

UNIVERSITÁ DEGLI STUDI DI SALERNO  
DOTTORATO IN MANAGEMENT & INFORMATION TECHNOLOGY



CURRICULUM: INFORMATION SECURITY & INNOVATION SYSTEMS

COORDINATORE: Ch.mo. Prof. Antonelli Valerio

Ciclo XVII N.S.

Novel tools for reproducible  
Next Generation Sequencing data analysis and integration

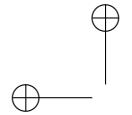
**Relatori**

Ch.mo. Prof. Tagliaferri Roberto  
Ch.mo. Prof. Angelini Claudia

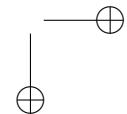
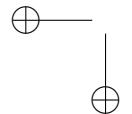
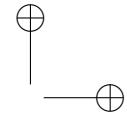
**Candidato**

Righelli Dario  
Matr. 8800800010

ANNO ACCADEMICO 2017/2018



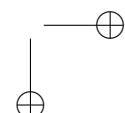
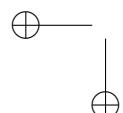
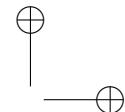
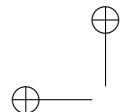
“Template” — 2018/9/10 — 23:00 — page 2 — #2

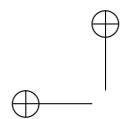


*How to reach a goal?*

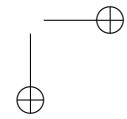
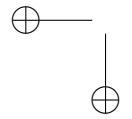
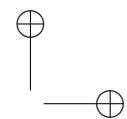
*Without haste but without rest*

*Goethe*





“Template” — 2018/9/10 — 23:00 — page 4 — #4



Add acknowledgements here

⊕

⊕

“Template” — 2018/9/10 — 23:00 — page 6 — #6

⊕

⊕

⊕

⊕

⊕

⊕

Write your abstract here



# Contents

<b>Acknowledgements</b>	<b>5</b>
<b>Abstract</b>	<b>7</b>
<b>1 Introduction</b>	<b>11</b>
1.1 Biological Background . . . . .	11
1.2 Sequencing Techniques . . . . .	11
1.2.1 RNA-Seq . . . . .	11
1.2.2 Atac-Seq . . . . .	11
1.3 Computational Aspects . . . . .	11
<b>2 TiCoRSe - Time Course RNA-Seq data analysis</b>	<b>13</b>
2.1 Introduction . . . . .	13
2.1.1 Time Course RNA-Seq . . . . .	13
2.2 Methods . . . . .	13
2.2.1 General Approach . . . . .	13
2.2.2 Time Course Methods . . . . .	13
2.2.3 Other Methods . . . . .	13
2.2.4 Additional Features . . . . .	13
2.3 Results . . . . .	13

<b>3 DEScan2 - Differential Enriched Scan 2</b>	<b>15</b>
3.1 Introduction . . . . .	15
3.2 Methods . . . . .	17
3.2.1 Peak Caller . . . . .	18
3.2.2 Peak Filtering and Alignment . . . . .	18
3.2.3 Peak Counts . . . . .	19
3.2.4 Additional Features . . . . .	20
3.3 Case Study . . . . .	20
<b>4 IntegrHO - Integration of High-Throughput Omics data</b>	<b>29</b>
4.1 Introduction . . . . .	30
4.2 Methods . . . . .	30
4.2.1 Single Omic Approach . . . . .	30
4.2.2 Multi Omic Approach . . . . .	30
4.3 Implementation Aspects . . . . .	30
4.4 Reproducible Computational Research . . . . .	30
4.5 Results . . . . .	30
<b>5 Conclusions &amp; Future Works</b>	<b>31</b>
<b>Appendices</b>	<b>33</b>
.1 R Language . . . . .	35
.2 R Markdown Language . . . . .	35
<b>6 Bibliography</b>	<b>37</b>
<b>List of Figures</b>	<b>41</b>
<b>List of Tables</b>	<b>43</b>

Chapter **1**

## Introduction

### 1.1 Biological Background

### 1.2 Sequencing Techniques

#### 1.2.1 RNA-Seq

#### 1.2.2 Atac-Seq

### 1.3 Computational Aspects



# Chapter 2

## TiCoRSe - Time Course RNA-Seq data analysis

### 2.1 Introduction

#### 2.1.1 Time Course RNA-Seq

### 2.2 Methods

#### 2.2.1 General Approach

#### 2.2.2 Time Course Methods

#### 2.2.3 Other Methods

#### 2.2.4 Additional Features

### 2.3 Results



# Chapter 3

## DEScan2 - Differential Enriched Scan 2

### ***few words on integration of epigenomic with transcriptomic***

To investigate and answer epigenetic biological questions we decided to create a useful instrument for analysing epigenomic data (such as *ChIP-Seq*, *Atac-Seq*, *Sono-Seq*). Very often the biological questions to be answered, as for the RNA-Seq, need the comparison of two or more different biological conditions. Starting from a set of already published [1] scripts, we designed *Differential Enriched Scan 2 (DEScan2)*, a software for helping the analysis of epigenomic data.

### 3.1 Introduction

The *Differential Enriched Scan 2* is an R [2] tool developed for detecting epigenomic signal in order to facilitate the Differential Enrichment of genomic regions (DERs) between two or more biological conditions.

The package has been implemented using Bioconductor [3] data structures and methods, and it is available on Bioconductor since version 3.7.

The tool is organized in three main steps. A peak caller, which is a standard moving scan window that compares the counts within a sliding window, to the counts in a larger region outside the window. It uses a Maximum Likelihood Estimator on a Poisson Distribution, providing a final score for each detected peak.

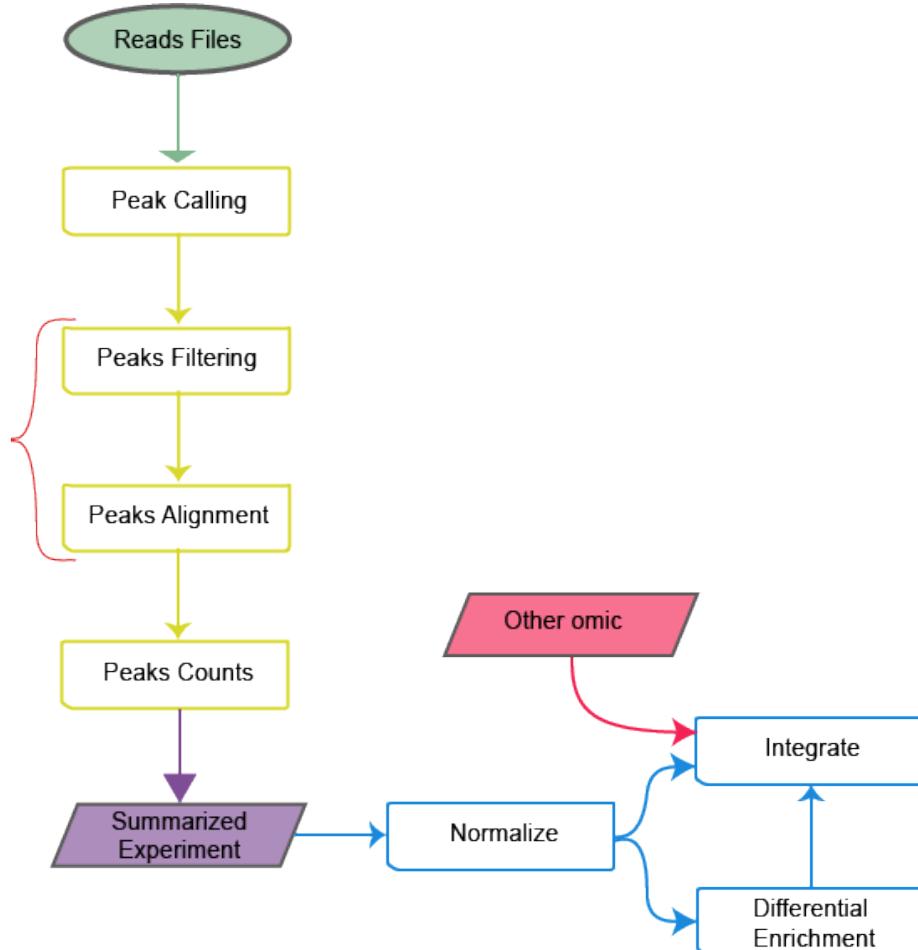
The filtering step is aimed to determine if a peak is a "true peak" on the basis of its replicability in other samples. This step is based on a double user-defined threshold, one on the peak's scores and one on the number of samples.

Finally, the third step produces a counts matrix where each column represents a sample and each row a peak. The value of each cell is the number of reads for the peak in the sample.

The so produced counts matrix, as illustrated in the figure 3.1.1, is useful both for doing differential enrichment between multiple conditions and for integrating the epigenomic data with other -omic data types.

### 3.2. METHODS

17



**Figure 3.1.1:** A differential enrichment flow representation. DEScan2 steps are highlighted in yellow.

## 3.2 Methods

The package is organized in three main steps, the peak caller in section 3.2.1, the filtering and alignment of the peaks in section 3.2.2 and the peak counting

described in section 3.2.3.

It offers some additional features described in 3.2.4.

### **3.2.1 Peak Caller**

The Peak Caller (`findPeaks` method) takes as input a set of alignment files (BAM [4] or BED format) with the code of the reference genome (i.e. *mm10* for Mus Musculus version 10) and several additional parameters, useful for the peak detection setup.

The alignment data are stored as GenomicRangesList [5], where each element represents a file. In order to facilitate the parallelization of the computations over the chromosomes, the list is re-arranged as a list of chromosomes of GenomicRangesList, where each element represents a file containing just the GenomicRanges of the specific chromosome (see section 3.2.4).

On this data structure the algorithm firstly divides each chromosome as bins of `binSize` parameter length (default is 50bp) and then computes the reads coverage on the bins with moving scan windows, spanning from `minWin` to `maxWin` parameters of `binSize` interval.

In order to be able to catch small and spread peaks the algorithm computes the coverage also using windows of two different lengths, that can be defined with `minCompWinWidth` and `maxCompWinWidth` (defaults are 5000 and 10000) parameters.

The so produced coverages are useful to compute a score for each detected region, applying a Maximum Likelihood Estimation (MLE) on the coverages between the sliding windows, assuming a Poisson distribution.

[PUT THE POISSON DISTRIBUTION AND THE LIKELIHOOD]

### **3.2.2 Peak Filtering and Alignment**

In order to filter out false positives peaks, we designed a filtering method (`finalRegions`) based on two different thresholds. A first threshold on the peaks score and a second threshold on the number of samples.

The filtering step is designed to take as input a list of peaks as GenomicRangesList, where each element represents a chromosome. This is the data

### 3.2. METHODS

19

structure produced by the peak caller, but, we also developed a method to load peaks produced also by other software like MACS [6], as described in section 3.2.4.

Firstly, using the threshold on the peak’s score (`zThreshold` parameter), the method filters out the peaks with a score lower than the user-defined threshold value.

Then it extends a 200bp window in both directions of the detected region, computing the overlaps between the samples using the `findOverlapsOfPeaks` method (with `connectedPeaks` parameter set as `merge`), defined by the ChIP-peakAnno [7] Bioconductor package.

Based on this idea, the filtering step is developed to filter out those peaks not present in at least a user-defined (`minCarriers` parameter) number of samples. In the light of this, the user can decide the minimum number of samples where each peak has to be detected. We suggest to set the threshold as a multiple of the number of replicates of the conditions.

#### 3.2.3 Peak Counts

The counting step (`countFinalRegions` method) is designed to take a *GenomicRanges* data structure as input, where for each peak additional features, as the score and the number of samples, are saved. Moreover, it requires also the path of the BAM/BED files where the reads are stored, in order to quantify the peaks given as input.

For each region the method counts the number of reads present in each sample. In so doing, it produces as result a matrix of the counts, where the rows and the columns represent, respectively, the regions and the samples.

In order to keep trace of all information associated to the regions, it produces a *SummarizedExperiment* [8] data structure, giving the possibility to retrieve the *GenomicRanges* peaks associated data structure and the count matrix, respectively, with `rowRanges` and `assays` method.

The choice to produce a count matrix is guided by the versatility of this data structure, useful not only for the differential enrichment of the regions between multiple conditions, but also for integrating the epigenomic data with other

-omics.

### 3.2.4 Additional Features

However, the package can work with any external peak caller returning results in terms of bed files, indeed the package provides additional functionalities to load BED files of peaks and handle them as GenomicRanges [5] structures.

Furthermore, our package provides several functionalities for GenomicRanges data structure handling. One over the others gives the possibility to split a GenomicRanges over the chromosomes to speed-up the computations parallelizing them over the chromosomes.

## 3.3 Case Study

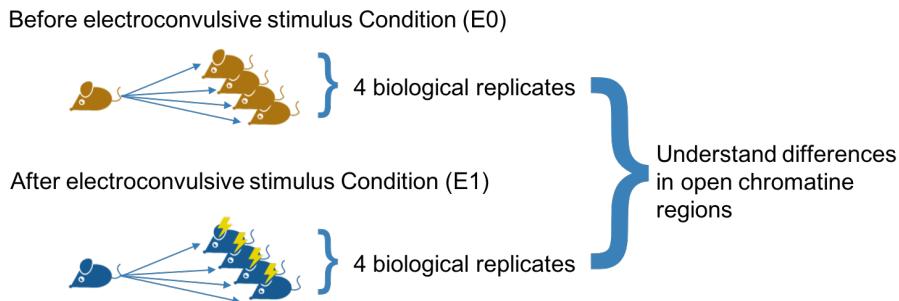
### A few words on epigenomic data

We illustrate the performances of DEScan2 using a dataset [9] that describes *in vivo* adult mouse dentate granule neurons before and after synchronous neuronal activation using Atac-Seq and RNA-Seq technologies (see sections 1.2.2 and 1.2.1 for a description of these sequencing techniques).

This dataset is organized in 62 samples of Atac-Seq and RNA-Seq, extracted at different time points, with four replicates at each time point. We chose to compare the differences at the first two stages, time 0 (E0) and 1 hour after neuronal induction (E1), in order to show a possible Atac-Seq workflow for Differential Enrichment, and how to integrate this data type with RNA-Seq. A general illustration of our dataset is represented in 3.3.1.

### 3.3. CASE STUDY

21



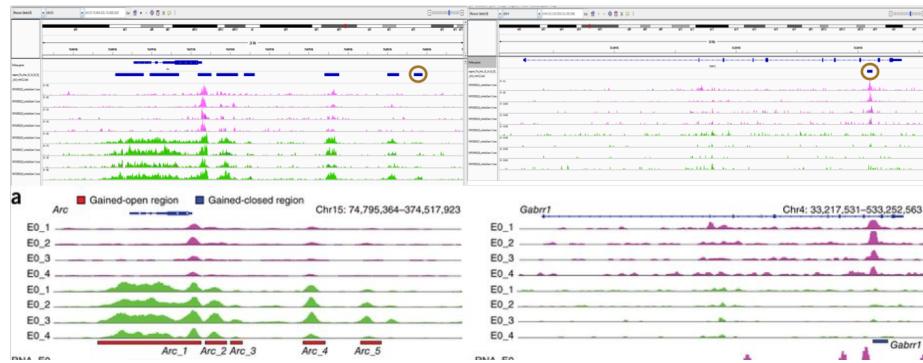
**Figure 3.3.1:** An illustration of our extraction of the [9] dataset.

We downloaded the data from *GEO* database [10, 11] with accession number GSE82015<sup>1</sup> and mapped raw data using *STAR* [12] with default parameter on *Mus Musculus* Genome ver.10 (mm10).

In order to detect the open chromatin regions we run our peak caller, cutting the genome in bins of 50bp and using running windows of minimum 50bp and maximum 1000bp. In such a way we are able to detect not just broad peak, but also smaller peaks.

To be confident with our results we compared the DEScan2 detected peaks with the same validated regions (Arc and Gabrr1) in the original work [9]. The lower part of figure 3.3.2 shows the detected and validated regions (in blue and red) resulting differentially enriched between the E0 (in pink) and E1 (in green) conditions, while the upper part shows DEScan2 peaks (in blue), highlighting a capability to catch not only the same regions of the published ones, but also (gold circles) to be more careful in the smaller peaks detection.

<sup>1</sup><https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE82015>

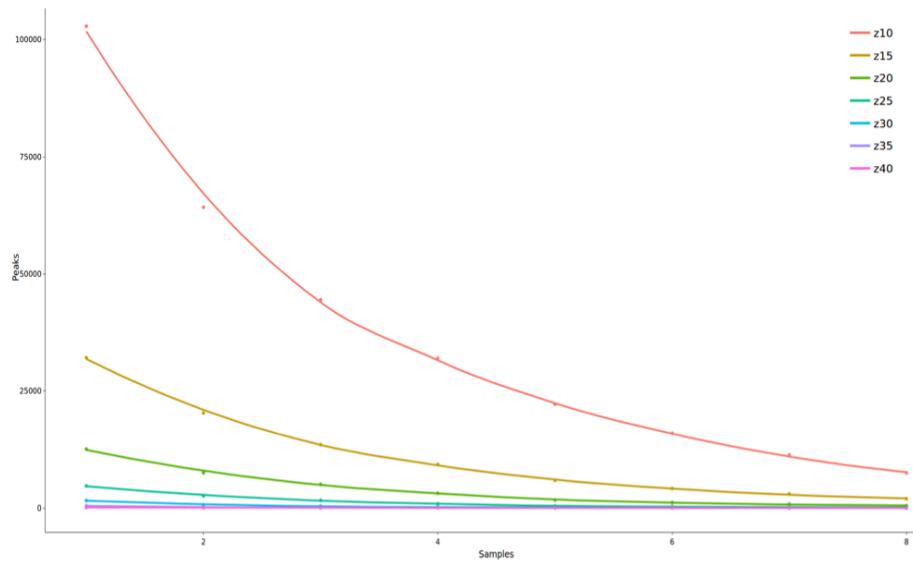


**Figure 3.3.2:** A comparison of DEScan2 detected peaks with validated peaks in article [9].

While it is very important to detect good peaks with a peak caller, it seems to be more relevant to detect reliable regions. Indeed, during the filtering step, the number of peaks depends not only by the peak score, but also by the number of replicates designed in the experiment. The figure 3.3.3 puts in relation these two relevant information. On the x-axis is represented the number of replicates, while on the y-axis is traced the number of peaks, and each curve represents a different threshold on the peaks score, showing that higher are the thresholds on the scores and the number of replicates, lower is the number of the detected peaks. Highlighting a proportional inversion between the number of the peaks and the combination of the number of samples and the detected regions score.

### 3.3. CASE STUDY

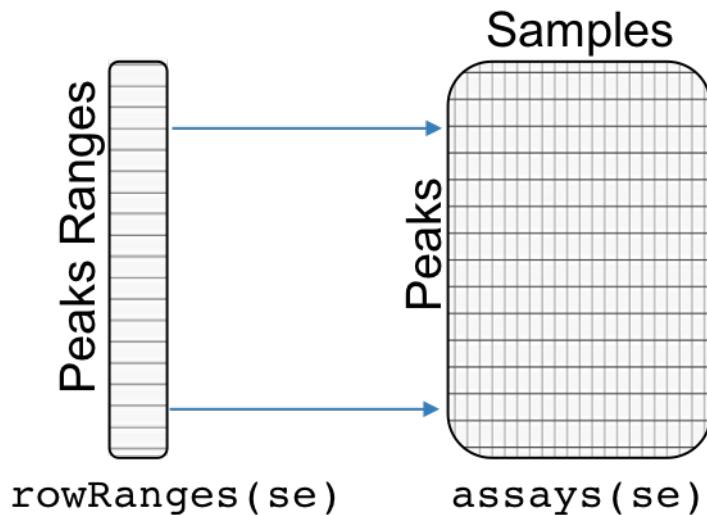
23



**Figure 3.3.3:** Filtering the detected regions with different thresholds on peak scores.

The filtered-in regions can be processed by DEScan2 in order to obtain a count matrix with samples on the columns and peaks on the rows. This type of data structure is very versatile, because it enables to perform several operations, like the differential enrichment of regions (DERs) and, if possible, the integration with other kind of omics, as RNA-Seq.

In order to preserve the information associated to the peaks, DEScan2 produces as output a *SummarizedExperiment* (figure 3.3.4) data structure, which enables to retrieve the count matrix with `assays` method, and to access the peaks information in *GenomicRanges* format with the `rowRanges` method.



**Figure 3.3.4:** An illustration of the `SummarizedExperiment` data structure produced by DEScan2.

Before to proceed to detect DERs, it is a good standard to normalize the data, also because without any kind of normalization we are not able to detect any DERs. The nature of the data, in count format, makes it possible to apply several well known RNA-Seq normalizations techniques, as *TMM*, *upper-quartile*, *full-quartile*, *RUV-Seq*, etc [13–15].

While the *TMM* and *upper-quartile* normalizations modify the data in a way that makes it impossible to detect DERs, other kind of normalizations and combinantions of them give good results.

The figure 3.3.5 sintetizes this concept very well, highlighting a relation between the number of DERs and the minumum number of samples used for filtering the data during the DEScan2 filtering step.

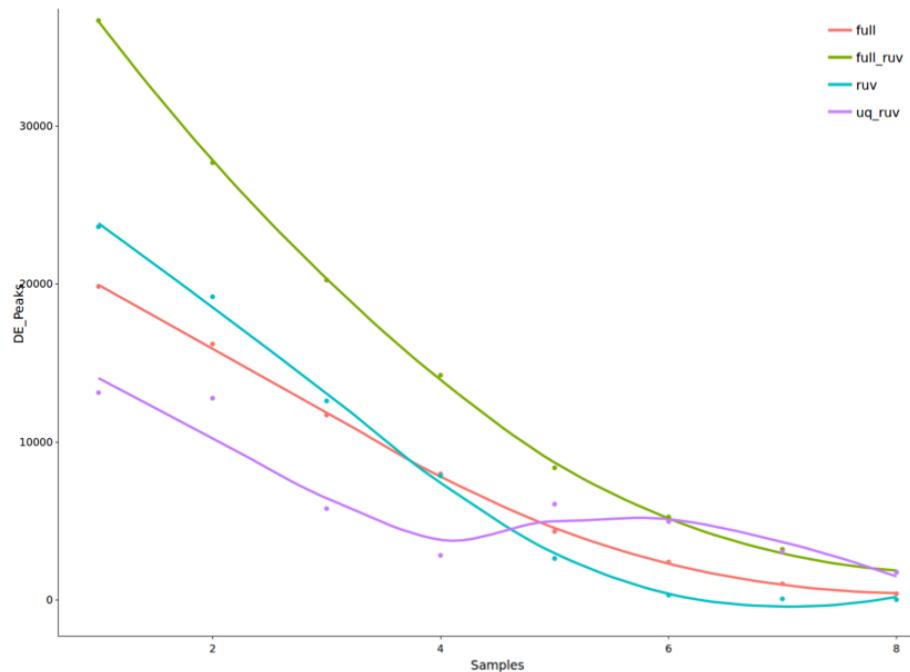
The plot shows that *upper-quartile*, even if combined with *RUV-Seq* normalization, is not able to linearly detect a good amount of DERs, while *full-quartile*, when combined with *RUV-Seq* seems to affect the data in a way that overde-

### 3.3. CASE STUDY

25

tect the number of DERs. When looking at the *full-quantile* and *RUV-Seq* by themself seem to perform better than the other normalizations. The first one has a downhill almost linear, while the second one has a very fast downhill with a regrowth when the number of samples is higher.

Even if these normalization methods show good performances with this type of epigenomic data, our investigations suggest that for sure more testing is required, but maybe an ad-hoc normalization method for these data has to be developed.



**Figure 3.3.5:** The figure shows the effects of different normalizations on the epigenomic differentially enriched regions.

To estimate the DERs any of the RNA-Seq methods can be applied, such as *DESeq2*, *edgeR*, *NOISEq*, etc [16–18].

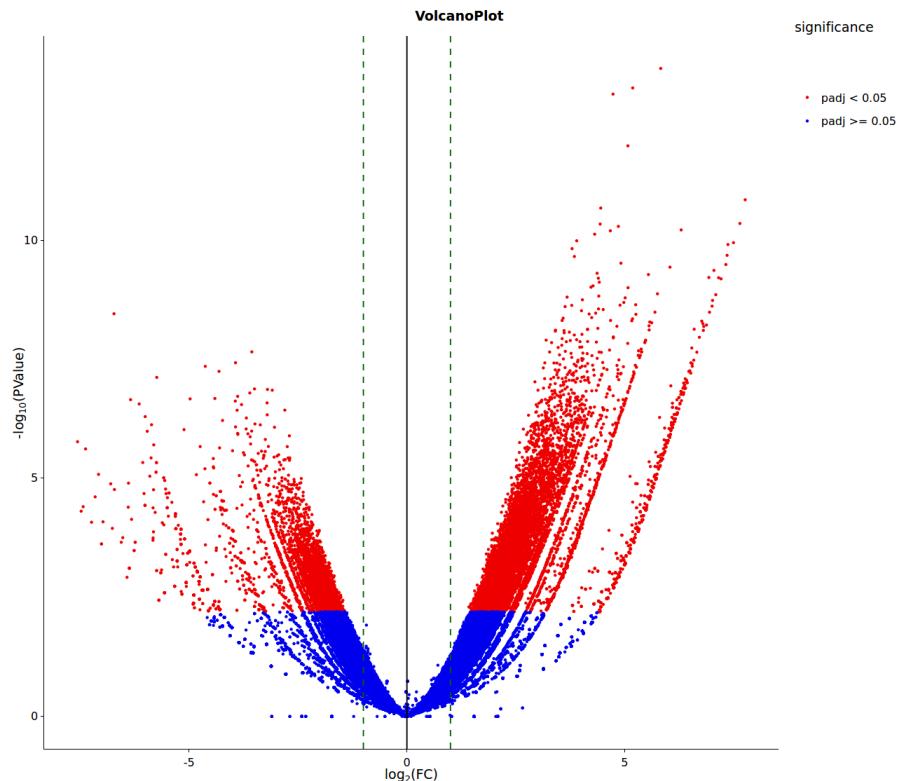
In this case, decided to use *edgeR* package, because of its wide range of avail-

able statistical methods and the possibility to better tune the design of the experiment. Indeed, in our case, we used the RUV-Seq normalized counts with `k` parameter set to 4. We modeled the design with the `model.matrix` function, adding to our model not only the experimental conditions, but also the RUV-Seq estimated weights. Then we used the resulting design to estimate the dispersion and to fit a Quasi-Likelihood test, as defined in edgeR.

The figure 3.3.6 shows a volcano plot of DERs between E0 and E1 conditions. Red dots highlights the regions with a False Discovery Rate (FDR) lower than 0.05, while blue dots highlight not significant regions.

### 3.3. CASE STUDY

27

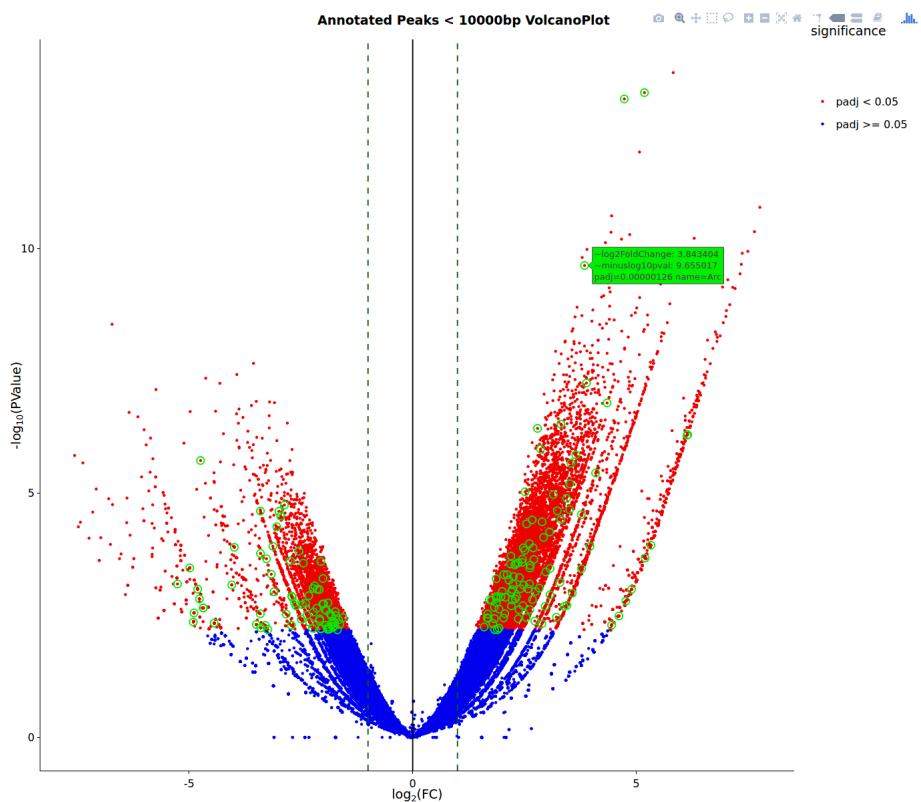


**Figure 3.3.6:** A volcano plot of Differential Enriched Regions. Blue dots represent the not significant DERs, while the red ones represent the significant DERs.

Next task is to integrate the obtained results with other omic data types, as RNA-Seq. Because of the low number of the samples, the easiest way to integrate the data is to annotate the DERs with differentially expressed genes resulting from the analysis of RNA-Seq.

For the Differential expression of the RNA-Seq data we firstly quantified the signal with featureCounts methods available in the Subread Bioconductor package. (cite) Then we filtered lowly expressed genes with the *proportion* as implemented in NOISeq package, and applied the noiseq method. We selected the

significant DE genes with a probability higher than 0.95, and used these genes to annotate the peaks with `annotatePeakInBatch` method of `ChIPpeakAnno`. Figure 3.3.7 illustrates with green circles the peaks with a gene annotated with a distance lower than 10000bp from the TSS of the gene. Realizing the plot with the `plotly` library it’s possible to enhance the names of the genes with a tip window.



**Figure 3.3.7:** A volcano plot of Differential Enriched Regions. Blue dots represent the not significant DERs, while the red ones represent the significant DERs. Green circles highlights the peaks with a Differential Enriched gene annotated.

Chapter **4**

# IntegrHO - Integration of High-Throughput Omics data

## 4.1 Introduction

## 4.2 Methods

### 4.2.1 Single Omic Approach

### 4.2.2 Multi Omic Approach

Low Level Itegration

High Level Itegration

## 4.3 Implementation Aspects

## 4.4 Reproducible Computational Research

## 4.5 Results

# Chapter 5

## Conclusions & Future Works



# Appendices

⊕

“Template” — 2018/9/10 — 23:00 — page 34 — #34

⊕

⊕

⊕

⊕

⊕

⊕

⊕

## **.1 R LANGUAGE**

---

**35**

### **.1 R Language**

### **.2 R Markdown Language**



# Chapter 6

## Bibliography

1. Koberstein, J. N. *et al.* Learning-dependent chromatin remodeling highlights noncoding regulatory regions linked to autism. *Science Signaling*. ISSN: 19379145. doi:10.1126/scisignal.aan6500 (2018).
2. Ihaka, R. & Gentleman, R. R: A Language for Data Analysis and Graphics. *Journal of Computational and Graphical Statistics*. ISSN: 15372715. doi:10.1080/10618600.1996.10474713. arXiv: arXiv:1011.1669v3 (1996).
3. Gentleman, R. *et al.* Bioconductor: open software development for computational biology and bioinformatics. *Genome Biology*. ISSN: 1465-6914. doi:10.1186/gb-2004-5-10-r80 (2004).
4. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics*. ISSN: 13674803. doi:10.1093/bioinformatics/btp352. arXiv: 1006.1266v2 (2009).
5. Lawrence, M. *et al.* Software for Computing and Annotating Genomic Ranges. *PLoS Computational Biology* **9** (ed Prlic, A.) e1003118. ISSN: 1553-7358 (2013).
6. Zhang, Y. *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biology*. ISSN: 14747596. doi:10.1186/gb-2008-9-9-r137 (2008).

7. Zhu, L. J. *et al.* ChIPpeakAnno: A Bioconductor package to annotate ChIP-seq and ChIP-chip data. *BMC Bioinformatics*. ISSN: 14712105. doi:10.1186/1471-2105-11-237 (2010).
8. Morgan M, Obenchain V, Hester J, P. H. SummarizedExperiment: SummarizedExperiment container. doi:<https://doi.org/doi:10.18129/B9.bioc.SummarizedExperiment> (2018).
9. Su, Y. *et al.* Neuronal activity modifies the chromatin accessibility landscape in the adult brain. *Nature Neuroscience*. ISSN: 15461726. doi:10.1038/nn.4494 (2017).
10. Edgar, R. Gene Expression Omnibus: NCBI gene expression and hybridization array data repository. *Nucleic Acids Research*. ISSN: 13624962. doi:10.1093/nar/30.1.207 (2002).
11. Barrett, T. *et al.* NCBI GEO: Archive for functional genomics data sets - Update. *Nucleic Acids Research*. ISSN: 03051048. doi:10.1093/nar/gks1193 (2013).
12. Dobin, A. *et al.* STAR: Ultrafast universal RNA-seq aligner. *Bioinformatics*. ISSN: 13674803. doi:10.1093/bioinformatics/bts635. arXiv: 1201.0052 (2013).
13. Risso, D., Ngai, J., Speed, T. P. & Dudoit, S. Normalization of RNA-seq data using factor analysis of control genes or samples. *Nature Biotechnology*. ISSN: 15461696. doi:10.1038/nbt.2931. arXiv: arXiv:1011.1669v3 (2014).
14. Robinson, M. D. & Oshlack, A. A scaling normalization method for differential expression analysis of RNA-seq data. *Genome Biology*. ISSN: 14747596. doi:10.1186/gb-2010-11-3-r25. arXiv: PMC2864565 (2010).
15. Dillies, M. A. *et al.* A comprehensive evaluation of normalization methods for Illumina high-throughput RNA sequencing data analysis. *Briefings in Bioinformatics*. ISSN: 14675463. doi:10.1093/bib/bbs046 (2013).

**BIBLIOGRAPHY**

**39**

16. Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics (Oxford, England)* **26**, 139–140. ISSN: <null> (2009).
17. McCarthy, D. J., Chen, Y. & Smyth, G. K. Differential expression analysis of multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids Research* **40**, 4288–4297. ISSN: 03051048 (2012).
18. Tarazona, S., García, F., Ferrer, A., Dopazo, J. & Conesa, A. NOIseq: a RNA-seq differential expression method robust for sequencing depth biases. *EMBnet.journal* **17**, 18. ISSN: 2226-6089 (2012).



## List of Figures

3.1.1 DEScan2 workflow . . . . .	17
3.3.1 DEScan2 dataset illustration . . . . .	21
3.3.2 DEScan2 peaks detection . . . . .	22
3.3.3 DEScan2 filtering step . . . . .	23
3.3.4 DEScan2 counts illustration . . . . .	24
3.3.5 Normalizations applied to detected regions . . . . .	25
3.3.6 Differential Enrichment Regions Volcano . . . . .	27
3.3.7 Annotated Differential Enrichment Regions Volcano . . . . .	28



## List of Tables