with the following columns **without headers**:

- Id (contains ENSEMBL IDs)
- siLuc (normalized and divided by median of transcripts length in kb)

    Hint: copy and paste normalized data with a special paste - **by value** – so that it doesn't copy the formula
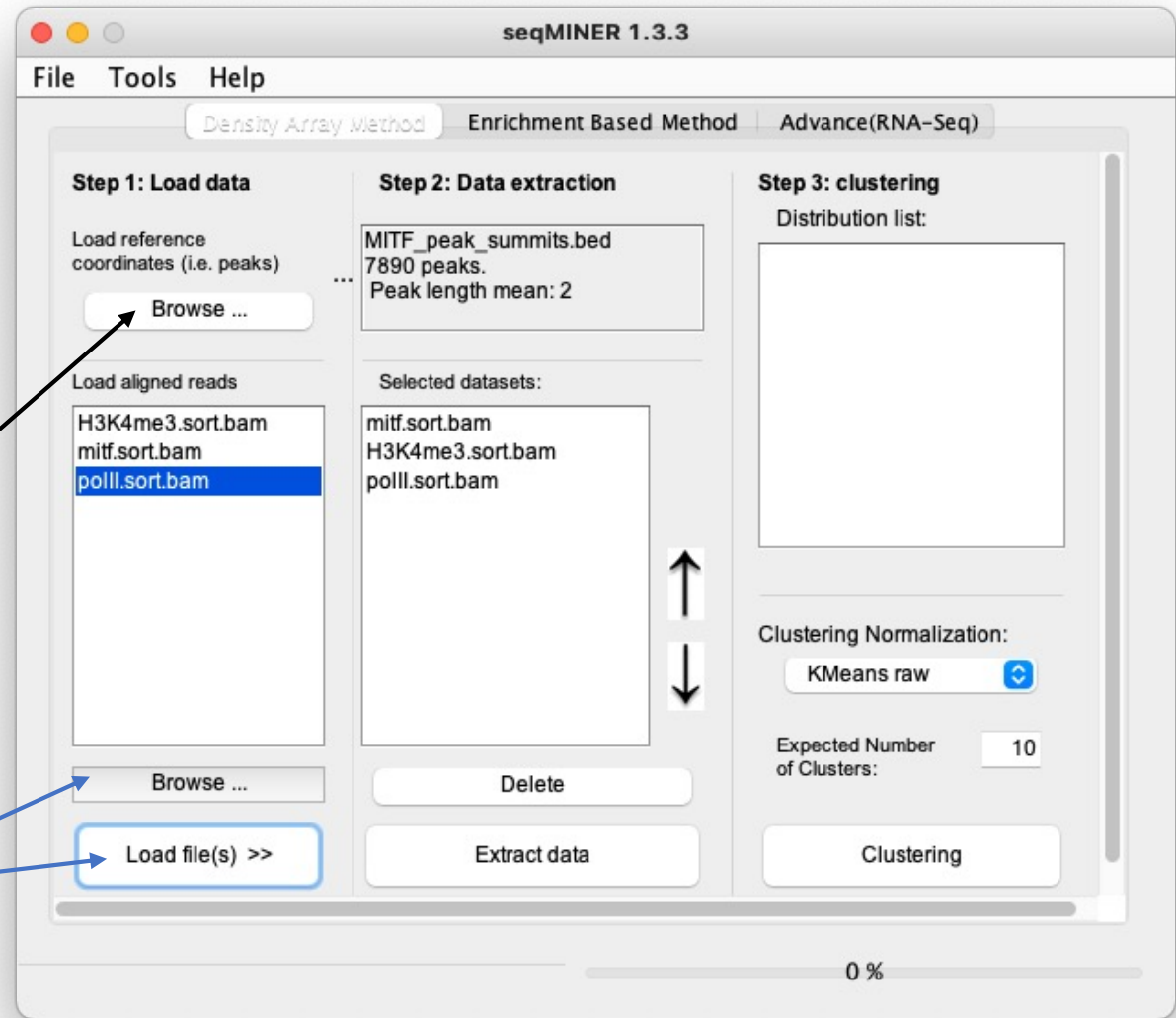
   Save the file as a Text (separator: tabs) (.txt) file named RNAseq_data_ready_for_seqMINER.txt

| A | B | C |
|---|---|---|
| ENSG0000000 | 26 | |
| ENSG0000000 | 0 | |
| ENSG0000000 | 62 | |
| ENSG0000000 | 5 | |
| ENSG0000000 | 20 | |
| ENSG0000000 | 0 | |
| ENSG0000000 | 1 | |

# Exercise 2

Use seqMINER to visualize at the same time chIP-seq data along with RNA-seq data

- Load MITF_peak_summits.bed (from your history "ChIP-seq data analysis") as reference coordinates.
- Load the 3 bam files (directory chipseq/mapping):
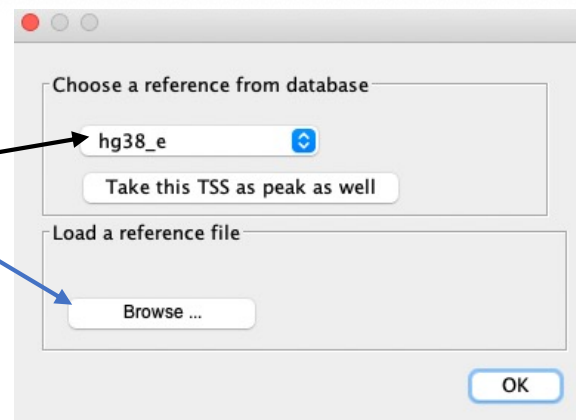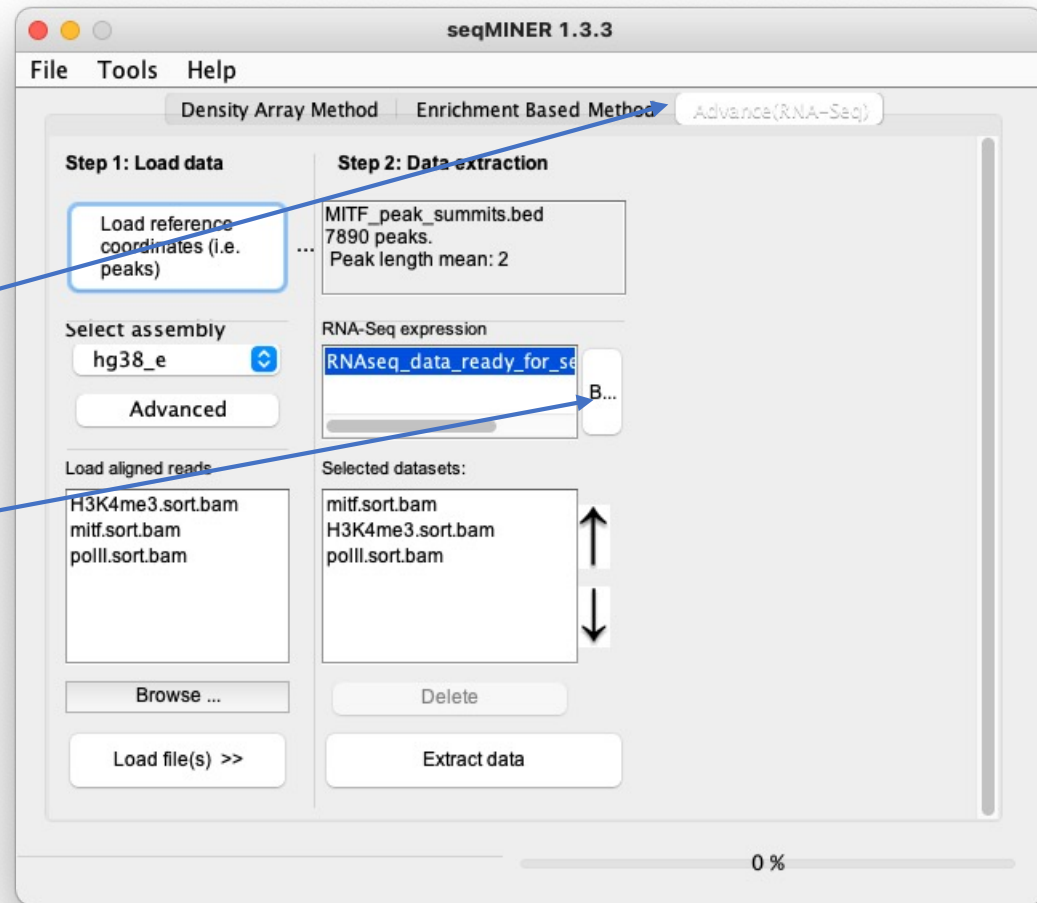  - mitf.sort.bam
  - H3K4me3.sort.bam
  - polII.sort.bam

# Exercise 2

In the Advance (RNAseq) tab:

1. Upload the file RNAseq_data_ready_for_seqMINER.txt:
   - 1st column contains Ensembl Gene IDs
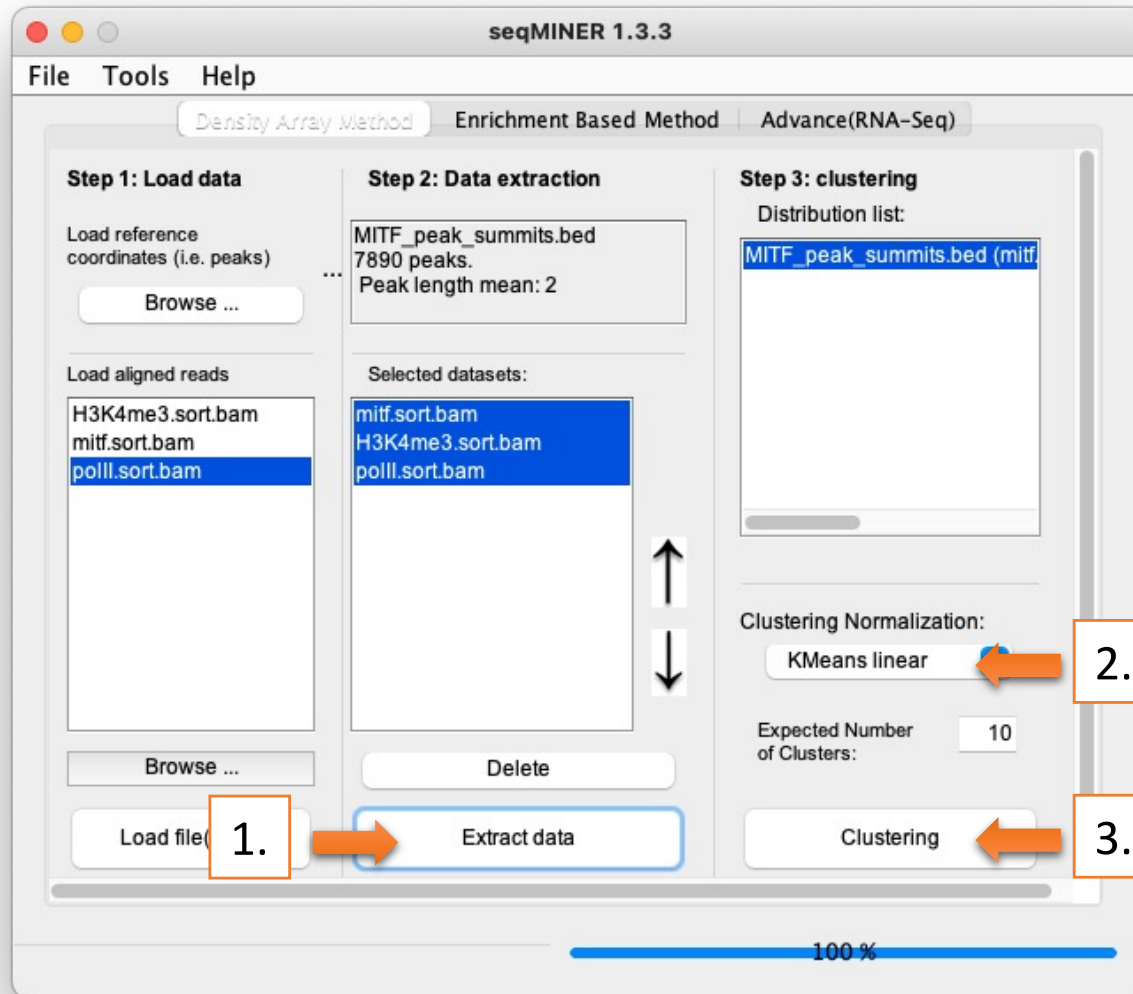   - 2nd column contains normalized read counts of siLuc divided by gene length in Kb.

2. In advanced:
   - Load a reference file: click on Browse and select the file hg38_ens105.bed.
   - Choose a reference from database: hg38_ens105.bed

# Exercise 2

Be careful, make sure that in Options > Gene profile, Gene profile analysis is not selected before clicking on « **Extract data** ».

# Exercise 2