Master Thesis

# Recommender System for

# Galaxy Tools and Workflows

## (Find similar tools and predict next tools in workflows)

Anup Kumar

Examiners:  Prof. Dr. Rolf Backofen and Prof. Dr. Wolfgang Hess

Adviser:    Dr. Björn Grüning

University of Freiburg
Faculty of Engineering
Department of Computer Science
Chair for Bioinformatics

July, 2018

**Thesis period**

10. 01. 2018 – 09. 07. 2018

**Examiners**

Prof. Dr. Rolf Backofen and Prof. Dr. Wolfgang Hess

**Adviser**

Dr. Björn Grüning

# Declaration

I hereby declare, that I am the sole author and composer of my thesis and that no other sources or learning aids, other than those listed, have been used. Furthermore, I declare that I have acknowledged the work of others by providing detailed references of said work.

I hereby also declare, that my Thesis has not been prepared for another examination or assignment, either wholly or excerpts thereof.

_____

Place, Date

_____

Signature

# Acknowledgment

# Abstract

The study explores two concepts to devise a recommendation system for Galaxy. One idea is to find similar tools for each tool and another is to predict a set of possible next tools in workflows.

To find similarities among tools, we need to extract information from each tool under multiple categories like a tool's name, description, input and output file types and helptext. We take into account these categories one by one and compute similarity matrices. Each row in a similarity matrix keeps similarity scores of one tool against all other tools. These similarity scores depend on the similarity measure used to compute the score between a pair of tools. We compute three such similarity matrices, one each for input and output file types, name and description and helptext. To combine these matrices, one simple solution is to compute an average. But, assigning equal importance weights to each matrix might be sub-optimal. To find an optimal combination, we use optimization to learn importance weights for the corresponding rows for each tool in the similarity matrices. To define a loss function, we use a true similarity value based on the similarity measures. We take an array of 1.0 as true value because we use jaccard index and cosine similarity as similarity measures which give a positive real number between 0 and 1.

Next task analyzes workflows to predict a set of next tools at each stage of creating workflows. While creating workflows, it would be convenient to leaf through a set of next possible tools as a guide. It can assist the less experienced (Galaxy) users in creating workflows when they are unsure about which tools can further be joined. In addition, it can curtail the time taken in creating a workflow. To achieve that, we need to learn the connections among tools in order to be able to predict the next possible ones. To preprocess each workflow, we compute all the paths bridging the starting and ending tools. These paths contain a set of connected tools. We follow a classification approach to predict the next tools. For the classification task, we use a variant of neural networks (long short-term memory). It performs well for learning long range dependencies (tools connections). We report the accuracy as precision.
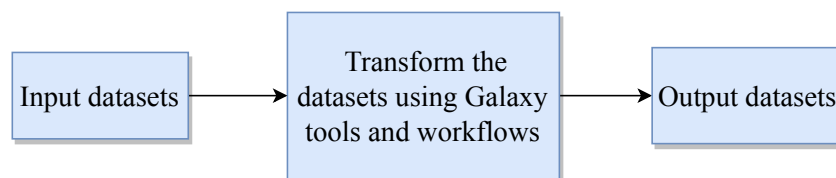
# Zusammenfassung

# Contents

# 1 Introduction

### 1.0.1 Galaxy

Galaxy [1] is an open-source biological data processing and research platform. It supports numerous types of extensively used biological data formats like FASTA, FASTAQ, GFF, PDB and many more. To process these datasets, Galaxy offers tools and workflows which either transform these datasets from one type to another or manipulate the data. A simple example of data processing is to merge two compatible datasets to make one. Another example can be to reverse complement a sequence of nucleotides. [2]



**Figure 1: Galaxy tools and workflows** This images shows the basic form in which tools and workflows carry out a task

A tool is a data-transforming entity which allows one or more types of datasets, transforms these datasets and releases output datasets in one or more types.

The workflows are data processing pipelines. A set of tools joined one after another constitute a workflow. A workflow can have one or more starting and ending tools. The tools which are connected

### 1.0.2 Galaxy tools

### 1.0.3 Tool categories

### 1.0.4 Motivation

---

[1] https://usegalaxy.eu/
[2] https://usegalaxy.eu/?tool_id=MAF_Reverse_Complement_1&version=1.0.1&__identifer=zmk9dx9ivbk

# 2 Approach

## 2.1 Data preprocessing

### 2.1.1 Tools attributes

### 2.1.2 Data extraction

### 2.1.3 Clean data

**Refine tokens**

## 2.2 Word embeddings

### 2.2.1 Latent semantic indexing

### 2.2.2 Paragraph vectors

## 2.3 Similarity measures

**Jaccard index**

**Cosine similarity**

## 2.4 Optimization

**Gradient descent**

**Backtracking line search**

# 3  Experiments

# 4 Results and Analysis

# 5 Conclusion

# 6 Future Work