

A Computational Platform for Gene Expression Analysis

Diogo Teixeira¹

Supervisors: Rui Camacho², Nuno Fonseca³

¹Check affiliation

²Check affiliation

³Check affiliation

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.
- We focused on two different areas:
 - differential expression analysis;
 - RNA-binding protein (RBP) discovery and analysis.

Domain Problem II

Introduction

Three distinct problems:

- Read alignment against a reference genome and differential expression analysis on the aligned data.
- RBP discovery, analysis and information enrichment.
- Further result analysis using data mining techniques.

Motivation and Objectives

Introduction

Tools are complex

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

Motivation and Objectives

Introduction

Tools are complex

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

Tasks are repetitive

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

Motivation and Objectives

Introduction

Tools are complex

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

Tasks are repetitive

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

Information is scattered

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

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.

Information is scattered

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

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.

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

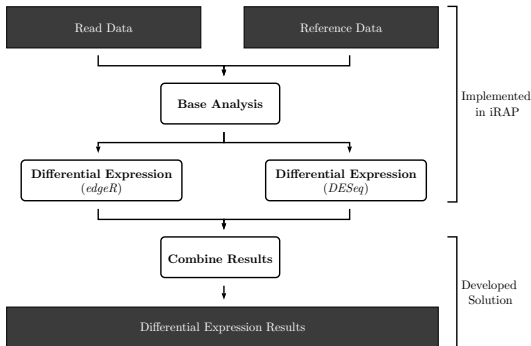
Developed Solution

- Two distinct problems warrant two different solutions.
- The developed system should be available anywhere, through the internet.
- The system should be as modular as possible, to allow future extensions.

RNA-Seq Analysis Pipeline I

Developed Solution

- Uses iRAP as the analysis pipeline.
- Conducts multiple differential expression analyses with different tools.
- Combines results from multiple tools.



RNA-Seq Analysis Pipeline II

Developed Solution

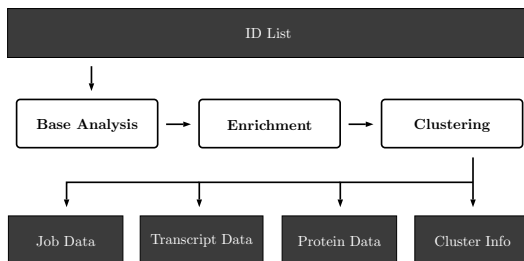
Additional features:

- Access to iRAP's web reports and gene browser, along with result visualization in the web interface.
- Synchronization with Ensembl's reference genome repositories.
- Graphical job configuration.
- Possibility to easily include other differential expression tools.

RBP Analysis Pipeline (PBS Finder) I

Developed Solution

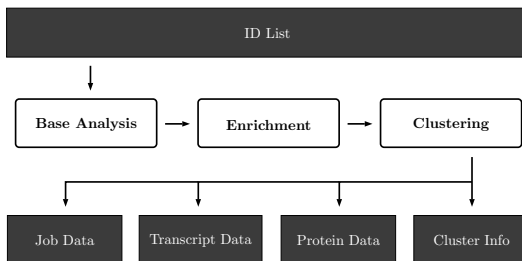
- Uses Ensembl and NCBI to identify gene species, obtain basic information and extract genetic sequences (5' UTR, 3' UTR, 3' UTR downstream).



RBP Analysis Pipeline (PBS Finder) II

Developed Solution

- Uses RBPDB to discovery RNA binding proteins based on the obtained 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.

RBP Analysis Pipeline (PBS Finder) IV

Developed Solution

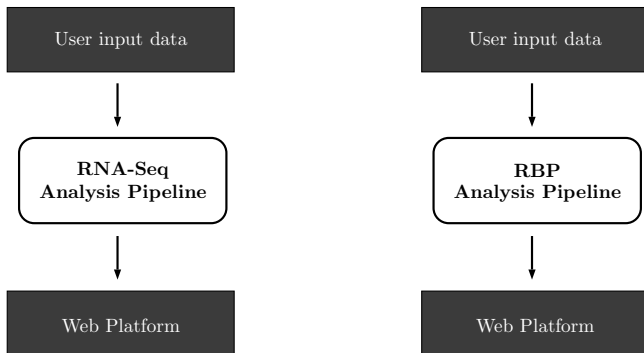
Additional features:

- User account and job management system.
- References to external platforms with relevant information based on context.
- Support for multiple identifier notations (Ensembl, Entrez, RefSeq and GenBank).
- Visualization of defining features for each cluster.
- Job completed notification system.

Integration

Developed Solution

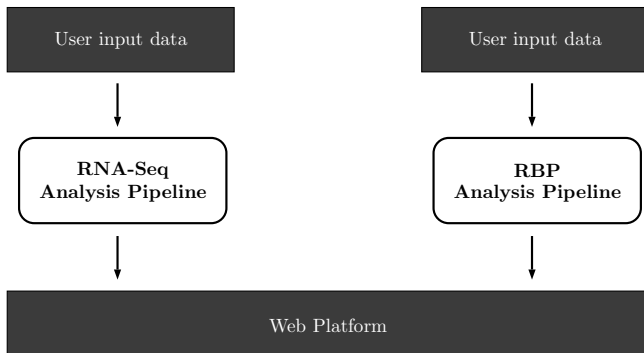
While focusing on aggregation and quick access to information, does it make sense to separate the results into two different platforms?



Integration

Developed Solution

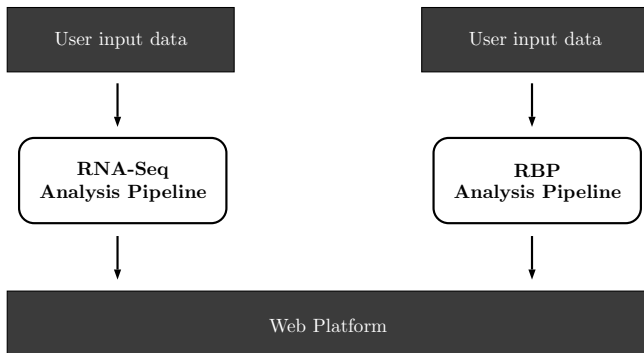
While focusing on aggregation and quick access to information, does it make sense to separate the results into two different platforms?



Integration

Developed Solution

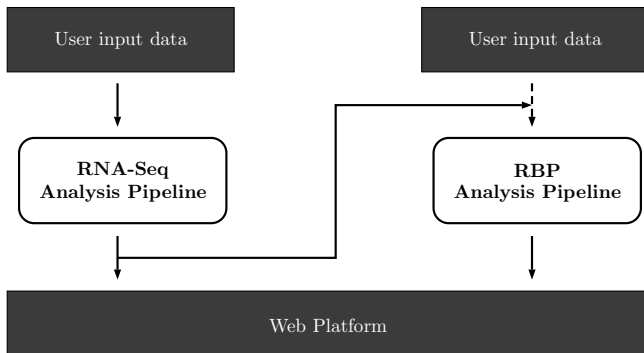
A list of differentially expressed genes is not very useful without further information about those genes. Does it make sense for a user to launch a new gene enrichment task by hand?



Integration

Developed Solution

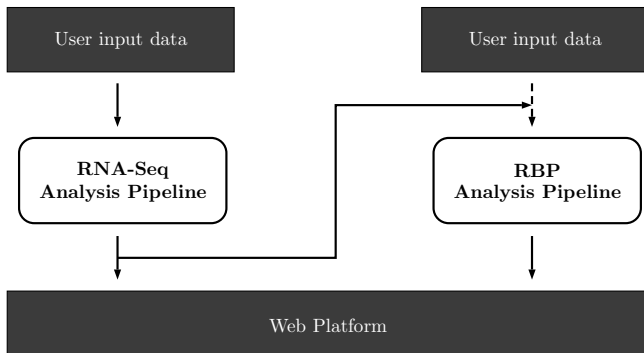
A list of differentially expressed genes is not very useful without further information about those genes. Does it make sense for a user to launch a new gene enrichment task by hand?



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

Case Studies

NOTES:

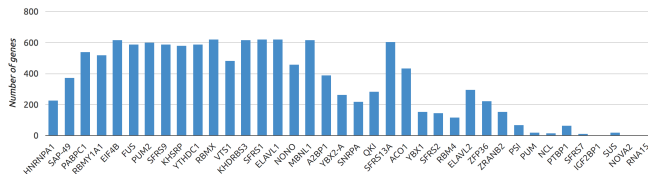
- Refer objectives, data and experimental method.
- Refer results.

RBP Analysis Pipeline (PBS Finder) I

Case Studies

Case study results viewed in PBS Finder

RBP FREQUENCY



Function	Proteins	Species	Gene (view original)	Transcript	Protein	HNRNP1	SAP-49	PABPC1	RBMY1A1	EIF4B	FUS	PUM2	SFRS9	KHSRP	YTHDC1	RBMX
		Rattus norvegicus (Rat)	ENSRNOG00000018923 (Npm1)	ENSRNOT00000025575 (Npm1-201)	P13084 (NPM1)											
		Rattus norvegicus (Rat)	ENSRNOG00000007345 (Amot)	ENSRNOT00000049482 (Amot-201)	D4A9Q2 (AMOT)											
		Rattus norvegicus (Rat)	ENSRNOG00000021373 (Gpd1)	ENSRNOT00000038757 (Cox14-201)	Q5XFV8 (COX14)											
		Cluster 4		ENSRNOG00000019213 (Gpd1)	ENSRNOT00000026199 (Gpd1-201)											
		DISTINCT PROTEINS NOVA2, RNA15														
		Rattus norvegicus (Rat)	ENSRNOG00000029512	ENSRNOT00000051134	N/A											

RBP Analysis Pipeline (PBS Finder) II

Case Studies

NOTES:

- Refer objectives, data and experimental method.
- Refer results.
- Show screenshot.
- Show table.

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.

- 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 Teixeira¹

Supervisors: Rui Camacho², Nuno Fonseca³

¹Check affiliation

²Check affiliation

³Check affiliation

July 2014