# GlycoPP V2.0 FRAMEWORK
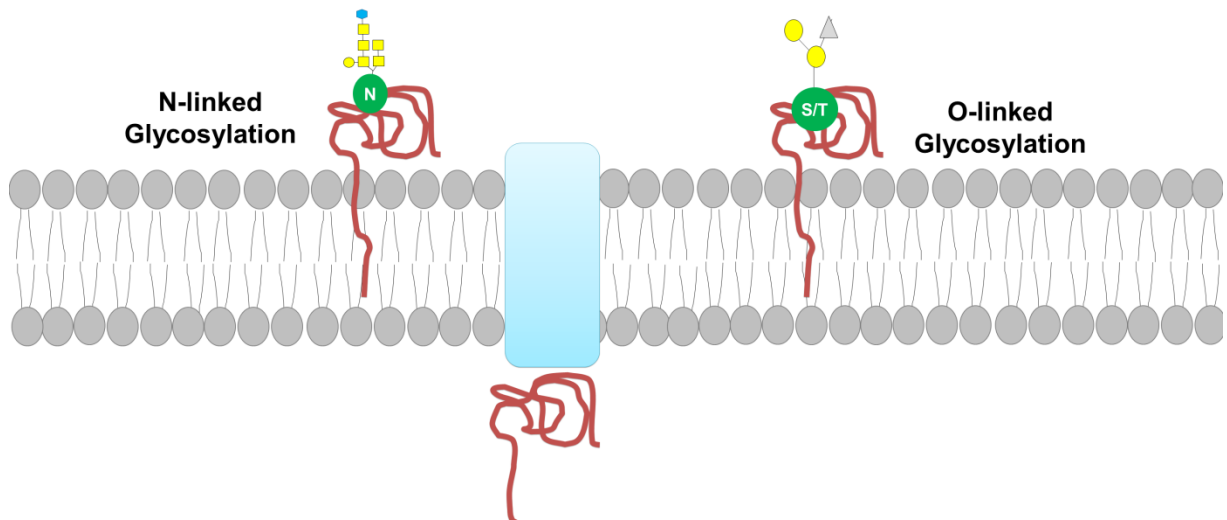


**Bioinformatics Centre, CSIR-Institute of Microbial Technology,**

**Chandigarh – 110036, India**

# Introduction

GlycoPP V2.0 is a highly accurate glycosylation prediction made available for the analysis of prokaryotic protein sequences on the web based Galaxy Platform. GlycoPP prediction programmes are trained on the largest available and an extensive dataset of N- and O-linked glycosites extracted from experimentally characterized glycoproteins of prokaryotes as obtained from ProGlycProt V2.0(http://www.proglycprot.org/).

GlycoPP V2.0 is an enhanced and updated version of our GlycoPP V1.0 (http://crdd.osdd.net/raghava/glycopp). The workflow system is implemented using the open source workflow architecture, Galaxy. GlycoPP2 is freely available and can be accessed at https://ab-openlab.csir.res.in/alkarao/glycopp2/#

There are three modules in GlycoPP2, described in this manual. All the modules are accessible without registering to the system. However, for maintaining user-sessions it is recommended that anyone who is interested in creating data or task intensive workflow should register (**Figure** ). The benefit of registration includes user sessions, saved histories, visualization, generation and execution of workflow and many others.



*Figure 1: GlycoPP V2.0 Main page*



**Figure 2:** *GlycoPP2 Login Page*

All the modules of GlycoPP2 and some default modules by Galaxy are accessible through a web-based interface which has following components (*Figure* ).

**Navigation Panel:** It provides the links to major components of the server like Tools Page (Analyze Data), Workflow System, Shared Libraries, Visualization, Help Section and User Login/Registration.

**Tool Panel:** This panel lists all the tools available in GlycoPP2 along with default utilities in Galaxy.

**Detail Panel (Canvas):** This panel displays the interface of all the tools along with Input Parameters required to run a tool. It also provides help and examples to run a tool. This panel also displays the Output of a tool after its execution when user clicks on the eye 👁 icon show in **History Panel**.

**History Panel:** This panel shows the information about the tools which are executed by a user. The information can include result after completion of a tool execution or error generated while running the tool. The workflow(s) are generated by extracting tasks from history panel.



*Figure 3:* *GlycoPP2 Homepage*

# I. Navigation Panel

## 1. Analyze Data

The data analysis page is where everything happens. There, you can run any available tools on the data, run complete workflows, browse or download a result, and share files with other users. It is

the default page when you open Galaxy in your browser, but you can also access it any time by clicking on "Analyze Data" in the Navigation Panel.



## 2. Workflow

Workflows are analyses that are intended to be executed (one or more times) with different user-provided input Datasets. Workflow can be reused over and over, not only reducing tedious work, but enhancing reproducibility by applying the same exact methods to all of your data. Workflow is nothing but creating pipeline, user can use it again and again or user can published it.

Workflow can be created through navigation panel or from tool panel. In workflow section user can create workflow or can upload or import the workflow. The canvas is where inputs, tools, and noodles are added and connected as you build and modify your workflow (Figure 4). Selecting Edit opens the workflow editor view (Figure 3). The navigator provides a full view of your workflow in a condensed format (Figure 4). Accessed by clicking on the gear icon on the right side of the center Workflow Canvas upper bar, the workflow editor menu (Figure 3) is for global editor actions. It consists of Save, Run, Edit Attributes, Auto re-layout, Close

1. Workflow



**Figure 5:** *Options for workflow*

The following example of workflow shows the "Prediction of N- and O-linked glycosylation prediction". The prediction workflow can use the four implemented svm model for N-linked glycosyation and six svm model for O-linked glycosylation

Run Workflow for N- and O-linked glycosylation



**Figure 6:** *Workflow overview*

## 3. Shared data library

Data libraries are collections of Datasets that are accessible from within a Galaxy instance. Libraries are designed for sharing datasets in between users or groups. The data library of GlycoPP2 consists of prokaryotic glycoproteins list protein list used in SVM model generation. Some of the actions that can be performed on data libraries are accessed by clicking the pop-up menu icon just right of the data library name.

- View Information –Shows the information about dataset.

- Import this dataset into your current history - this creates an item in your current

history on which you can perform analysis. The item is a pointer to the library dataset disk file, so the file is not copied on disk.

- Download this dataset - this allows you to download a local copy of the dataset.

1. Shared Data library



*Figure 6: Viewed for shared library*

## 4 Help

The help section of Galaxy consists of Support, Search, Mailing List, Videos, Wiki and How to cite Galaxy. User can fine user manual in help section.

## 5. User

Login option and register option can get in user section. It is recommended that user register their account before using framework. Although unregistered users have access to tools available but their history is stored temporarily. On the other hand, registered users can save and retrieve their results in history panel later too.

# Tool Panel – GlycoPP2 tools

The user can get tools in tool panel. There are two categories of tools galaxy inbuilt tool and GlycoPP2 tools. Galaxy tools consists of Data importing, Manipulation, Filtering, Sorting, Format conversion etc. GlycoPP2 tool are specific for finding N- and O-linked glycosite in prokaryotic protein sequence

## 1. Importing data to the GlycoPP2

A user can upload the data using the Galaxy tool **Get Data.** The uploaded data can be used for GlycoPP V2.0 Analysis.The following figure shows the file upload method.

**1. Get data / Upload file**



*Figure 8: Get Data screen*

## 2. Examples of Glycopp V2.0 Analysis

These have example fasta and multifasta file reloded for the prediction analysis of N- and O-linked glycosylation SVM model. BPP, BPP+ASA, BPP+SS and BPP+ASA+SS for N-linked Glycosylation and CTD, PAAC, SER, CPP+SS, DPC+SS and DPC+ASA for O-linked Glycosylation



*Figure 9: Example of GlycoPP V2.0 Analsis Screen*

## a) Example of N-linked Glycosylation

### BPP example Binary Profile Pattern based prediction example



**Figure-10 :** *Example of GlycoPP V2.0 Analsis overview*



**Figure11:** *Example of GlycoPP V2.0 Analsis output for BPP (Binary Profile Pattern) forN-linked glycosylation*

## BPP+ASA example (Binary Profile Composition+ Accessible Surface Area) based prediction example



**Figure 12:** *Example of GlycoPP V2.0 Analsis overview*



**Figure 13:** *Example of GlycoPP V2.0 Analsis output for BPP+ASA (Binary Profile Pattern+ Accessible Surface Are) forN-linked glycosylation*

## BPP+ASA+SS example (Binary Profile Composition+ Accessible Surface Area + Secondary Structure) based prediction example



**Figure 14:** *Example of GlycoPP V2.0 Analsis overview*



**Figure 15: Example of GlycoPP V2.0 Analsis output for BPP+ASA+SS example (Binary Profile Composition+ Accessible Surface Area + Secondary Structure) forN-linked glycosylation**

## BPP+SS example (Binary Profile Composition+ Secondary Structure) based prediction example



**Figure 16:** *Example of GlycoPP V2.0 Analsis overview*



**Figure 17:** *Example of GlycoPP V2.0 Analsis output for BPP+ASA+SS example (Binary Profile Composition+ Accessible Surface Area + Secondary Structure) forN-linked glycosylation*

## b) Example of O-linked Glycosylation

### CTD_example Conjoint Triad Descriptors Example



**Figure 18:** *Example of GlycoPP V2.0 Analsis overview*



**Figure 19: Example of GlycoPP V2.0 Analsis output for CTD (Conjoint Triad Distributor) for O-linked glycosylation**

## PAAC_example Pseudo Amino Acid Composition based prediction example



**Figure 21:** *Example of GlycoPP V2.0 Analsis overview*



**Figure 22:** *Example of GlycoPP V2.0 Analsis output for PAAC (Pseudo Amino Acid Composition) forO-linked glycosylation*

## SER example Shannon Entropy of Residues based prediction example



**Figure 23:** *Example of GlycoPP V2.0 Analsis overview*



**Figure 24:** *Example of GlycoPP V2.0 Analsis output for **SER (Shannon Entropy of Residues)** for O-linked glycosylation*

## CPP+SS example (Composition Profile of Patterns + Secondary Structure) based prediction example



**Figure 25:** *Example of GlycoPP V2.0 Analsis overview*



**Figure 26: Example of GlycoPP V2.0 Analsis output for CPP+SS (Composition Profile of Patterns + Secondary Structure) forO-linked glycosylation**

## DPC+SS example (Dipeptide Composition + Secondary Structure)based prediction example



**Figure 27:** *Example of GlycoPP V2.0 Analsis overview*



**Figure 28:** *Example of GlycoPP V2.0 Analsis output for DPC+SS(Dipeptide Composition + Secondary Structure)forO-linked glycosylation*

## DPC+ASA example (Dipeptide Composition + Accesible surface Area)based prediction example



**Figure 29:** *Example of GlycoPP V2.0 Analsis overview*



**Figure 30:** *Example of GlycoPP V2.0 Analsis output for DPC+ASA example (Dipeptide Composition + Accesible surface Area ) forO-linked glycosylation*

## 2. Glycopp V2.0 Analysis

A user can upload the fasta file using the Galaxy tool **Get Data as we shown in figure 8** this fasta files can be used for GlycoPP V2.0 Analysis.The following figure shows the file upload method.

Get Data Galaxy tool for uploading the fasta protein sequence for prediction of N- and O_linked glycosite



**Figure31 :** *GlycoPP V2.0 Analsis overview*

## a) N-linked Glycosylation

### BPP  Binary Profile Pattern based prediction



**Figure-32 :** *GlycoPP V2.0 Analsis overview for BPP*



**Figure33:** *GlycoPP V2.0 Analsis output for BPP (Binary Profile Pattern) forN-linked glycosylation*

## BPP+ASA  (Binary Profile Composition+ Accessible Surface Area) based prediction



**Figure 34:** *GlycoPP V2.0 Analsis*



**Figure 35:** *GlycoPP V2.0 Analsis output for BPP+ASA (Binary Profile Pattern+ Accessible Surface Are) forN-linked glycosylation*

## BPP+ASA+SS  (Binary Profile Composition+ Accessible Surface Area + Secondary Structure) based prediction



**Figure 36:** *GlycoPP V2.0 Analsis overview*



**Figure 37** *GlycoPP V2.0 Analsis output for BPP+ASA+SS  (Binary Profile Composition+ Accessible Surface Area + Secondary Structure) forN-linked glycosylation*

## BPP+SS (Binary Profile Composition+ Secondary Structure) based prediction



**Figure 38:** *GlycoPP V2.0 Analsis overview*



**Figure 39** *GlycoPP V2.0 Analsis output for BPP+ASA+SS (Binary Profile Composition+ Accessible Surface Area + Secondary Structure) forN-linked glycosylation*

## b)  O-linked Glycosylation

## CTD Conjoint Triad Descriptors



**Figure 40:**  *GlycoPP V2.0 Analsis overview*



**Figure 41:**  *GlycoPP V2.0 Analsis output for CTD (Conjoint Triad Distributor) forO-linked glycosylation*

## PAAC Pseudo Amino Acid Composition based prediction



**Figure 42:** *GlycoPP V2.0 Analsis overview*



**Figure 43:** *GlycoPP V2.0 Analsis output for PAAC (Pseudo Amino Acid Composition) forO-linked glycosylation*

## SER  Shannon Entropy of Residues based prediction



**Figure 44:** *GlycoPP V2.0 Analsis overview*



**Figure 45:  GlycoPP V2.0 Analsis output for SER (Shannon Entropy of Residues) forO-linked glycosylation**

## CPP+SS (Composition Profile of Patterns + Secondary Structure) based prediction



*Figure 46:* *GlycoPP V2.0 Analsis overview*



*Figure 47:* *GlycoPP V2.0 Analsis output for CPP+SS (Composition Profile of Patterns + Secondary Structure) forO-linked glycosylation*

## DPC+SS (Dipeptide Composition + Secondary Structure)based prediction



**Figure 48:** *GlycoPP V2.0 Analsis overview*



**Figure 49:** *GlycoPP V2.0 Analsis output for DPC+SS(Dipeptide Composition + Secondary Structure)forO-linked glycosylation*

## DPC+ASA  (Dipeptide Composition + Accesible surface Area)based prediction



**Figure 50:** *GlycoPP V2.0 Analsis overview*



**Figure 51:  GlycoPP V2.0 Analsis output for DPC+ASA  (Dipeptide Composition + Accesible surface Area ) forO-linked glycosylation**

# III. Detail Panel

This panel displays the interface of all the tools along with Input Parameters required to run a tool. It also provides help and examples to run a tool. This panel also displays the Output of a tool after its execution when user clicks on the eye ◉ icon show in **History Panel**.

# IV. History Panel

When data is uploaded from your computer or analysis is done on existing data using Galaxy, each output from those steps generates a dataset. These datasets (and the output datasets from later analysis on them) are stored by Galaxy in **Histories**.

Users that have registered an account and logged in can have many histories and the history panel allows switching between them and creating new ones.