

# MASTER'S THESIS IN BIOINFORMATICS

## T-CELL CLONOTYPING FROM BULK WHOLE TRANSCRIPTOME RNA-SEQ DATA

Sergio Fernández Adán and Sergio Galera Alquegui  
Bioinformatics Research Centre (BiRC),  
Aarhus University



Master's thesis' supervisors:  
Christian Storm Pedersen, BiRC - Aarhus University  
Paula Tataru Bjerg, QIAGEN - Aarhus



# TABLE OF CONTENTS

<b>0. Abstract .....</b>	<b>5</b>
<b>1. Introduction.....</b>	<b>7</b>
1.1. Immune system intro + T-Cell Receptor background .....	7
1.2. Sequencing .....	9
1.3. Software for benchmarking .....	12
1.3.1. Immune repertoires extraction .....	12
1.3.2. Immune repertoires comparison .....	15
1.4. Optimising results .....	16
<b>2. Material/Datasets and methods.....</b>	<b>17</b>
2.1. Dataset selection .....	17
2.2. Data pre-processing .....	18
2.2.1. Software and tools .....	19
2.2.2. Reads structure discovery .....	24
2.2.3. Reads pre-processing for optimal yield .....	28
2.3. Immune Repertoire Analysis .....	31
2.4. Benchmarking.....	33
2.5. Optimization and Filtering.....	36
2.5.1. De Novo assembly .....	36
2.5.2. Map reads to reference/RNA seq analysis .....	38
2.5.3. Our implementation .....	40
<b>3. Results .....</b>	<b>45</b>
3.1 Sensitivity.....	45
3.2 CDR3 amino acid sequence.....	48
3.3 Clonotypes abundance.....	49
3.3.1 Clone frequency correlation: CLC Workbench vs MiXCR.	51
3.3.2 Quantitative comparison of shared clonotypes.....	53
3.4 Overlap between repertoires.....	54
<b>4. Conclusions .....</b>	<b>59</b>
<b>5. References.....</b>	<b>61</b>
<b>6. Appendices.....</b>	<b>62</b>



## 0. *Abstract*

Our immune system is a compound of different types of cells, tissues, and organs that help us to get protection against harmful infections and diseases. T-cells are one of the cells that take a role in this system, and it is defined as a cellular receptor associated with an intrinsic enzymatic response that acts in diverse intracellular pathways. Usually, we require targeted sequencing data to extract those sequences that have been classified as TCR data to start analysing the immune system of a sample. In this report, we present the *CLC Genomics Workbench* as a tool for immune repertoire extraction out of bulk RNA-seq data. We benchmarked our results against *MiXCR*, open-source software that recalls the correct results with high accuracy and precision and that has been widely used by scientists to extract immune data. We have obtained plenty of results from filtering data by applying different methods after correctly pre-processing all the data. In the process of creating a set of results, we have also developed an algorithm that will be explained in further detail in this report. Observed results show that most of the methods used have a proper performance relative to *MiXCR* protocols while others did not correctly extract all the repertoires.