# GlyGen

## *Computational and Informatics Resources for Glycosciences*

*René Ranzinger*

THE GEORGE WASHINGTON UNIVERSITY
WASHINGTON, DC

**PI: Raja Mazumder**
*The George Washington University*

Glycoscience

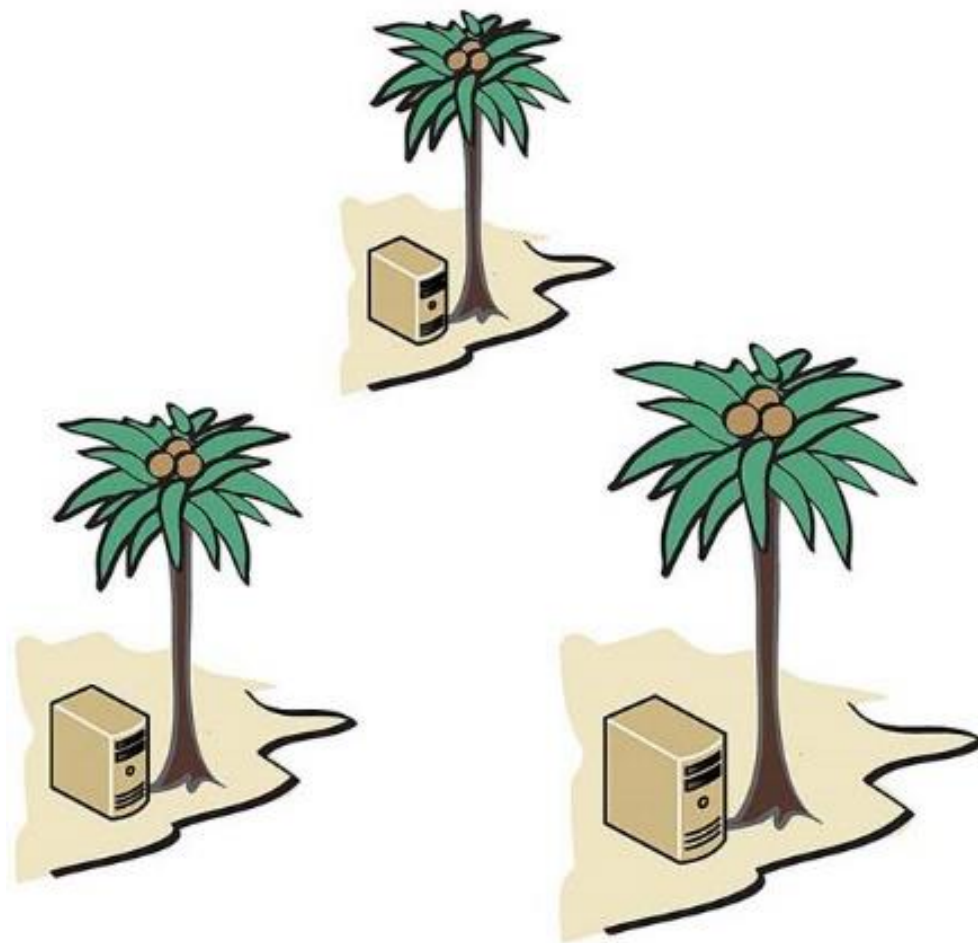**NIH Common Fund Glycoscience Program**
1U01GM125267-01

UNIVERSITY OF GEORGIA

**PI: William York**
CCRC*, University of Georgia*

**GlyGen.org**

https://twitter.com/gly_gen

**GlyGen**

- Get all **enzymes** that my of been involved in *synthesis of the glycans* on <u>protein X</u>.
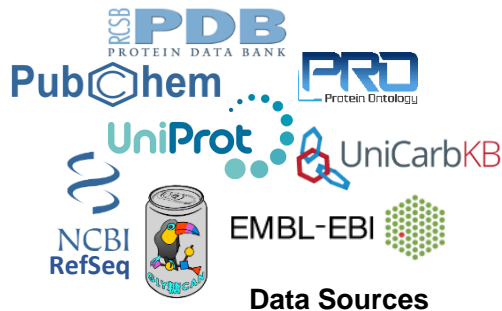
- **Integration of glycobiology-related information** from diverse domains/databases

- **Creation of an intuitive web portal** to browse and search for knowledge in glycobiology

- Provide **free and standardized access** to the integrated datasets

- **Developing essential new information resources**, including:
  - An open, comprehensive **Glycan microarray data repository**
  - A **Glycan Naming Ontology (GNOme)** that facilitates interpretation of incomplete structural information in the context of biological function
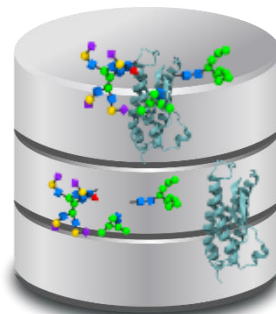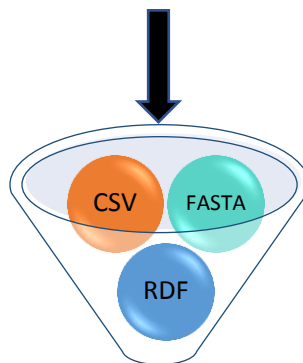
- Data integration

- Data access

- Data sharing

# Architecture



www.glygen.org

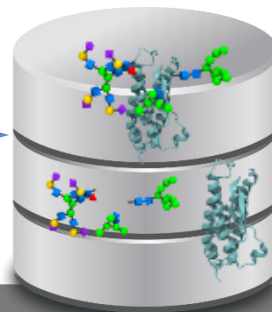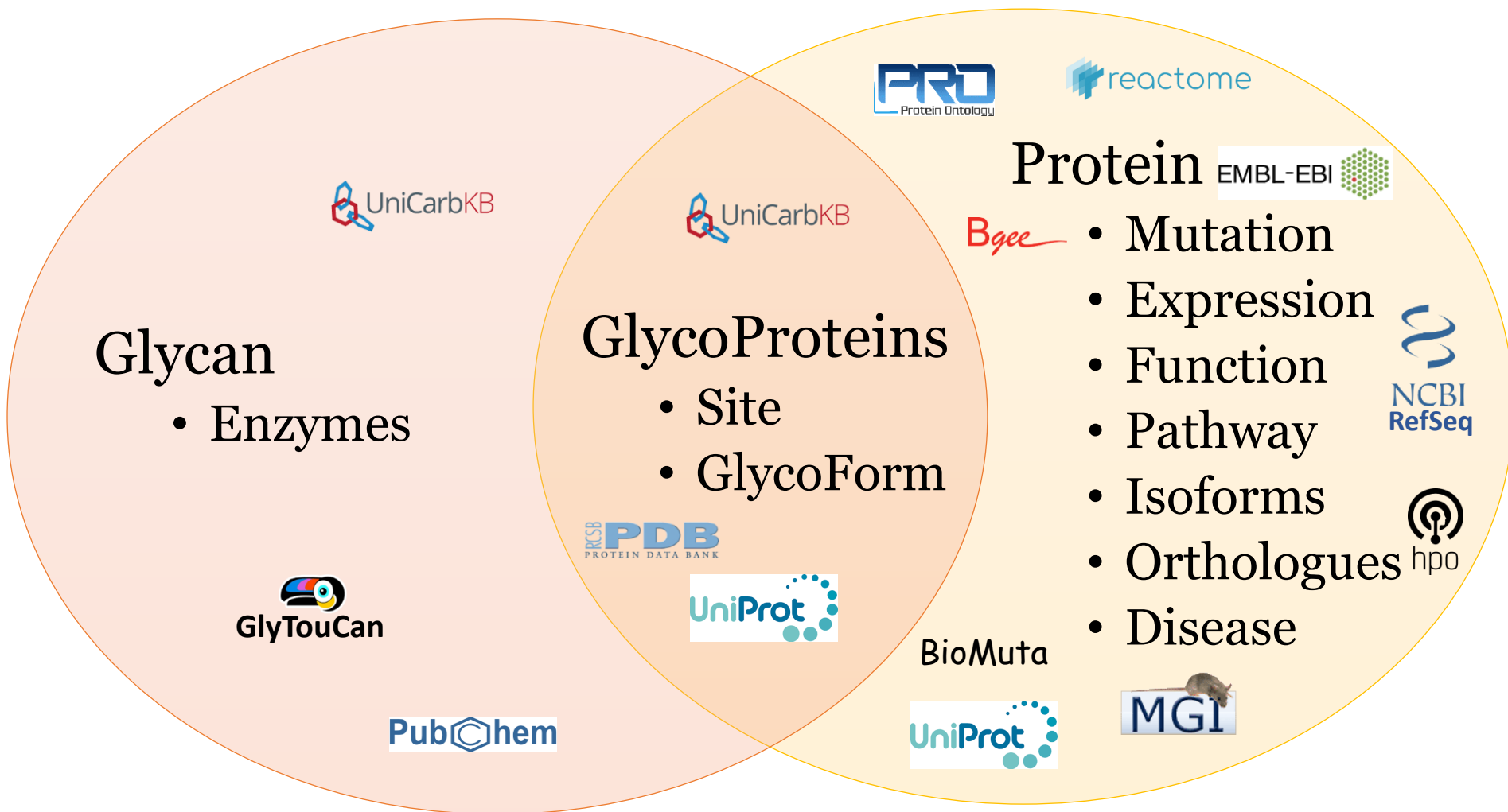sparql.glygen.org

api.glygen.org

data.glygen.org

UniProt
Pub Chem
MGI
hpo
DISEASE ONTOLOGY
Onco MX
PRO Protein Ontology
UniCarbKB
PDB PROTEIN DATA BANK
NCBI
RefSeq
reactome
Data Generators
EMBL-EBI
GlyTouCan
oma BROWSER
Human Phenotype Ontology

**GlyGen**



Glycan
- Enzymes

GlycoProteins
- Site
- GlycoForm

Protein
- Mutation
- Expression
- Function
- Pathway
- Isoforms
- Orthologues
- Disease

PRO Protein Ontology

reactome

EMBL-EBI

Bgee

NCBI RefSeq

hpo

BioMuta

UniProt

MGI

UniCarbKB

UniCarbKB

RCSB PDB PROTEIN DATA BANK

UniProt

GlyTouCan

PubChem

- BioCompute is a community driven project to build a framework to standardize bioinformatics computations and analyses communication.

- **The dataset readme provides the following info** –
  - Provenance, Authoring and Versioning details
  - Usability and description of the dataset
  - Detailed steps required to create the dataset
  - Input and Output files and their formats
  - Software and computational platform requirements
  - Tools, scripts and codes required for creating the dataset
  - Content description of the dataset
  - Dataset statistics

# GlyGen  ≡  **BioCompute Object example**

**BioCompute** Objects

**Provenance Domain**
  **ID: http://data.glygen.org/DSBCO_000038**
  **Name:** human_proteoform_glycosylation_sites_unicarbkb_glytoucan.csv
  **Title:** Glycosylation Sites [UniCarbKB]
  **Version:** 1.0
  **Created:** 2018-02-21T14:46:55-5:00
  **Created by:** Rahi Navelkar [rsn13@gwu.edu; Jeet Vora [jeetvora@gwu.edu]
  **Digital Signature:** RYFNNKE22594E007JKV457
  **Review status:** Approved
  **Contribution:** Matthew Campbell[contributedBy], Robel Kahsay [curatedBy], Rahi Navelkar [createdBy]
  **License:** Data – Attribution 4.0 International CC BY 4.0 [https://creativecommons.org/licenses/by/4.0/]
          Scripts – GNU General Public License v3.0 [https://www.gnu.org/licenses/gpl-3.0.en.html]

| **Provenance, Authoring, Versioning and License Information** |

**Usability Domain**
  List of human [taxid:9606] proteins with information on glycosylation sites from
  UniCarbKB database. The file also includes GlyTouCan accessions and UniCarbKB structure ids for associated glycan structures.
  https://academic.oup.com/nar/article/42/D1/D215/1052197 ,https://doi.org/10.1093/nar/gkt1128]

| **Describes Intended Use Cases** |

**Description Domain**
  **Keywords:** protein, canonical, glycosylation, glycan
  **Pipeline Steps:**
    Step 1: The input file was retrieved directly from source.
    Step 2: The UniProtKB protein accessions in the input file were mapped to UniProt canonical accessions …..
    Step 3: The glycosylation type [linkage type] was retrieved through UniCarbKB structure webpage using scripts
            [make-proteoform_glycosylation_sites_unicarbkb_glytoucan-csv-step2a.py,………
    Step 4: The file was processed for quality check using a python Records which fall under one or
            more following criteria's are flagged and eliminated [eliminated records can be accessed using log
            file]

| **Human Readable Description of Process** |

**Execution Domain:**
  **Script Access Type:** Text
  **Scripts:** make-proteoform_glycosylation_sites_unicarbkb_glytoucan-csv-step2a.py, make-
          proteoform_glycosylation_sites_unicarbkb_glytoucan-csv-step2b.py,
  **Script Location:** https://github.com/glygener/glygen-backend/blob/master/integration/
  **Script Driver:** manual
  **Platform:** CentOS7

| **Explicit Computational Inputs and Processes** |

**Software Prerequisites:**
  **Name:** Python
  **Version:** 2.7.13

| **Software Requirements** |

(……………)

https://data.glygen.org                    GlyGen.org

**BioCompute Objects**

```
I/O Domain
  Input Subdomain:
    uri: https://data.glygen.org/ln2wwwdata/source/human_glytoucan_140918_2018_10_31_02_17_32.txt
    filename: human_glytoucan_140918_2018_10_31_02_17_32.txt

    uri: https://data.glygen.org/ln2wwwdata/reviewed/human_protein_all.fasta
    filename : human_protein_all.fasta
```

**Information on Input files required in the pipeline steps to create the dataset**

```
Output Subdomain:
    mediatype: csv
    uri: https://data.glygen.org/ln2wwwdata/reviewed/human_proteoform_glycosylation_sites_unicarbkb_glytoucan.csv
    filename: human_proteoform_glycosylation_sites_unicarbkb_glytoucan.csv

    mediatype: csv
    uri: https://data.glygen.org/.../reviewed/human_proteoform_glycosylation_sites_unicarbkb_glytoucan.stat.csv
    filename: human_proteoform_glycosylation_sites_unicarbkb_glytoucan.csv

    Content:
      Column Headers:
        uniprotkb_canonical_ac: Accession assigned to the protei
                             UniProtKB database
        glycosylation_site: Site on the protein sequence where g
        evidence: NCBI PubMed Id (PMID) as evidence for the entry
        unicarbkb_id: UnicarbKB database identifier
        glytoucan_ac: Unique accession assigned to the registered glycan structure in GlyTouCan database
        amino_acid: Three letter abbreviation code of the amino acid
        glycosylation_type: Type of glycosylation (N/O/C/S linked glycosylation)

    Statistics [Unique Value]:
                  uniprotkb_canonical_ac: 92
                  glycosylation_site: 223
                  evidence: 163
                  uckb_id: 984
                  glytoucan_acc: 824
                  amino_acid: 3
                  glycosylation_type: 3
```
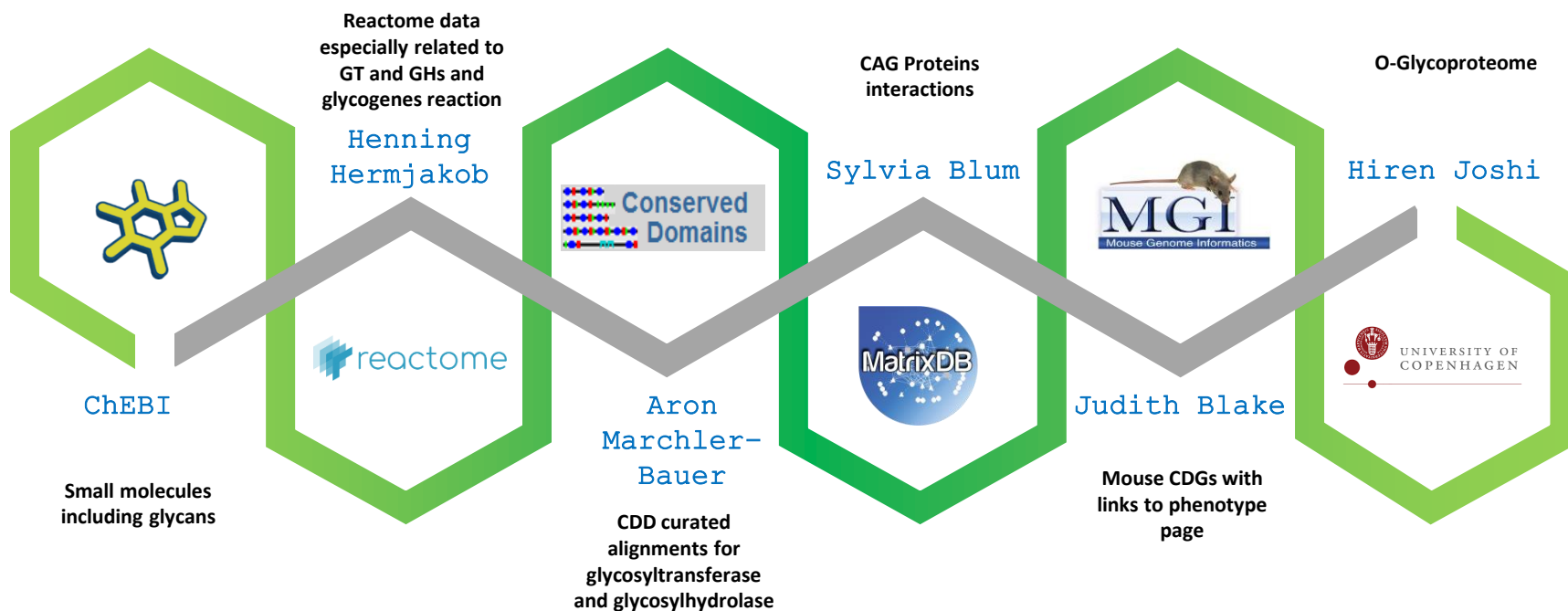
**Information on the output files Log files provide information on discarded/obsolete entries. The stat csv provides description of the column headers and the unique value statistics of the dataset.**
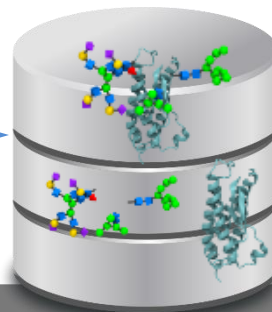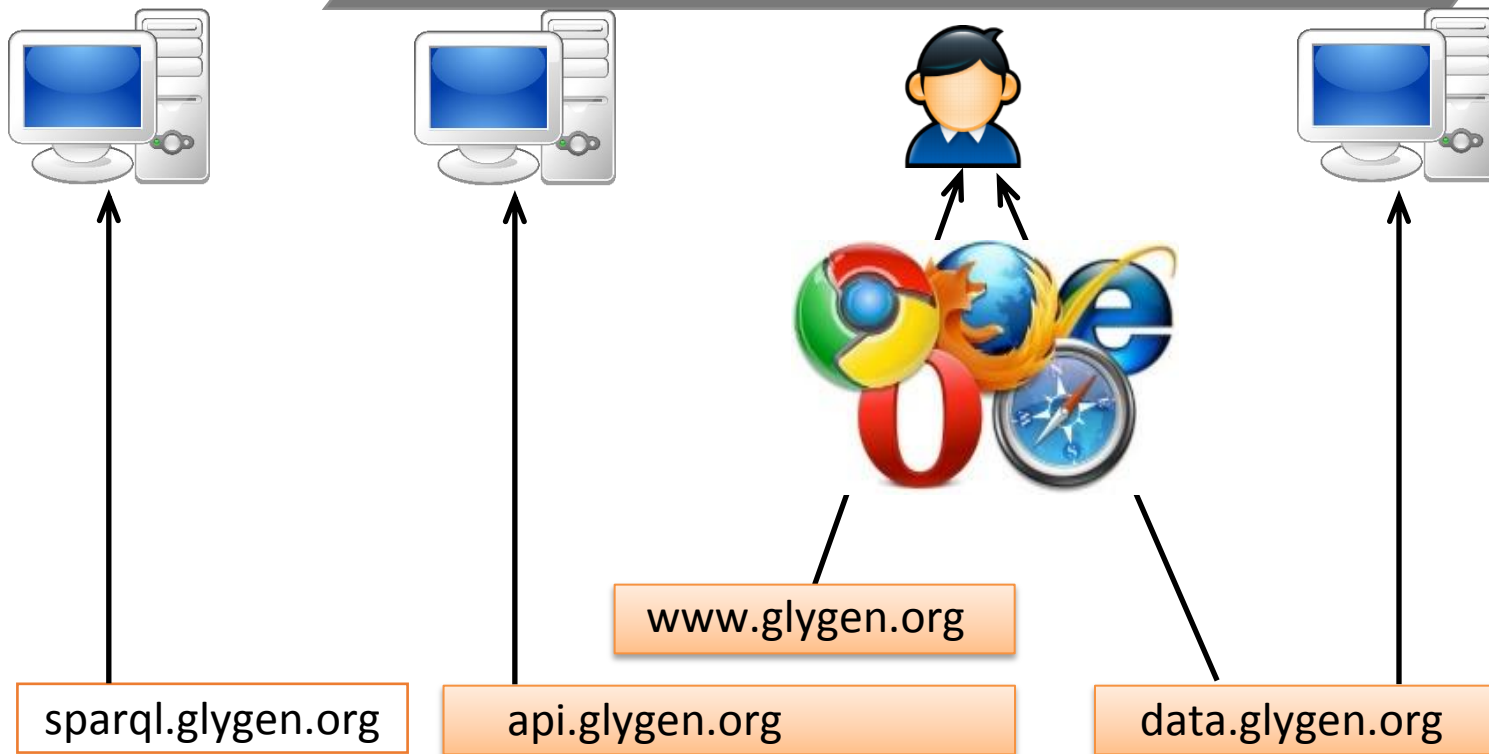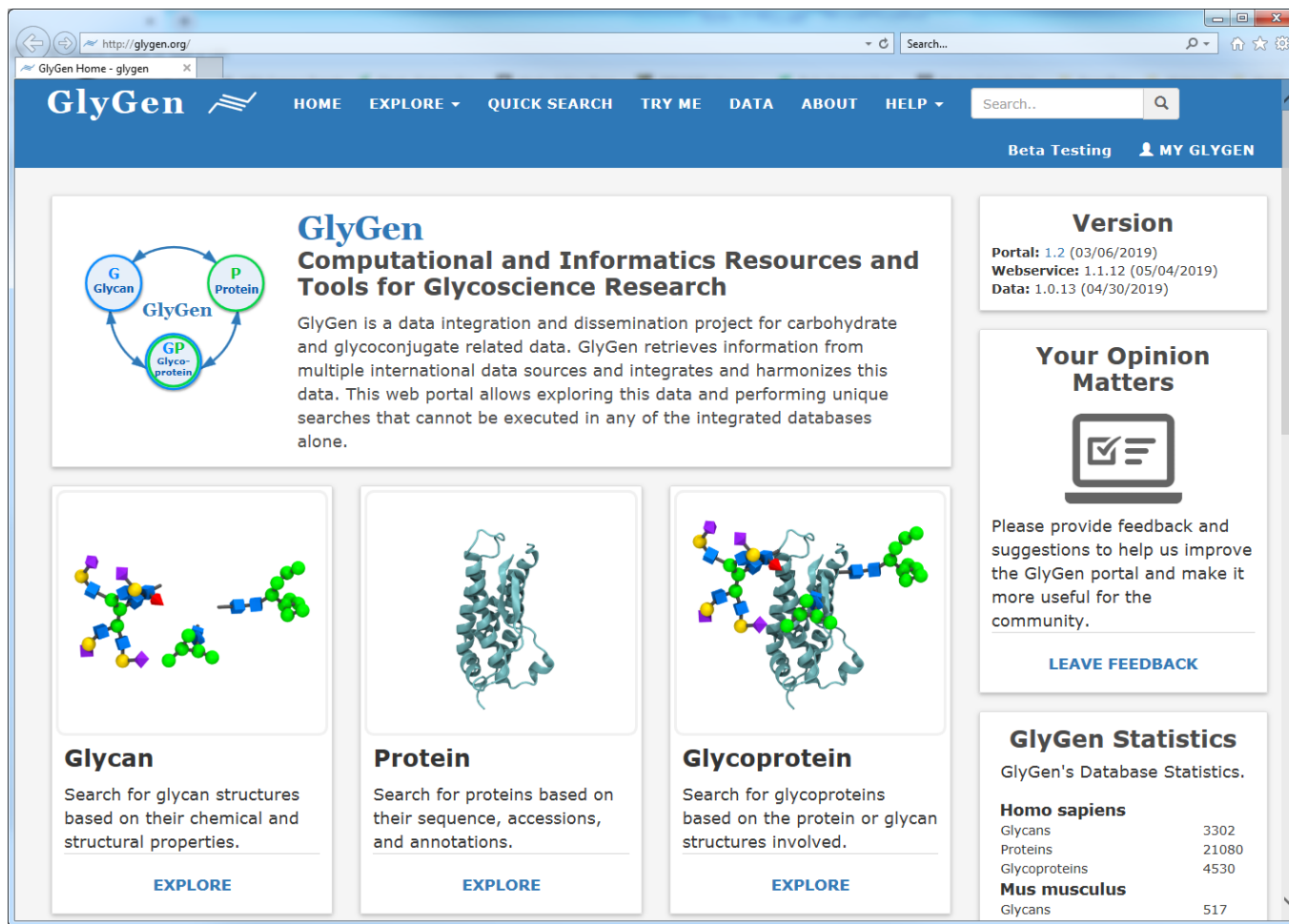
- Homo sapiens
  - Glycans          3302
  - Proteins         21080
  - Glycoproteins    4530
- Mus musculus
  - Glycans          517
  - Proteins         22289
  - Glycoproteins    3792

## https://data.glygen.org/

Reactome data especially related to GT and GHs and glycogenes reaction

**Henning Hermjakob**

CAG Proteins interactions

O-Glycoproteome

**Sylvia Blum**

**Hiren Joshi**

**ChEBI**

**Aron Marchler-Bauer**

**Judith Blake**

Small molecules including glycans

CDD curated alignments for glycosyltransferase and glycosylhydrolase

Mouse CDGs with links to phenotype page

- Data integration

- Data access

- Data sharing

sparql.glygen.org

api.glygen.org

www.glygen.org

data.glygen.org

UniProt
Pub Chem
MGI
hpo
DISEASE ONTOLOGY
Onco MX
PRO Protein Ontology
UniCarbKB
PDB PROTEIN DATA BANK
NCBI RefSeq
reactome
Data Generators
GlyTouCan
EMBL-EBI
oma BROWSER
Human Phenotype Ontology

# GlyGen ⫤ http://www.glygen.org

- GlyTouCan ID
- Mass
- Number of Monosaccharides
- Organism
- Glycan Type
- Protein Id / name
- Motif
- Enzyme

- Accession (UniProt, RefSeq)
- Mass
- Organism
- Protein name
- Gene name
- Attached glycan
- Glycosylated amino acid
- Sequence
- Pathway

General

Species

Function

**Glycosylation**

Sequence

Cross References

Pathway

Isoform

Orthologs

Disease

Mutation

Expression Tissue

Expression Disease

Publications



2keto

| UniCarbKB 1 | G77252PU | N-linked | Asn653 |

2keto

PubMed 2

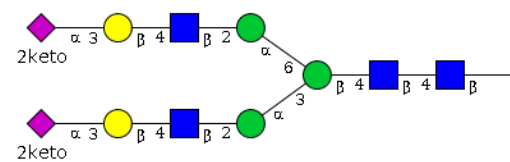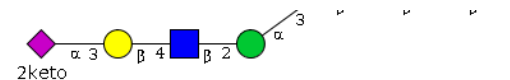2keto

Showing 1 to 9 of 9 rows

## Sequence

```
              +10        +20        +30        +40        +50
               |          |          |          |          |
    1 MWVTKLLPAL LLQHVLLHLL LLPIAIPYAE GQRKRRNTIH EFKKSAKTTL IKIDPALKIK
   61 TKKVNTADQC ANRCTRNKGL PFTCKAFVFD KARKQCLWFP FNSMSSGVKK EFGHEFDLYE
  121 NKDYIRNCII GKGRSYKGTV SITKSGIKCQ PWSSMIPHEH SFLPSSYRGK DLQENYCRNP
  181 RGEEGGPWCF TSNPEVRYEV CDIPQCSEVE CMTCNGESYR GLMDHTESGK ICQRWDHQTP
  241 HRHKFLPERY PDKGFDDNYC RNPDGQPRPW CYTLDPHTRW EYCAIKTCAD NTMNDTDVPL
  301 ETTECIQGQG EGYRGTVNTI WNGIPCQRWD SQYPHEHDMT PENFKCKDLR ENYCRNPDGS
  361 ESPWCFTTDP NIRVGYCSQI PNCDMSHGQD CYRGNGKNYM GNLSQTRSGL TCSMWDKNME
  421 DLHRHIFWEP DASKLNENYC RNPDDDAHGP WCYTGNPLIP WDYCPISRCE GDTTPTIVNL
  481 DHPVISCAKT KQLRVVNGIP TRTNIGWMVS LRYRNKHICG GSLIKESWVL TARQCFPSRD
```

☑ N-linked Glycosylation 4

☑ O-linked Glycosylation 1

☑ Mutation 2

## General

### General



- **GlyTouCan Accession:** G77252PU
- **Chemical Mass:** 2,222.78 Da
- **Glycan Type/Subtype:** N-Glycan complex

### Species

**Mus musculus:** GlycomeDB 1    UniCarbKB 1

## Species

## Motif

## Found Glycoproteins

## Cross References

## Biosynthetic Enzymes

## Digital Sequence

• Get all **enzymes** that my of been involved in *synthesis of the glycans* on <u>protein X</u>.

# GlyGen

## Releases

- Simple and Advance search options
- Example input for searches
- Download of data on all pages

- Feedback feature on all pages
- Orthologous for proteins
- Improved sequence display
- Improved navigation

**September 2018**

**March 2019**

**1.0**  **1.1**  **1.2**  **1.3**

**December 2018**

**June 2019**

- Human and Mouse data
- Basic search for glycan, proteins and glycoproteins
- Information pages for glycan, proteins and glycoproteins

- Improved Evidence display
- Motif pages
- New help feature
- Global search

- Data from new sources

- New species (Rat, HIV)

- Improved help system for users

- More use cases

- 3D structures for glycans, proteins and glycoproteins

- Better widgets to search for and display data

**GlyGen**



- Try the webpage – give feedback

- What are the questions you would like to ask

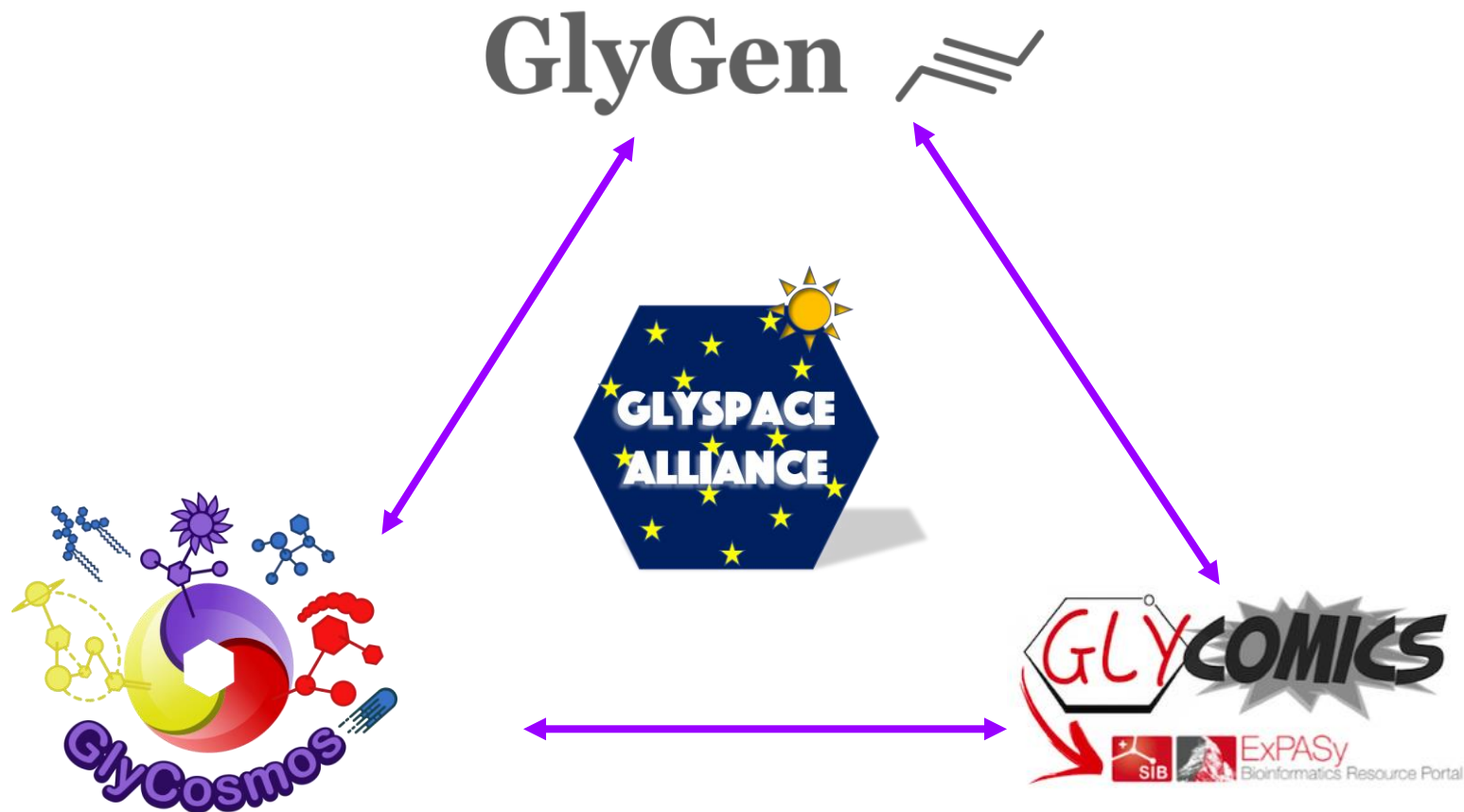- Enriching the database with own data

# Table of content

- Data integration

- Data access

- Data sharing

www.glygen.org
data.glygen.org
api.glygen.org

https://github.com/glygener

**Attribution 4.0 International**

Established August, 2018 @ Warren Workshop

http://glyspace.org

# FAIR data

**GlyGen**

- Data and material have sufficient rich metadata
- Use presistent identifier
- Data and metadata are indexed in a searchable resource

**Findable**

- Metadata and data can be accessed by humans and machines
- Use standarised communication protocols

**Accessible**

**Interoperable**

- Use open formats
- Consistent vocabulary
- Common metadata standards

**Reusable**

- Clear usage license
- Accurate information on provenance

https://www.go-fair.org/fair-principles/

**University of Georgia**
René Ranzinger
Michael Pierce
Robert Woods
Rupali Mahadik
Tatiana Williamson
Sena Arpinar
Sanath Bhatt
Sujeet Kulkarni
Sandeep Nakarakommula

**EMBL-EBI**
Maria Martin
Leyla Jael Garcia Castro
Preethi Vasudev

**NCBI**
Kim Pruitt
Evan Bolton

**The George Washington University**
Raja Mazumder
Robel Kahsay
Jeet Vora
Rahi Navelkar
Reza Mousavi
Nagarajan Pattabiraman
Xavier Holmes
Brian Fochtman

**Georgetown University**
Nathan Edwards
Radoslav Goldman
Darren Natale
Karen Ross
Wenjin Zhang

**Harvard University**
Richard Cummings

**The Jackson Laboratory**
Judith Blake

**Soka University**
Kiyoko Aoki-Kinoshita

**The Griffith University**
Matthew Campbell

**Imperial College London**
Ten Feizi

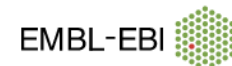**Macquarie University**
Nicki Packer

**NIH-NCI**
Jefferey Gildersleeve

**NIH Grant - U01 GM125267-01**

*Web portal:* *https://glygen.org*
*Data store:* *https://data.glygen.org*
*WS API:* *https://api.glygen.org*

**CONTACT**
**Will York** will@ccrc.uga.edu
**Raja Mazumder** mazumder@gwu.edu

# GlyGen.org