

# A Computational Platform for Gene Expression Analysis

Diogo Teixeira

Supervisors: Rui Camacho<sup>1</sup>, Nuno Fonseca<sup>2</sup>

<sup>1</sup>LIAAD INESC, Porto & DEI FEUP, Universidade do Porto, Porto

<sup>2</sup>EMBL-EBI, Cambridge, UK

July 2014

- 1 Introduction
  - Domain Problem
  - Motivation and Objectives
- 2 Developed Solution
  - Overview
  - RNA-Seq Analysis Pipeline
  - RBP Analysis Pipeline (PBS Finder)
  - Integration
- 3 Case Studies
  - RNA-Seq Analysis Pipeline
  - RBP Analysis Pipeline (PBS Finder)
- 4 Conclusions
  - Objective Fulfilment
  - Future Work

# Domain Problem I

## Introduction

- Molecular biology is a young field of study, with a lot of unknowns and partial knowledge.
- Studying gene expression is crucial to understand the mechanisms that control living organisms.
- Two problems sparked the interest of biologists:
  - understanding differences in gene expression between individuals of the same species;
  - understanding interactions between genes and RNA binding proteins that bind with them.

# Domain Problem II

## Introduction

Studying these problems involves:

- being able to take sequencing reads, align them against a reference genome and perform differential expression analysis;
- being able to take a list of gene identifiers, cross reference them between multiple online platforms and discovering their potential RBPs (as well as additional relevant information);
- being able to uncover implicit relationships in the produced information that might be useful to biologists, using data mining techniques.

# Motivation and Objectives

## Introduction

### **Tools are complex**

Tools for biological data analysis often require a very technical set of skills.

### **Create simpler tools**

Any user should be able to use the tools, with little to no training.

# Motivation and Objectives

## Introduction

### **Tools are complex**

Tools for biological data analysis often require a very technical set of skills.

### **Create simpler tools**

Any user should be able to use the tools, with little to no training.

### **Tasks are repetitive**

Analysing high quantities of data can be repetitive, especially if executed manually.

### **Automate tasks**

Automated systems should perform repetitive tasks, so that users can focus on their work.

# Motivation and Objectives

## Introduction

### **Tools are complex**

Tools for biological data analysis often require a very technical set of skills.

### **Create simpler tools**

Any user should be able to use the tools, with little to no training.

### **Tasks are repetitive**

Analysing high quantities of data can be repetitive, especially if executed manually.

### **Automate tasks**

Automated systems should perform repetitive tasks, so that users can focus on their work.

### **Information is scattered**

Information is easy to acquire, but is often scattered through multiple platforms, services and institutions.

### **Gather information**

Information should be contextually aggregated, allowing for quick access of relevant information.

# Overview

## Developed Solution

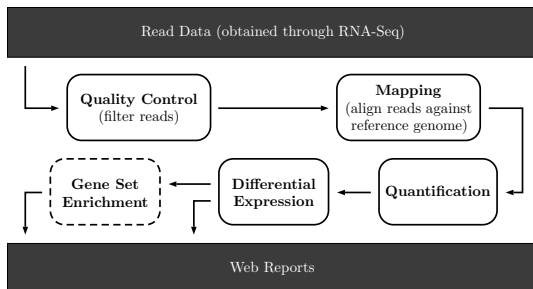
- Develop two distinct pipelines: one for differential expression analysis and another for RBP discovery and analysis.
- Both pipelines should be available through web platforms. These platforms should be able to manage user accounts, user jobs (analysis tasks), results viewing and exporting, etc..
- These platforms should be able to overcome the three problems previously mentioned.



# RNA-Seq Analysis Pipeline I

## Developed Solution

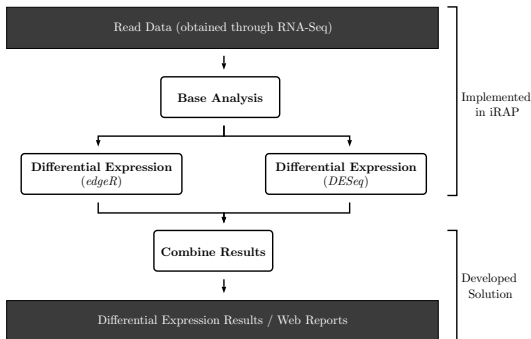
- Uses iRAP (an RNA-Seq analysis pipeline) as its base.
- iRAP is flexible and allows users to choose what tool should be used in each step of the process.
- But, if users have little experience with these tools, how can they choose the best one for the job?



# RNA-Seq Analysis Pipeline II

## Developed Solution

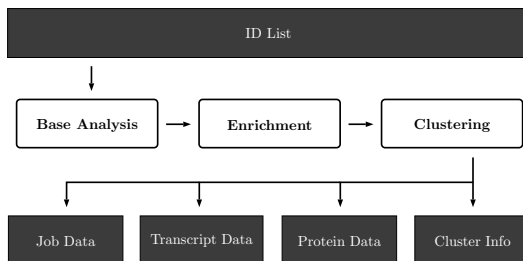
- Uses multiple differential expression analysis tools (edgeR e DESeq).
- Combines the results of those tools.



# RBP Analysis Pipeline (PBS Finder) I

## Developed Solution

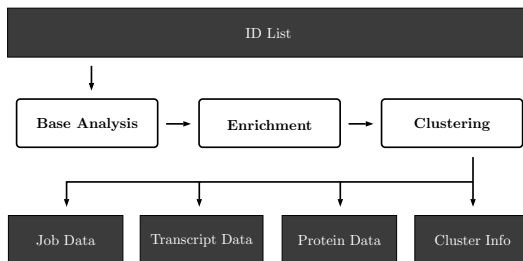
- Uses DB2DB, Ensembl and NCBI to convert identifiers, identify gene species, obtain basic information and extract genetic sequences (5' UTR, 3' UTR, 3' UTR downstream).



# RBP Analysis Pipeline (PBS Finder) II

## Developed Solution

- Uses RBPDB to discover RNA binding proteins (based on the obtained genetic sequences).
- Uses UniProt to enrich the obtained results and performs clustering analysis on those results.



# RBP Analysis Pipeline (PBS Finder) III

## Developed Solution

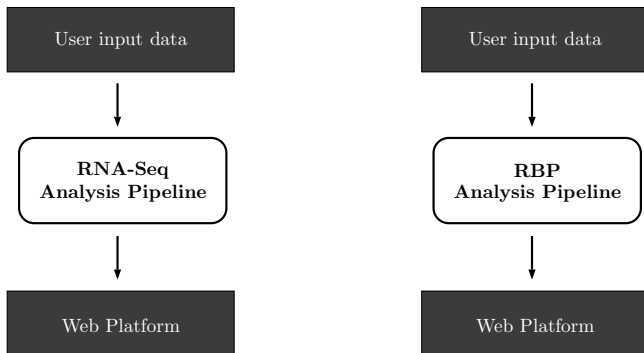
Clustering analysis:

- Uses  $k$ -medoids and hierarchical clustering, both with Jaccard and binary distance matrices.
- Executes every possible combination of clustering setups (alternates algorithms, distance matrices, used features, etc.).
- Results are filtered (acceptable solutions must have a minimum percentage of entries per cluster, clusters must have defining features, etc.).
- Solution quality internally determined based on the average silhouette.

# Integration

## Developed Solution

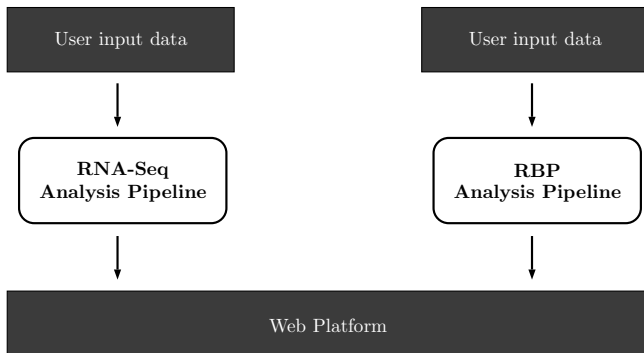
While focusing on aggregation and quick access to information, does it make sense to separate the results into two different platforms?



# Integration

## Developed Solution

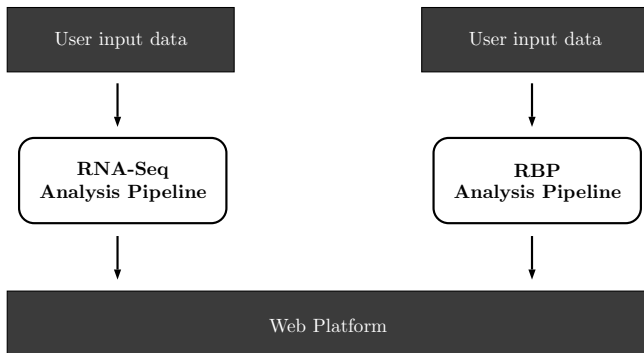
While focusing on aggregation and quick access to information, does it make sense to separate the results into two different platforms?



# Integration

## Developed Solution

A list of differentially expressed genes is not very useful without further information about those genes. Does it make sense for a user to launch a new gene enrichment task by hand?

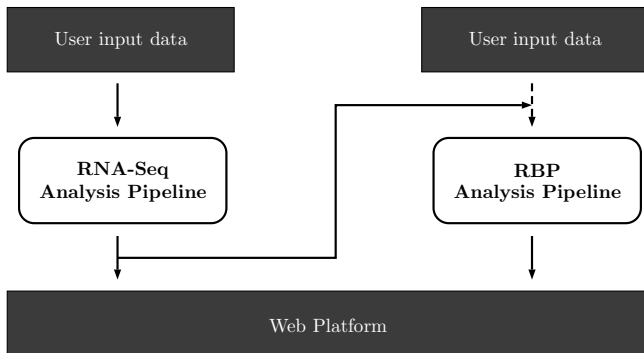




# Integration

## Developed Solution

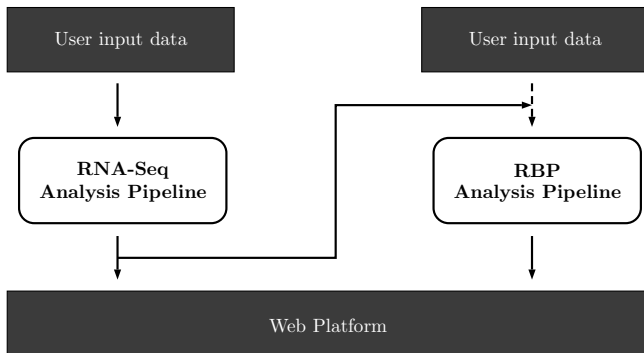
A list of differentially expressed genes is not very useful without further information about those genes. Does it make sense for a user to launch a new gene enrichment task by hand?



# Integration

## Developed Solution

A fully integrated solution: the analysis pipelines can be used separately or automatically executed in sequence; result visualization for both pipelines is isolated.



# RNA-Seq Analysis Pipeline I

## Case Studies

### Objective

- Ascertain if combining the results of multiple tools has impact on the set of differentially expressed genes.

### Data set

- Reproduction of ArrayExpress experiment E-GEOD-48829 (*Escherichia coli*).
- Reference genome obtained from Ensembl Genomes and read data obtained from ENA Sequence Read Archive.

# RNA-Seq Analysis Pipeline II

## Case Studies

Results (number of differentially expressed genes)

	<i>Raw results</i>	<i>Filtered results</i>	<i>Combined results</i>
<i>edgeR</i>	4494	386	191
<i>DESeq</i>	4494	204	

## Conclusions

- Combining results impacts the final differentially expressed gene list by reducing its size.
- The combined results will hopefully give researchers an higher confidence in the experimental results.

# RBP Analysis Pipeline (PBS Finder) I

## Case Studies

### Objectives

- Assess the general usefulness of PBS Finder.
- Compare PBS Finder with the existing techniques of manual analysis.
- Assess the impact of differences in hardware performance in the overall performance of the platform.

### Data set

- 23 genes from the *RhoGTPase* family (*Rattus norvegicus*) provided by IBMC.

# RBP Analysis Pipeline (PBS Finder) II

## Case Studies

Results (expert estimation of 30 minutes per gene analysed)

<i>Number of IDs</i>	<i>Machine1</i>	<i>Machine2</i>	<i>Manual method</i>
100	9m 56s	11m 1s	$\approx 50h$
500	41m 47s	55m 51s	$\approx 250h$
900	1h 33m 32s	2h 7m 4s	$\approx 450h$

## Conclusions

- PBS Finder can reproduce the results an expert would get.
- Months worth of an expert's manual work can be accomplished in a few hours.
- While hardware performance has a significant impact on analysis time, the platform achieves satisfactory performance on personal computer-level hardware.

# Objective Fulfilment

## Conclusions

- RBP analysis pipeline and web platform (PBS Finder) implemented and tested. PBS Finder has been in production for several months; during this time it was thoroughly tested by IBMC experts.
- RNA-Seq analysis pipeline implemented and tested (iRAP deployed and result consolidation tool implemented).
- Integration of both tools could not be accomplished.

# Future Work

## Conclusions

- Fully integrate the RNA-Seq analysis pipeline with the web platform (automatic job configuration, result visualization, etc.).
- Study the requirements for deploying the platform in large scale, and assess the feasibility of making it available internet-wide.



- 1 Introduction
  - Domain Problem
  - Motivation and Objectives
- 2 Developed Solution
  - Overview
  - RNA-Seq Analysis Pipeline
  - RBP Analysis Pipeline (PBS Finder)
  - Integration
- 3 Case Studies
  - RNA-Seq Analysis Pipeline
  - RBP Analysis Pipeline (PBS Finder)
- 4 Conclusions
  - Objective Fulfilment
  - Future Work

# A Computational Platform for Gene Expression Analysis

Diogo Teixeira

Supervisors: Rui Camacho<sup>1</sup>, Nuno Fonseca<sup>2</sup>

<sup>1</sup>LIAAD INESC, Porto & DEI FEUP, Universidade do Porto, Porto

<sup>2</sup>EMBL-EBI, Cambridge, UK

July 2014