**Breakout Session 3: Track A**

**Systems Biology of Glycosylation: Extending Mechanistic Analysis Toward AI**

Dr. Sriram Neelamegham
*Professor/PI, University at Buffalo, State University of New York*

Dr. Rudiyanto Gunawan
*Associate Professor, State University of New York - Buffalo*

# *Systems Biology of Glycosylation: Extending Mechanistic Analysis Toward ML/DL*

**Sriram Neelamegham**

**Rudiyanto Gunawan**

**Changyou Chen**

Chemical & Biological Engineering, and Computer Science and Engineering

State University of New York, Buffalo

**University at Buffalo**
*The State University of New York*

**AI Supplement Program PI meeting**
**Breakout 3/Day 1**
**[2:45-3:45]**

# Summary of the project and project goals

**Overarching goal**: To conduct Systems Biology experimental and computational investigations in order to understand how gene expression and cellular epigenetics regulate glycan biosynthesis at the single cell level

**Aim 1. To enhance the depth of single-cell multi-omics studies in order to collect sufficient data for ML/DL.**

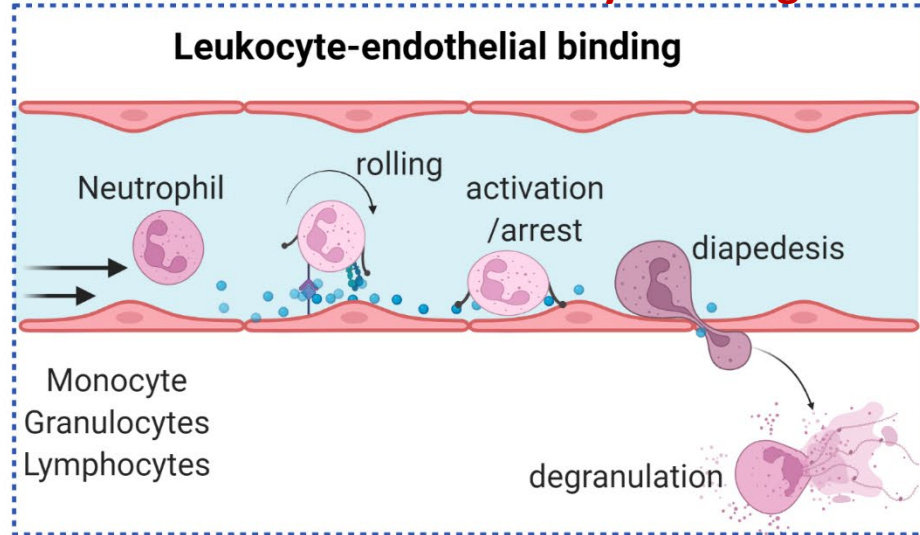**Aim 2. To process and normalize data for ML/DL applications.**

**Aim 3. To demonstrate the use of the transformed data in a DL/ML application.**
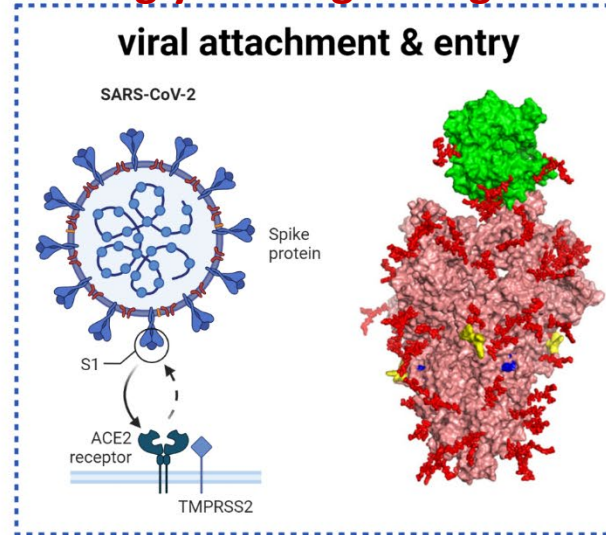
# Highlights: Research outputs and shared data

- Papers (other manuscripts are in preparation):
  a. P. Chrysinas, C. Chen, and R. Gunawan. CrossTx: Cross-cell-line transcriptomic signature predictions. *Processes*, 12:332, 2024.
  b. Cell and tissue-specific glycosylation pathways and transcriptional regulation informed by single-cell transcriptomics. *bioRxiv*, 559616, 2023.

- Website:
  a. glycoCARTA: Single-cell transcriptome of glycosylation. http://vgdev.cedar.buffalo.edu/glycocarta/
  b. glycoTF: Transcriptional factors of glycosylation. http://vgdev.cedar.buffalo.edu/glycotf/

# Ubiquitous in nature and relevant to biotechnology

**Role of selectins in leukocyte rolling**

**O- and N-glycans regulating SARS-CoV-2**



Leukocyte-endothelial binding

Neutrophil · rolling · activation/arrest · diapedesis

Monocyte Granulocytes Lymphocytes · degranulation



viral attachment & entry

SARS-CoV-2 · Spike protein · S1 · ACE2 receptor · TMPRSS2



host-microbiome

mucin covered gut



Glycans on rat capillary endothelial cells
(Essentials of Glycobiology; George Palade)

Vascular Lumen



B. subtilis bacterium, with hair-like glycocalyx (Wikipedia)



biotechnology

# Most common glycans are blood group antigen:
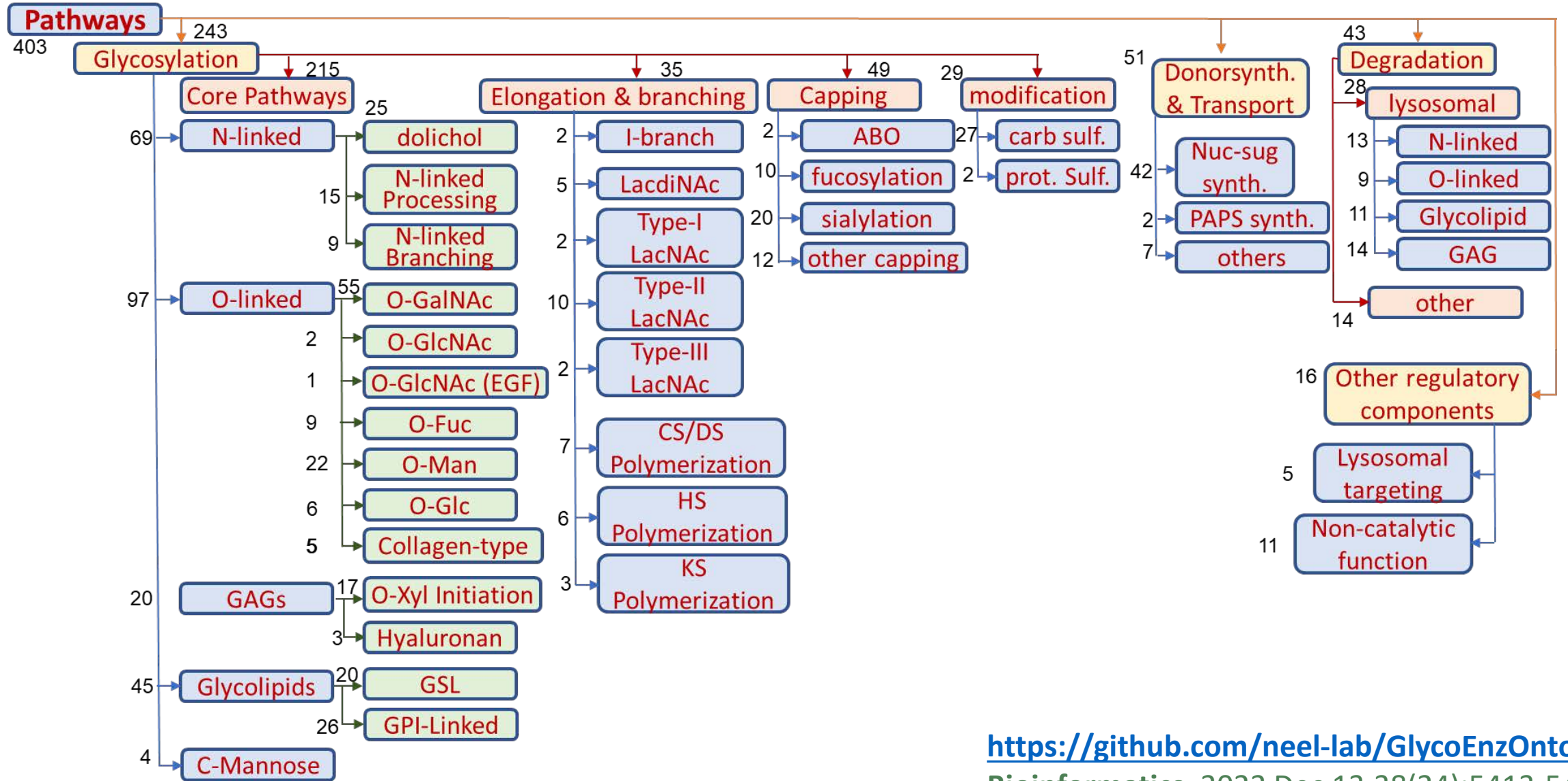# Expressed on RBCs and critical for transfusion medicine



**9 human monosaccharides**

| Monosaccharide | Symbol | Donor |
|---|---|---|
| Galactose (Gal) | ○ | → UDP-○ |
| N-Acetylgalactosamine (GalNAc) | □ | → UDP-□ |
| Glucose (Glc) | ● | → UDP-● |
| N-Acetylglucosamine (GlcNAc) | ■ | → UDP-■ |
| Mannose (Man) | ● | → GDP-● |
| Fucose (Fuc) | ▲ | → GDP-▲ |
| N-Acetylneuraminic acid (Neu5Ac) | ◆ | → CMP-◆ |
| Xylose (Xyl) | ☆ | → UDP-☆ |
| Glucuronic acid (GlcA) | ◇ | → UDP-◇ |

| Genotypes | O/O | A/O | A/A | B/O | B/B | A/B |
|---|---|---|---|---|---|---|
| Phenotypes | O | A | A | B | B | AB |
| Antigens | H | A | A | B | B | AB |
| Antibodies | anti-A anti-B | anti-B | anti-B | anti-A | anti-A | none |

# GlycoEnzOnto: An ontology for human glycosylating enzymes

# Multimodal measurements and their data integration



**Single-cell multiOMICs**

Fluorescent lectin → Spectral flow cytom.

transcript
Oligonucleotide-Tagged lectin
CRISPR/Cas9 Edit (sgRNA) → 10X plus NGS sequencing

Glycan structure → Mass spec. (Lumos Fusion)

BV605

Intensity

Channel

**Spectral flow cytometry**
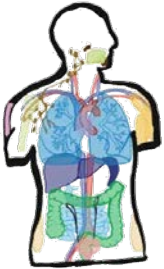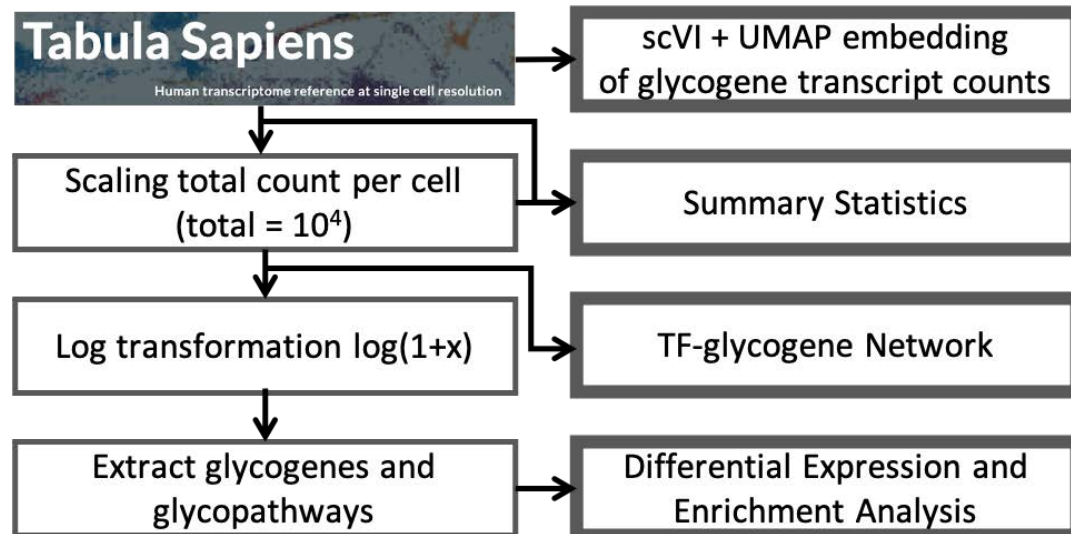
Glycopeptides

N-glycans

**Mass spectrometry**

# Single-cell Analysis of Glycosylation

**Tabula Sapiens Single-cell RNA-seq**
- ~500K cells, 400 cell types
- 24 organs,15 normal subjects
  - Sex: 9 male/6 female
  - Age: 22-74y
  - 6-White/6-Hispanic/2-Black/ 1-Asian

DecontX counts, 10X



- Establish baseline glycogene single-cell expressions in human
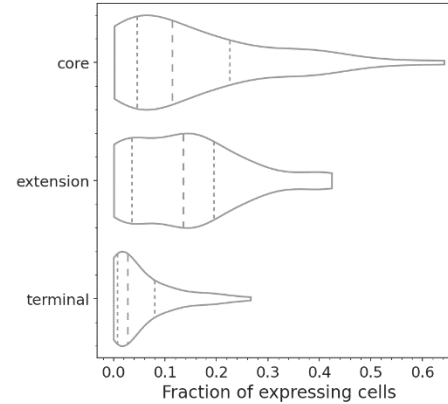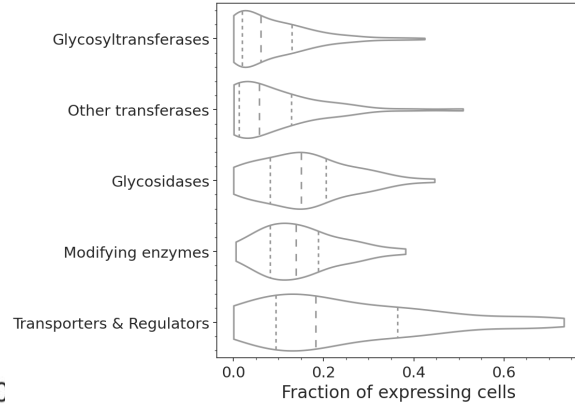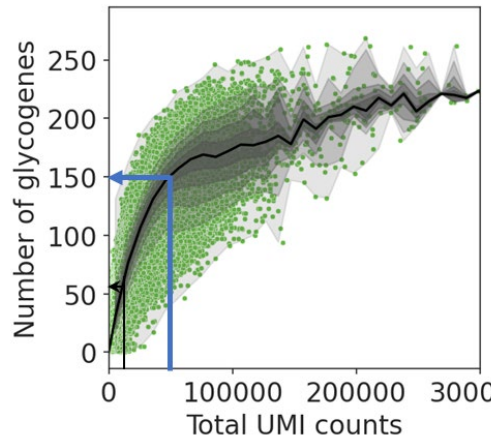- Establish data processing pipeline for ML/DL modeling

**How does the expression of glycogenes vary with cell/tissue types?**

**How prevalent are the glycogenes?**

**Computational prediction of TFs of Glycosylation?**

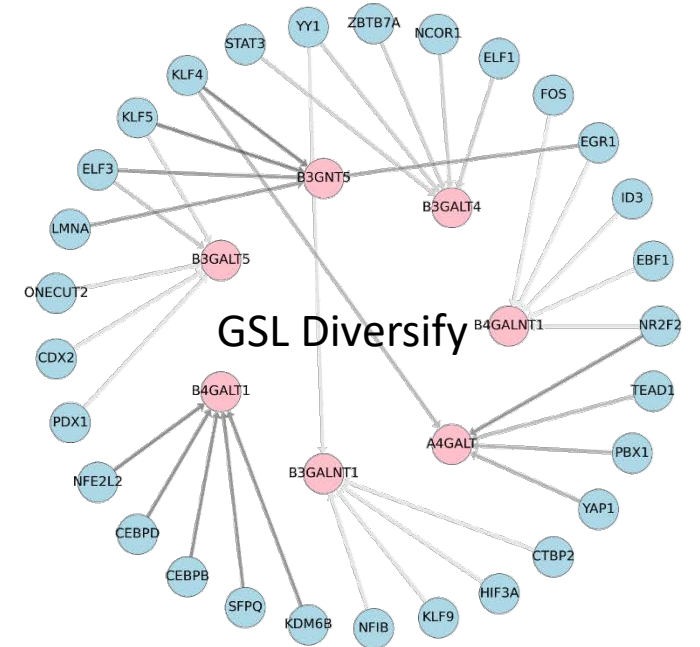**Glycosylation pathway variation across cell type and tissues ?**

https://tabula-sapiens-portal.ds.czbiohub.org/
*Science* 376(6594):eabl4896 2022.

# Single-cell Analysis of Glycosylation

GSL Diversify

glycoCARTA

- Glycogenes are as commonly expressed as other protein coding genes.

- At 50K-70K reads/cell, on average ~60 glyco-genes are detected (a max. of ~220 genes).

- Core pathways are expressed at higher levels than extension and terminal pathways.

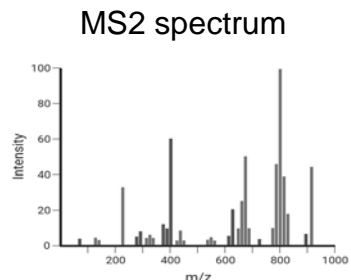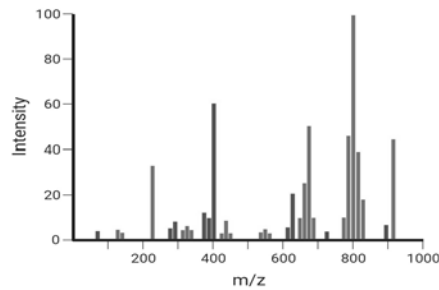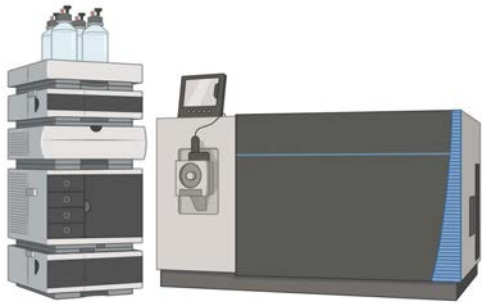# Large Language Model for Mass Spectrometry

# Large Language Model for Mass Spectrometry

| Level | Test Accuracy (%) | |
|---|---|---|
| | **GlycoBERT** | **CandyCrunch** |
| Mass | 99.75 | 98.49 |
| Composition | 99.57 | 97.7 |
| Topology | 96.73 | 89.8 |
| Structure | 95.33 | 87.18 |



- Trained on MS2 data from glycomics (~480K spectra)

- Transformer-based LLM is a powerful architecture for analyzing MS data of glycomics profiling.

- GlycoBART is capable of generating *de novo* glycan structure prediction.

- Metadata (glycan type, experimental parameters) are highly informative.

- A promising framework for building foundational models of mass spectra

# Challenges and future work

- Sparsity of glycosylation-specific data in literature:
  - Develop glycan specific tools, e.g. focused transcriptomics on glycogenes
  - Streamlined quantitative analysis of glycoproteins using MS

- Data: lack of labelled data and imbalanced dataset
  - Employ i*n silico* data, self-supervised model

- Generative AI (glycoBART) can hallucinate.
  - Incorporate postprocessing of predictions

- Translation to better patient healthcare and treatment
  - Incorporate constraints / structures informed by biological knowledge