

# ***Glycoinformatics tools to analyze and curate large scale experimental datasets***

***Sriram Neelamegham***

**Departments of Chemical & Biological Engineering,  
Biomedical Engineering and Medicine**

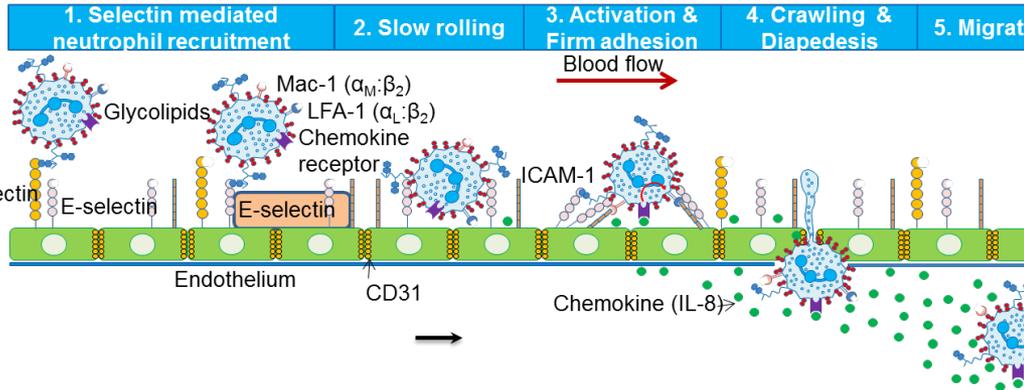
**State University of New York, Buffalo, NY**

9:10-9:50am March 6, 2018  
Tokyo, Japan

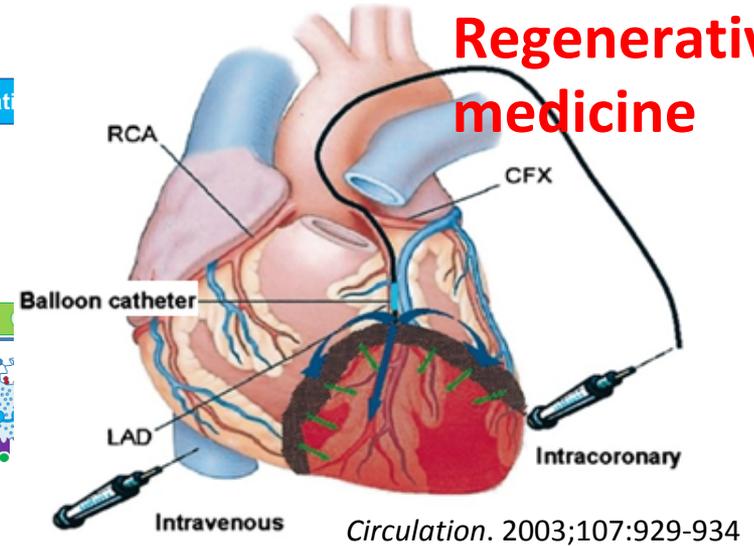


# Overview of research interests

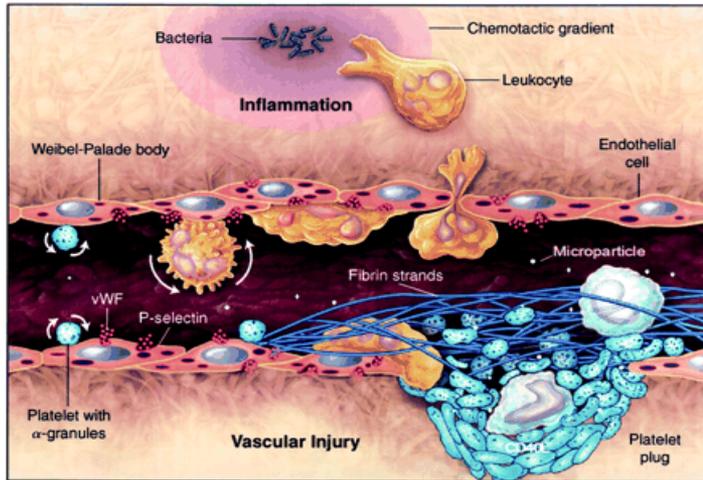
## Inflammation



## Regenerative medicine



## Thrombosis



*ATVB, 25:1321, 2005.*

**Systems Biology**  
Input-output  
relationship

# Input-output response

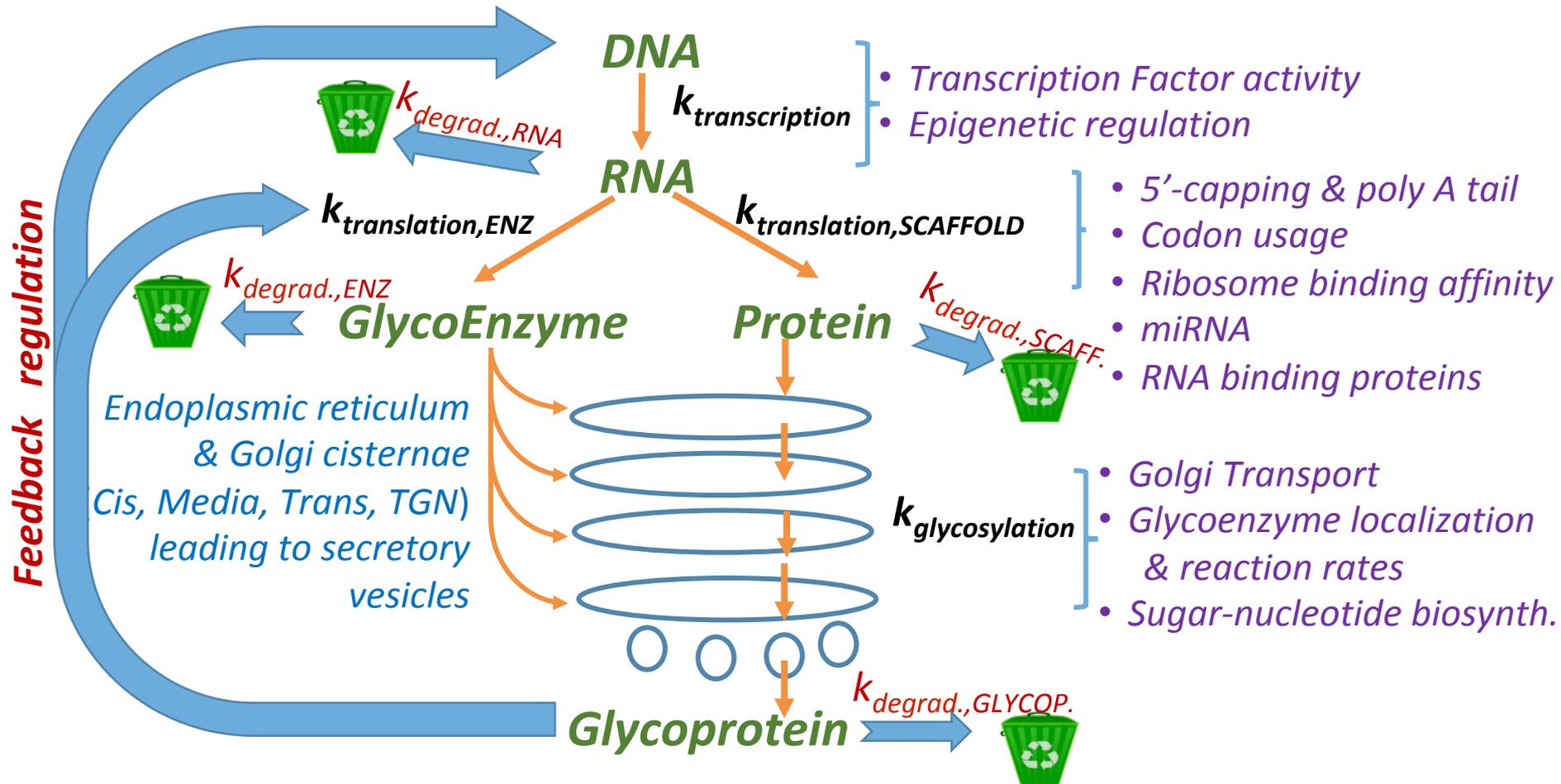
- **Their generation**

*Wet-lab: Next Generation Seq., LC-MS with CRISPR-Cas9 perturbations*

- **Their visualization, analysis and simulation**

*Dry lab: LC-MS data analysis programs,  
Pathway maps*

# VirtualGlycome.org: Systems level view of glycosylation



# Open-source integration of knowledge across scales

- **GNAT-Web**: Glycosylation network analysis toolbox
- **DrawGlycan-SNFG**: Simple tool to convert IUPAC strings to SNFG sketches
- **GlycoPAT**: High-throughput analysis of LC-MS<sup>n</sup> data, with focus on glycoProteomics

# 1. GNAT-Web

## Glycosylation Network Analysis Toolbox

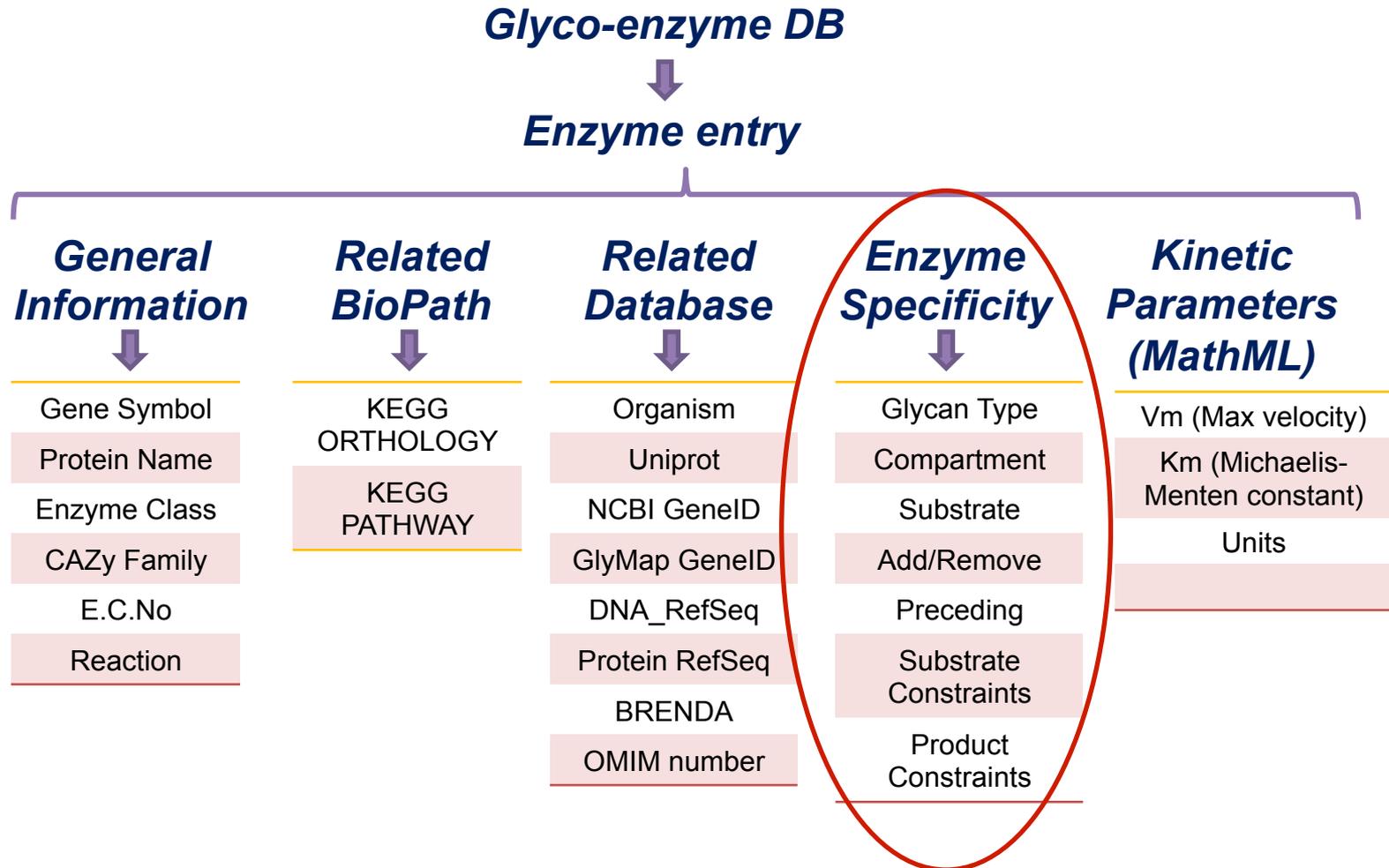


Yusen Zhou

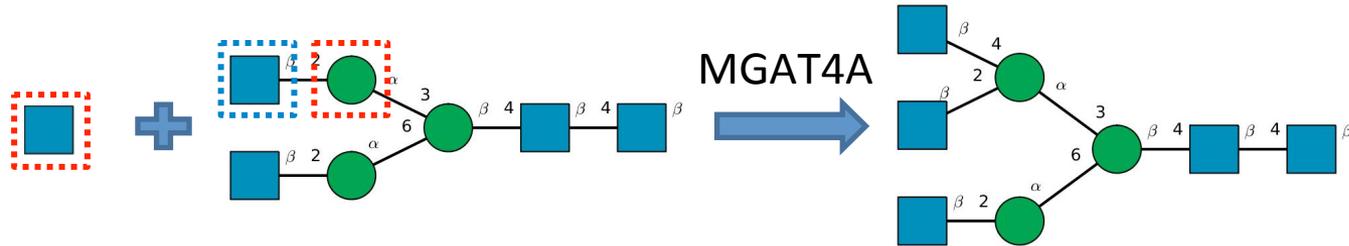
- Define glycoEnzymes *in silico*
- Develop reaction network from RNA-Seq and MS data processing
- Eventually, simulate reaction networks to bridge data across scales

Liu G, et al.  
*Bioinformatics*. 24(23):2740-7, 2008;  
*Glycobiology*. 21(12):1541-53, 2011;  
*Bioinformatics*. 29(3):404-6, 2013;  
*PLoS One*. 9(6):e100939, 2014.

# XML based glycoenzyme definition



# Enzyme specificity (e.g. MGAT 4)



**add** GlcNAc( $\beta$ 1-4) to **substrate** '^Man(a1-3)', provided it is **preceded** by GlcNAc( $\beta$ 1-2)

^: Caret is space for inclusion of arbitrary branches

**Constraint:** MGAT4 acts before addition of:

- Galactose (Gal), i.e. Gal cannot exist in string or **Gal#0**
- Bisecting MGAT3, i.e. **GlcNAc( $\beta$ 1-4)^Man( $\beta$ 1-4)#0**

Maximum # = 0 indicates **NOT**; But it could be **any other number as well**

Enzyme rule	Value
Add	GlcNAc( $\beta$ 1-4)
Substrate	^Man(a1-3)
Preceding	GlcNAc( $\beta$ 1-2)

References:

a. Bennun et al. *GalNAc6S Com* *Comput Biol.* 9(1):e1002813, 2013;

SubstConstraints	Value
MaxSubsubst	Gal#0&GlcNAc( $\beta$ 1-4)^Man( $\beta$ 1-4)#0

# Enzyme specificity (e.g. MGAT 4)

```

- <GeneralInfo>
  <GeneSymbol>MGAT4A</GeneSymbol>
  <ProteinName>Alpha-1,3-mannosyl-glycoprotein 4-beta-N-acetylglucosaminyltransferase A</ProteinName>
  <EnzymeClass>Glycosyltransferase</EnzymeClass>
  <GlycanType>N_linked</GlycanType>
  <CAZy>GT54</CAZy>
  <ECNo>2.4.1.145</ECNo>
  <TissuesName/>
  <OrgansName/>
</GeneralInfo>

```

```

  <Add>GlcNAc(b1-4)</Add>
  <Remove/>
  <Substrate>^Man(a1-3)</Substrate>
  <Preceding>GlcNAc(b1-2)</Preceding>

```

^: C  
of a

# = 0  
**NOT;**  
ld be

```

- <BioPath>
  <KO>K00738</KO>
  <Pathway>ko00510/ko00513</Pathway>
</BioPath>

```

```

  </SubstConstraints>
  - <ProdConstraints>
    <MaxProd/>
    <MinProd/>
    <MaxProdsubst/>
    <MinProdsubst/>
  </ProdConstraints>
</EnzSpecificity>

```

MaxSubsubst

```

- <EnzKinetics>
  <Vm>1e-5</Vm>
  <Km>3.4e9</Km>
  <Units>pM</Units>
  - <EnzDistribution>
    <ER/>
    <Cis>0.15</Cis>
    <Medial>0.45</Medial>
    <Trans>0.3</Trans>
    <TGN>0.1</TGN>
  </EnzDistribution>
</EnzKinetics>

```

3;

Gal#0&GlcNAc(b1-4)Man(b1-4)#0



# Custom database generation

Obtain from existing databases

Enzyme specificity & kinetics

## GNAT-WEB

Glycosylation Networks Analysis Toolbox

**Custom database generation** Database Visualization Create glyco-pathway

Database Edit: **Homo Sapiens** Export enzyme table Input Database Name... Load database Go

FUT1	FUT2	FUT3	FUT4	FUT5	FUT6
FUT7	FUT8	FUT9	FUT10	FUT11	ST3GAL2
ST3GAL1	ST3GAL3	ST3GAL4	ST3GAL5	ST3GAL6	ST6GAL1
ST6GAL2	ST6GALNAC1	ST6GALNAC2	ST6GALNAC3	ST6GALNAC4	ST6GALNAC5
ST6GALNAC6	ST8SIA1	ST8SIA2	ST8SIA6	ST8SIA3	ST8SIA4
ST8SIA5	B4GALT1	B4GALT2	B4GALT3	B4GALT4	B4GALT5
B4GALT6	B4GALT7	B3GALT1	B3GALT2	B3GALT4	B3GALT5
B3GALT6	UGT8	A4GALT	C1GALT1	B3GALNT1	B3GALNT2
B4GALNT1	B4GALNT2	B4GALNT3	B4GALNT4	GBGT1	ABO(A)

**General Info:**

Uniprot #:  Add

Tissues:

Organs:

Glycan class:  N-glycan  O-glycan  Glycosphingolipid

Compartment:  ER  Cis-Golgi  Medial-Golgi  Trans-Golgi  TGN (Trans Golgi Network)

**Enzyme Specificity:**

Add  or Remove

to/from

preceding

	Substrate Constraint	Product Constraint
Max Structure	<input type="text"/>	<input type="text"/>
Min Structure	<input type="text"/>	<input type="text"/>
Sub-struct, max#	<input type="text"/>	<input type="text"/>
Sub-struct, min#	<input type="text"/>	<input type="text"/>

**Simulation Parameters:**

Michaelis-Menten parameters

Vm  pM/cell/hr Km  pM

Relative concentration of enzymes in compartment (sum will be normalized to 1).

ER	Cis-	Medial-	Trans-	TGN
<input type="text"/>				

**Comments:**

# View database elements



GNAT-WEB  
Glycosylation Networks Analysis Toolbox

Custom database generation **View database elements** Create glyco-pathway

Customized Database: **GeneDB**

## Glyco-Enzyme List

Search... **GeneSymbol**

Gene Symbol	Pathway
Enzyme Family MGAT1	all
	<input type="checkbox"/> N_linked
	Dolichol Pathway
	Branching Pathway
	Complex Pathway
	Type-1/2 LacNAc
	Blood Group i&l
	ABO Blood Group
	Sd <sup>a</sup> Antigen & GM2
	$\alpha$ 2-3 and $\alpha$ 2-6 NeuAc
	$\alpha$ 2-8 on N_glycan
	Gal $\alpha$ 1-3Gal Antigen
	<input type="checkbox"/> O_linked
	Core1/2 Pathway
	Core3/4 Pathway

### Pathway Map

## MGAT1

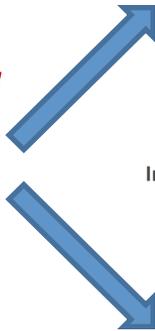
### General Information

Gene Symbol: **MGAT1**  
Protein name: **Alpha-1,3-mannosyl-glycoprotein 2-beta-N-acetylglucosaminyltransferase**  
Enzyme Class: **Glycosyltransferase**  
CAZy Family: **GT13**  
E.C. No: **2.4.1.101**  
Tissues:  
Organs:

Legend: n, +, .., ..

# Create pathways : forward

**Specify starting material and enzymes**



View examples

**Enzyme**

Enzyme List:  N\_linked

- MAN1A1
- MAN1B1
- MAN2A1
- MGAT1
- MGAT2
- MGAT3
- MGAT4A
- MGAT5
- B4GALT1
- B3GNT2
- FUT8
- ST6GAL1
- B4GALT1
- B4GALT6
- B4GALNT3
- MGAT2
- MGAT4C
- FUT10
- ST6GAL2
- B4GALT2
- B3GALT1
- B4GALNT4
- MGAT3
- B3GNT2

Input Glycans:

Compartments:

Group in bracket:

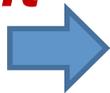
Choose mechanism:

- Termination Step
- Maximum M:
- # of sub-struct:
- Final Gly:

Simulation Parameters

**Run**

**Constraints to limit network size**



# Create pathways :reverse



View examples

Enzyme List:

N\_linked ▾

- FUT1
- ST3GAL3
- ST8SIA2
- B4GALT3
- B3GALT2
- ABO(A)
- MGAT4A

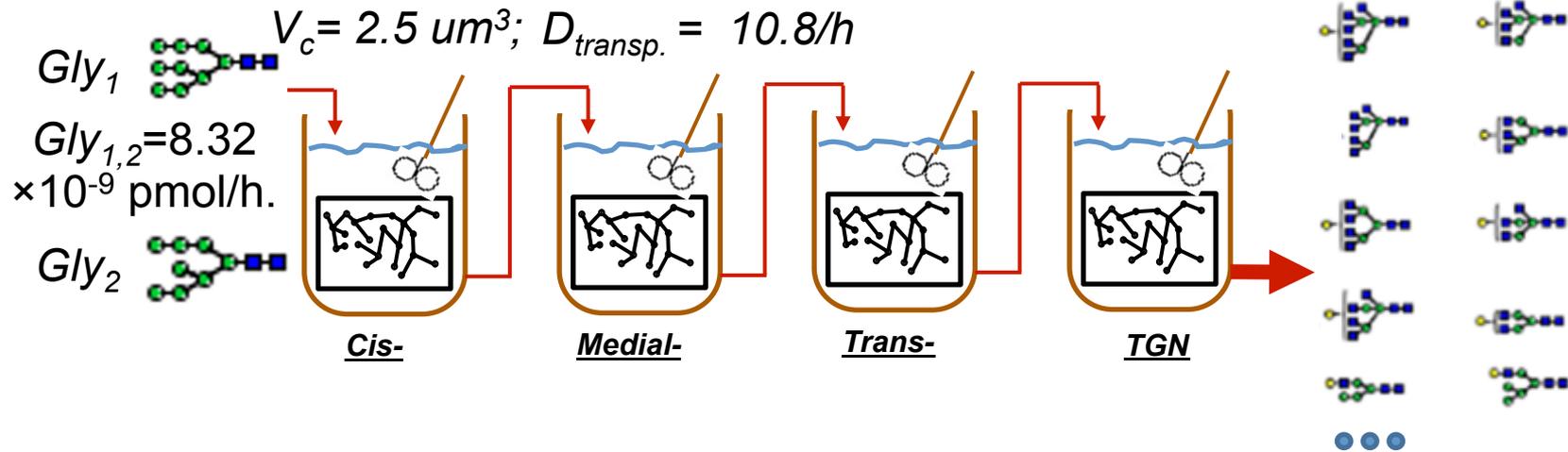
**Enzyme**  
MAN1A1 MAN1B1  
MAN2A1 MGAT1  
MGAT2 MGAT3  
MGAT4A MGAT5  
B4GALT1

- FUT10
- ST6GAL2
- B4GALT2
- B3GALT1
- B4GALNT4
- MGAT3
- B3GNT2

Specify ending  
glycans &  
enzymes



# Pathway: 4 compartment CSTR



## Species balance equation:

$$\frac{d[Gly_{i,j}]}{dt} = D_{\text{transp.}} \times [Gly_{i,j-1}] - V_{i,j} \times [Gly_{i,j}] / V_c \times K_{m_{i,j}} \times (1 + \sum k_{i,j} \times [Gly_{k,j}] / (K_{m_{k,j}})) - D_{\text{transp.}} \times [Gly_{i,j}]$$

0  
(CSTR)

↑  
in

↑  
Reaction

↑  
out

$i=36$  glycans  
 $j=4$  compartment

# *In silico simulation*

- **Deterministic**
- **Stochastic**

	High mannose	Bi-	Tri-	Tetra-	Bisecting
Cis-	11.62%	54.64%	4.61%	0.06%	48.38%
Medial-	2.73%	80.26%	8.85%	0.07%	70.33%
Trans-	2.44%	82.61%	9.12%	0.07%	71.11%
TGN	2.43%	82.73%	9.13%	0.07%	71.11%



# Pathway generation times: **short!**

Forward inference		
# of species	# of reactions	Time
144	300	~8s
405	692	~30s
916	3330	~70s
166	452	~11s
Reverse inference		
160	526	~12s
181	611	~14s
356	1246	~37s
96	267	~6s

# What else could this be useful for?

- Mapping RNA-Seq and Glycomics data to construct pathway maps.
- GlycoMir: The glycogene microRNA targets
- *Functional integration with other databases*

hsa-miR-342-3p  
hsa-miR-186-5p  
hsa-miR-5192  
hsa-miR-1256  
hsa-miR-4281  
hsa-miR-203b-3p  
hsa-miR-7843-3p  
hsa-miR-186-5p  
hsa-miR-583  
hsa-miR-4694-3p  
hsa-miR-3184-3p  
hsa-miR-383-5p.2  
hsa-miR-139-5p  
hsa-miR-4330  
hsa-miR-6758-3p  
hsa-miR-2116-3p  
hsa-miR-6837-3p  
hsa-miR-4493  
hsa-miR-3150b-3p

O-linked  
GSL

Selection parameters  
Transcript: GGGCCTTGTT CCGCTGCTCG ACACCACGGG GATGAACCCA GCCAGGI

Gene Symbol: **FUT8**  
miRNA Name: **hsa-miR-449b-5p**

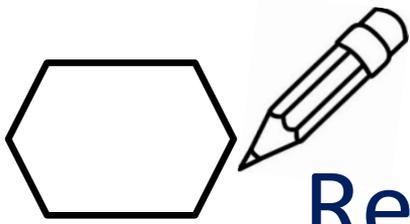
\* Collaboration with Lara Mahal (NYU)



*Kai Cheng*

## 2. DrawGlycan-SNFG

- From IUPAC to Symbolic Nomenclature for Glycans (SNFG)
- Draw glycopeptides
- Draw glycan and peptide fragmentation
- *Other features...*

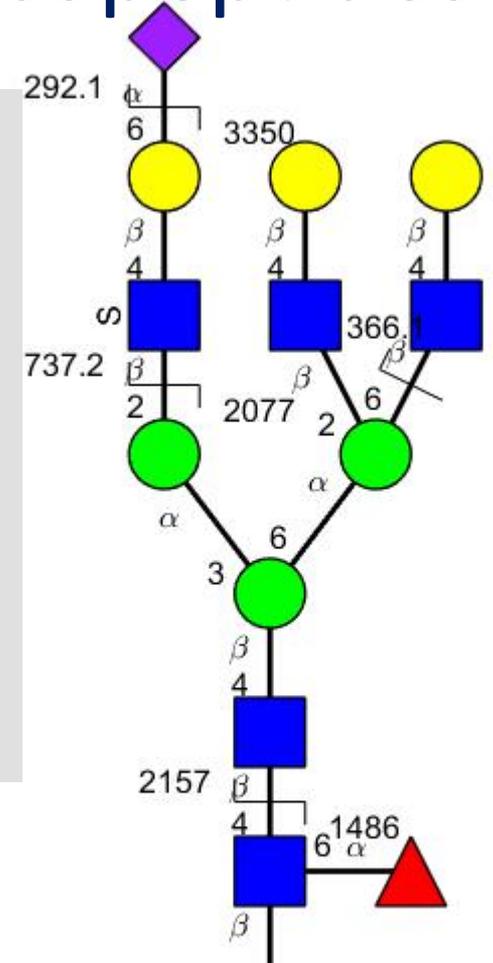


# DrawGlycan-SNFG: Render glycans and glycopeptides

```

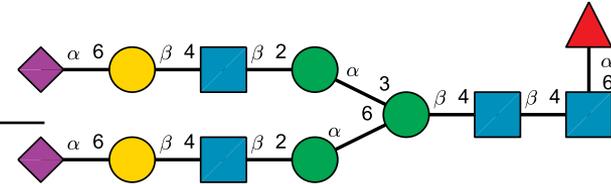
IGEADFN[Gal(b1-4)GlcNAc(b1-6 -NR "366.1")
[Gal(b1-4)GlcNAc(b1-2)]Man(a1-6)
[Neu5Ac(a2-6 -NR "292.1" -R
"3350")Gal(b1-4)GlcNAc(b1-2 -NR "737.2" -R
"2077" -U
"S")Man(a1-3)]Man(b1-4)GlcNAc(b1-4 -NR
[Gal(b1-4)GlcNAc(b1-2)]Fuc(a1-6)GlcNAc(b1-4 -NR]RSK
"2157" -R "1486")Fuc(a1-6)GlcNAc(b?-?)RSK
  
```

**Advantages:** straightforward, easy to read & write, adequate for common use

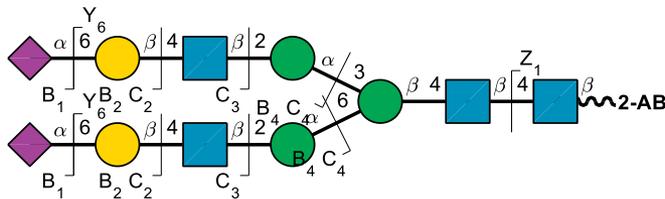
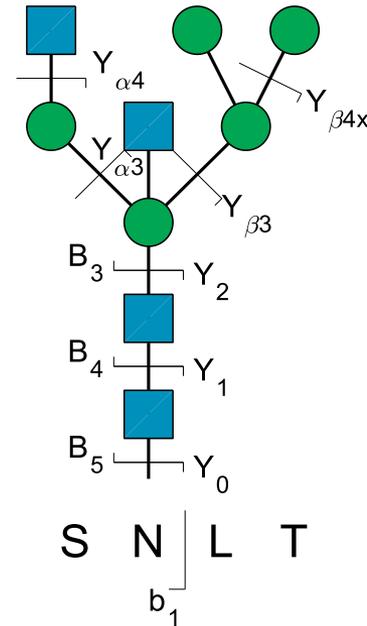


IGEADFNRSK  
b2 171.1

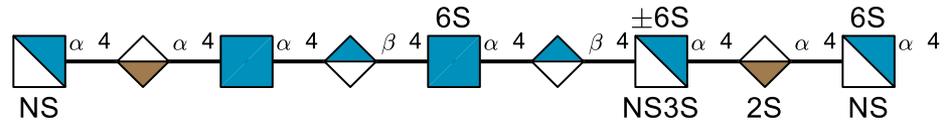
# Fragmentation options



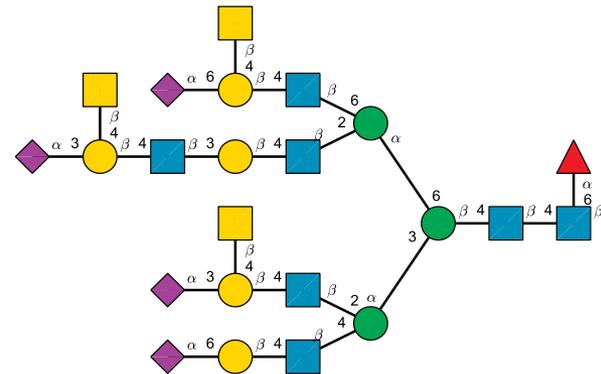
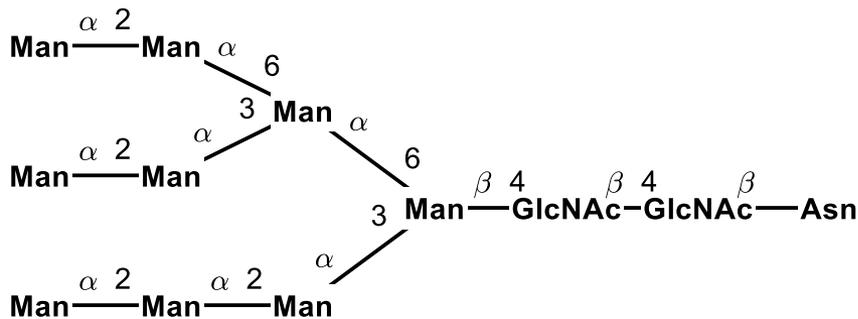
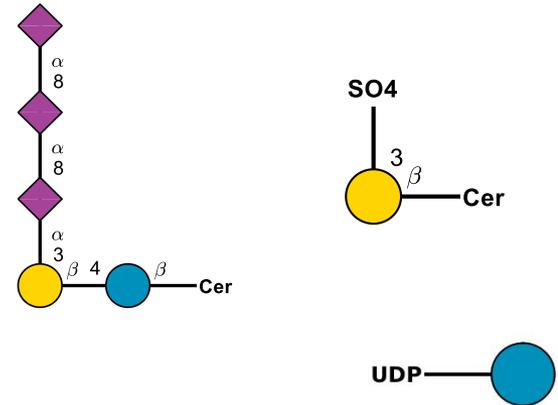
-Option	Representation
1 -R	Glycan reducing end
2 -NR	Glycan non-reducing end
3 -N	Peptide backbone N-terminus
4 -C	Peptide backbone C-terminus



# Monosaccharide options

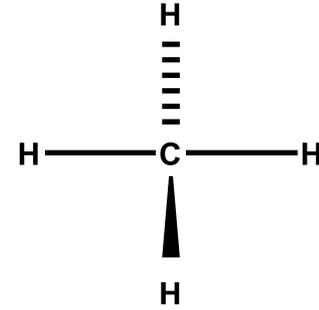
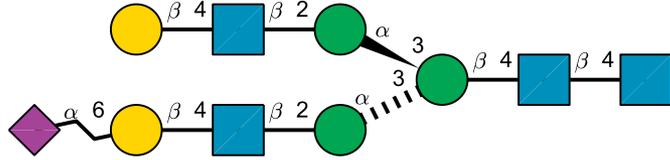
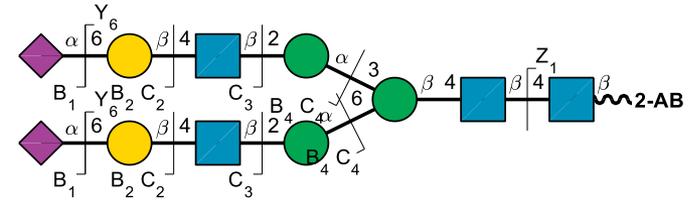
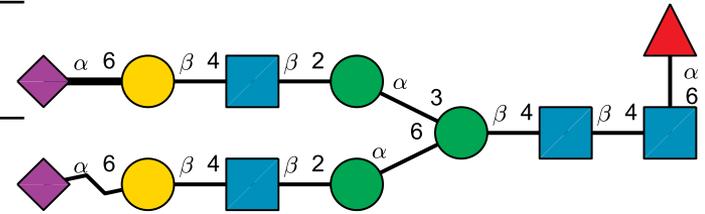


-Option	Representation
1 -U	Annotate above monosac.
2 -D	Annotate below monosac.
3 -P	Identify a perpendicular monosac.
4 -CHAR	Introduce arbitrary text or present monosac. in text form

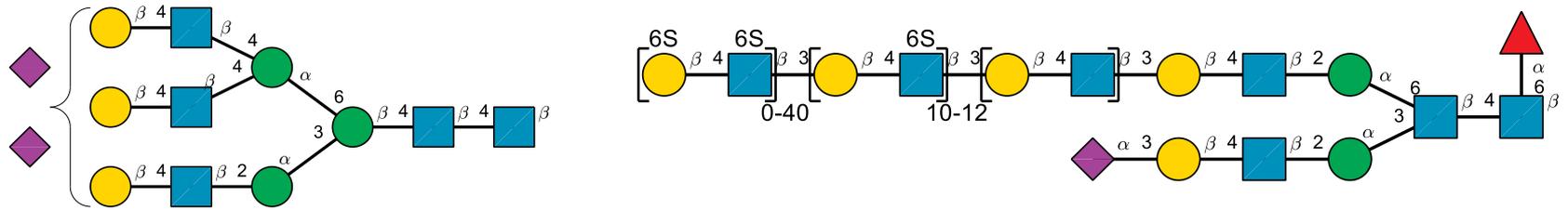


# Bond options

	-Option	Representation
1	-BOLD	Paint glycosidic bond bold
2	-ZIG	Paint glycosidic bond zigzag
3	-WAVY	Paint glycosidic bond wavy
4	-DASH	Paint glycosidic bond dashed
5	-WEDGE	Paint glycosidic bond wedge



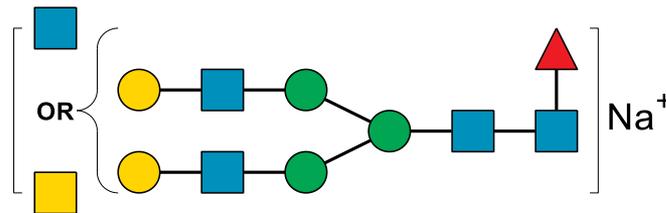
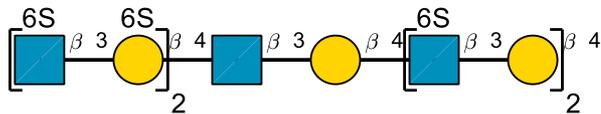
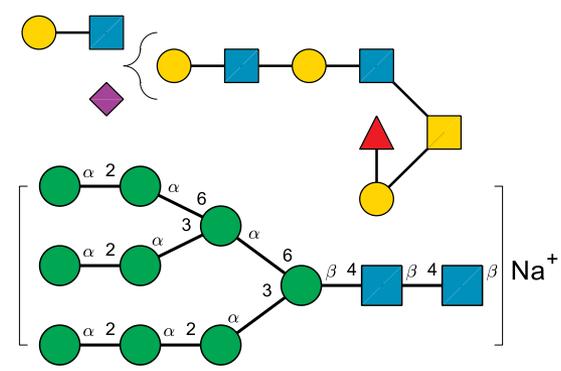
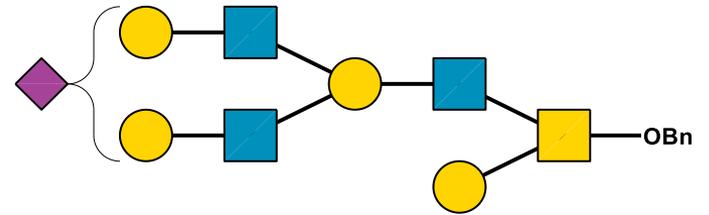
# Repeats, adducts and fuzzy options



**-Option**

**Representation**

- |   |         |  |
|---|---------|--|
| 1 | -RS     | Repeating unit start                   |
| 2 | -RE     | Repeating unit ends                    |
| 3 | -ADDUCT | Add glycan adduct                      |
| 4 | -CURLY  | Ambiguous assignments/fuzzy structures |



# DrawGlycan: Web, GUI & Command-line version

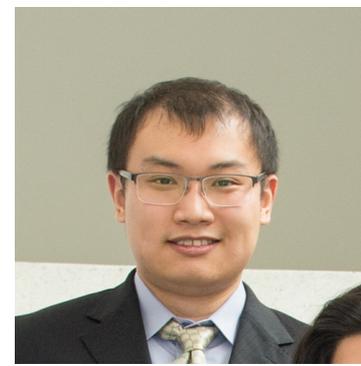


The image shows a screenshot of the DrawGlycan-SNEFG web application interface. The top navigation bar includes the application logo (a hexagon with a pencil), the title "DrawGlycan-SNEFG", and links for "Virtual Glycome", "Release notes", and "FAQ". Below the navigation bar, there is a window titled "drawglycangui" with standard window controls. In the foreground, a "Command Window" is open, displaying the following command:

```
fx >> drawglycan('Gal()GlcNAc()[Gal()GlcNAc()]Man()[Neu5Ac()Gal()GlcNAc()Man()])Man()GlcNAc()[Fuc()]GlcNAc()','perpendicularmonosac',{})
```

[VirtualGlycome.org/DrawGlycan](http://VirtualGlycome.org/DrawGlycan)

# 3. GlycoPAT

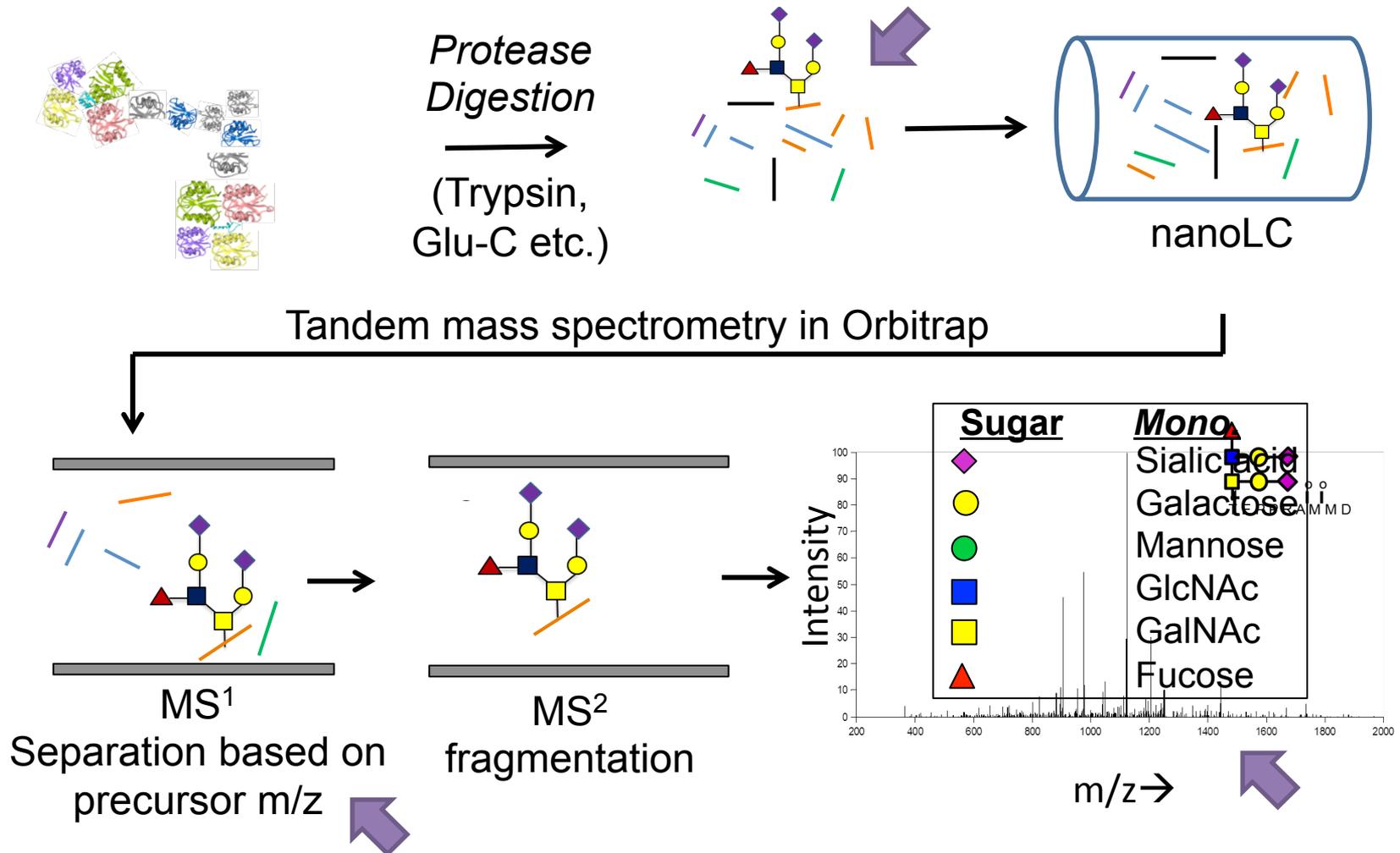


*Kai Cheng*

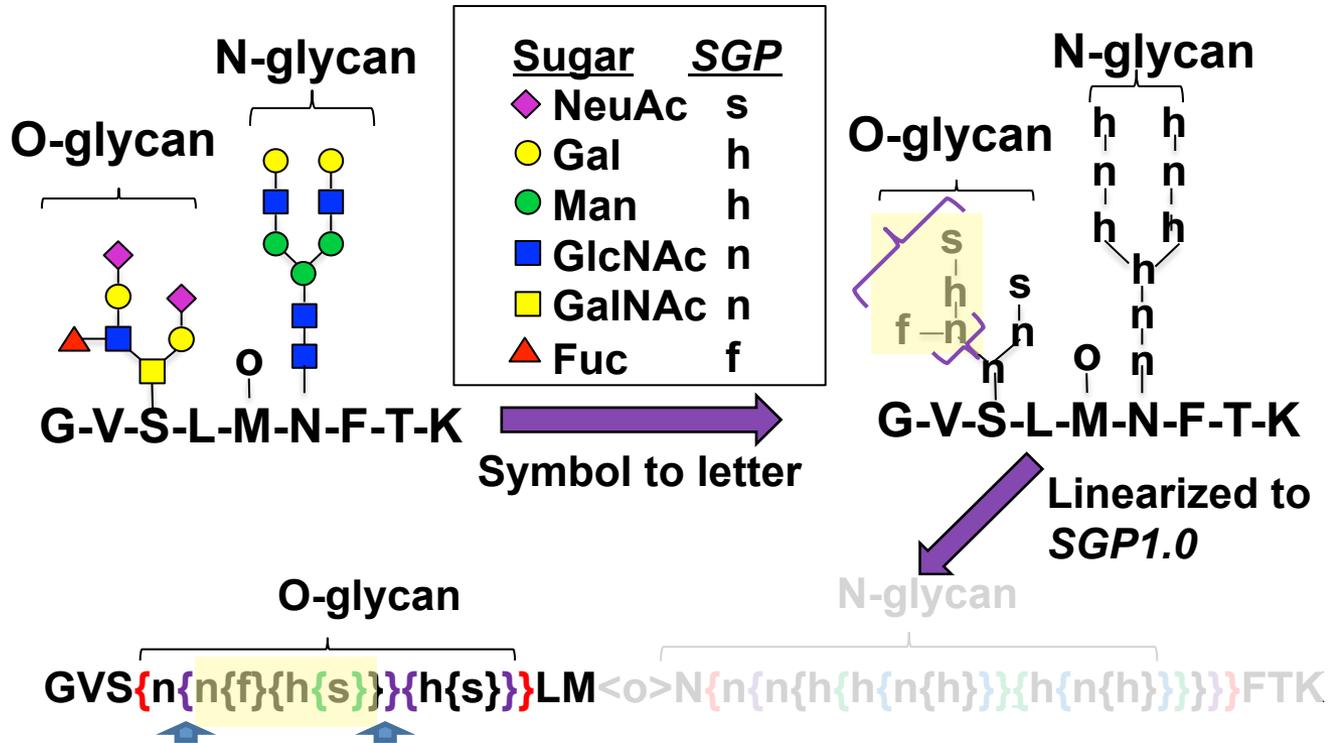
- Analyze high-throughput glycoproteomics data
- Comprehensive scoring and false discovery rate (FDR) calculation algorithm for MS<sup>n</sup> data analysis

*Mol Cell Proteomics.16:2032-47, Nov 2017.*

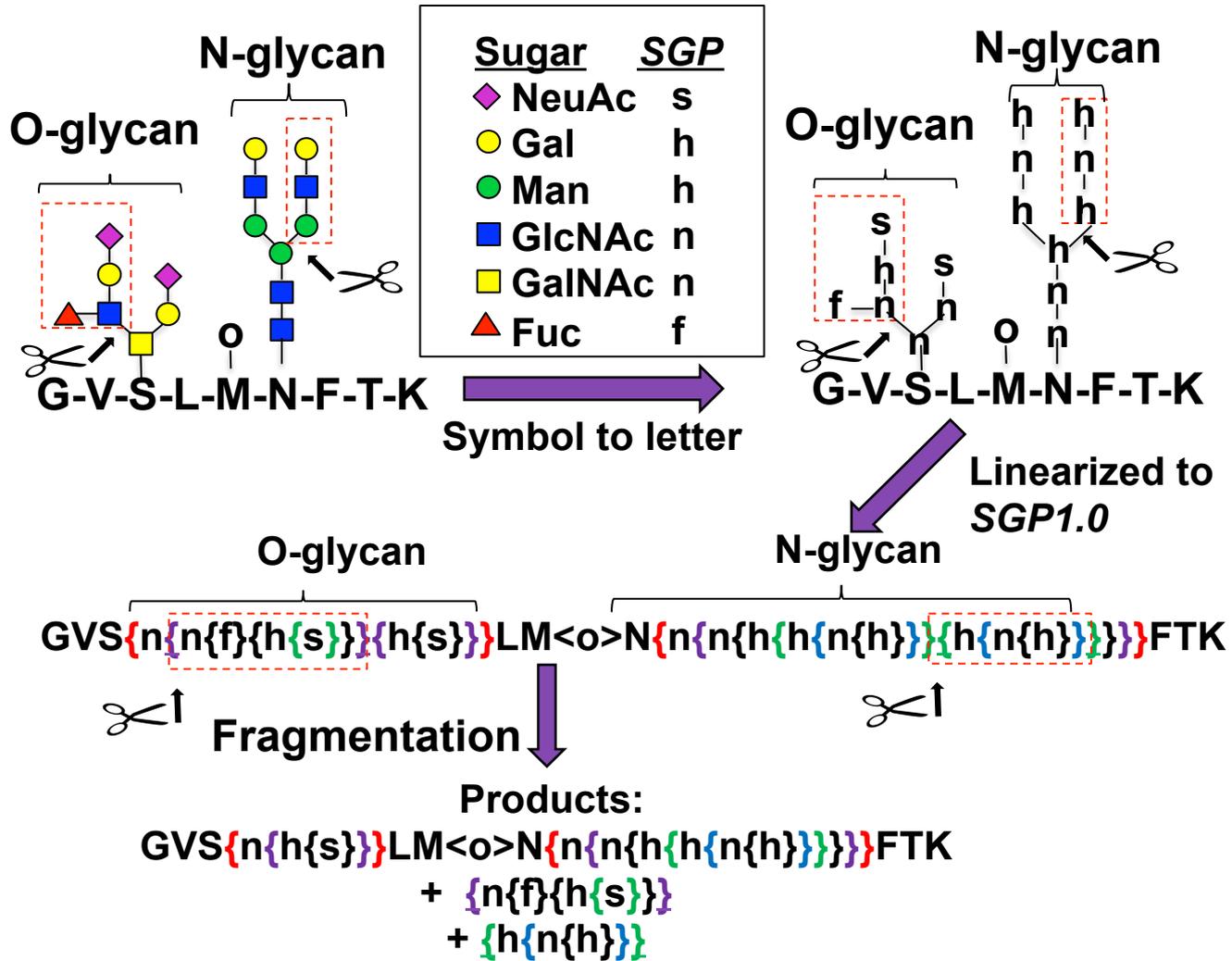
# GlycoPAT: High-throughput glycoproteomics analysis



# SmallGlyPep : The *minimal* representation of glycopeptide for MS

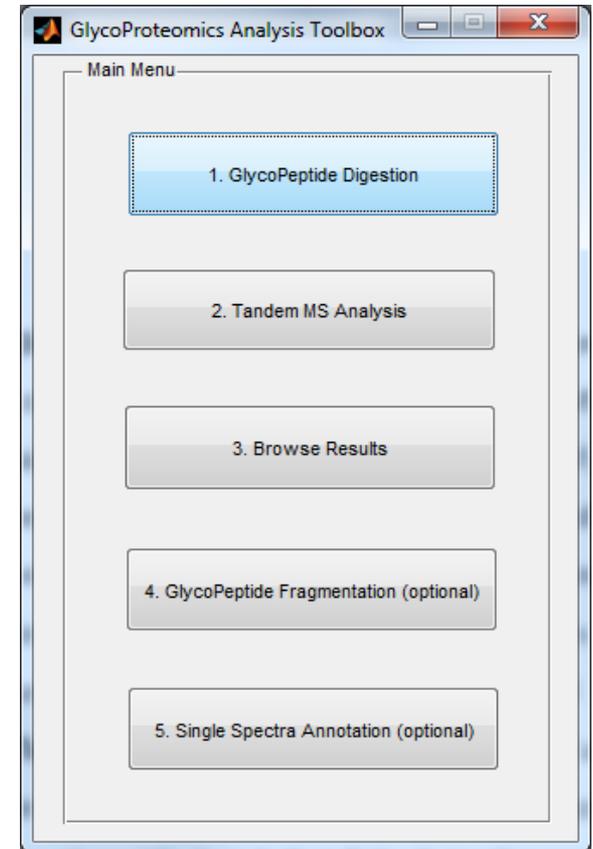
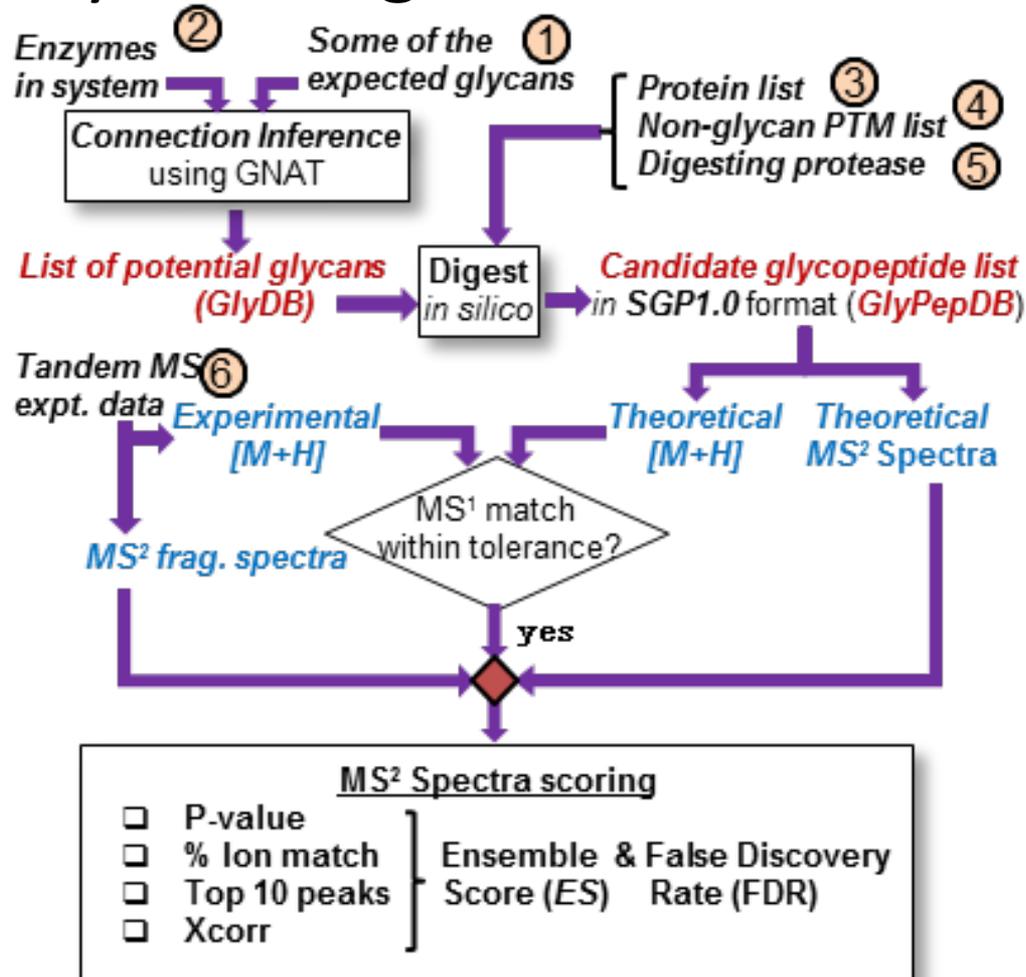


# SmallGlyPep : The *minimal* representation of glycopeptide for MS



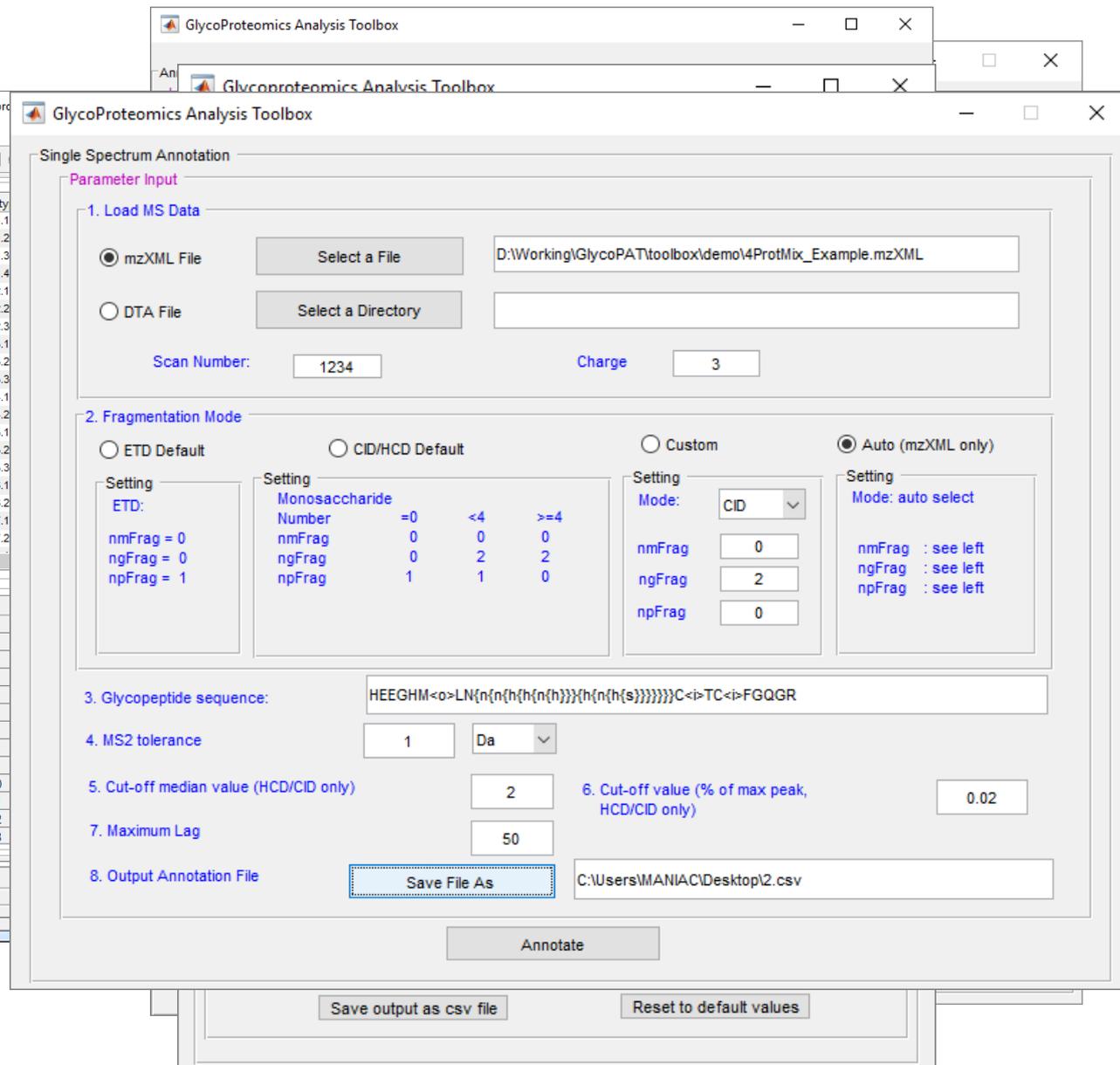
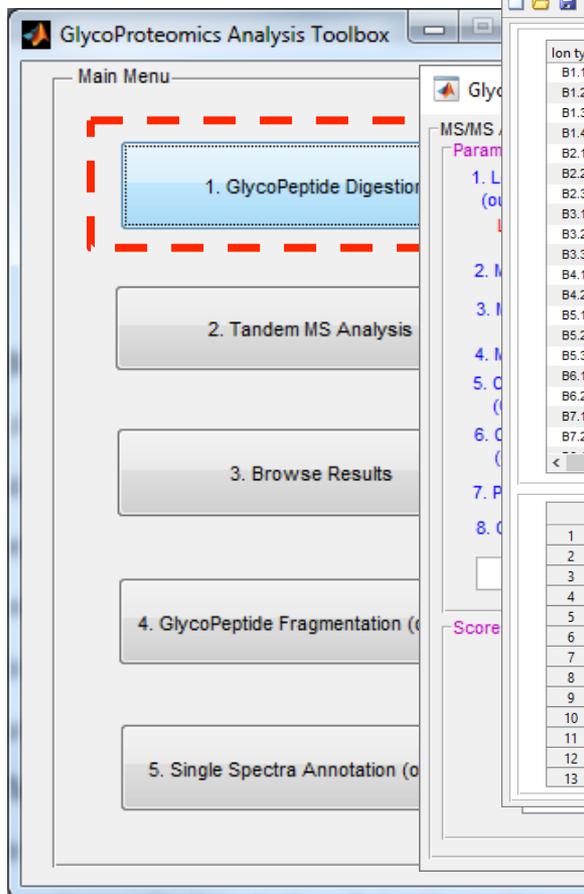
# Methods

## GlycoPAT: general workflow

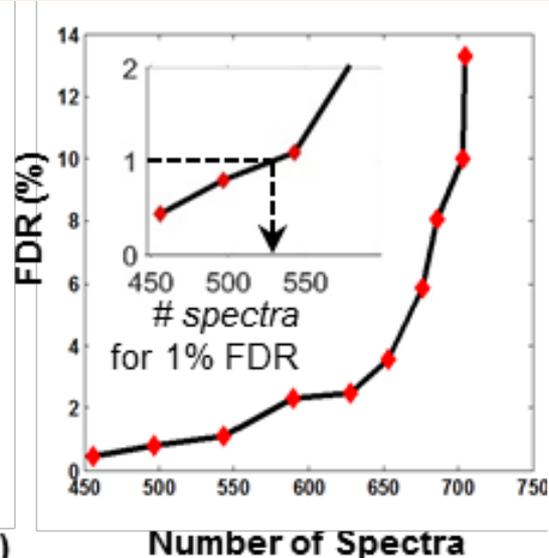
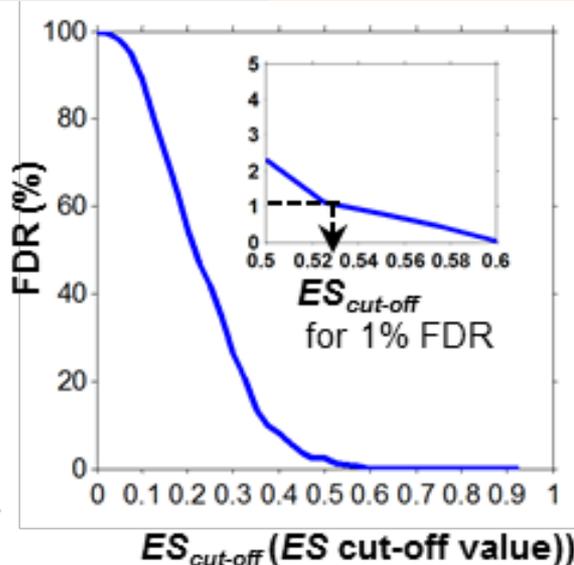
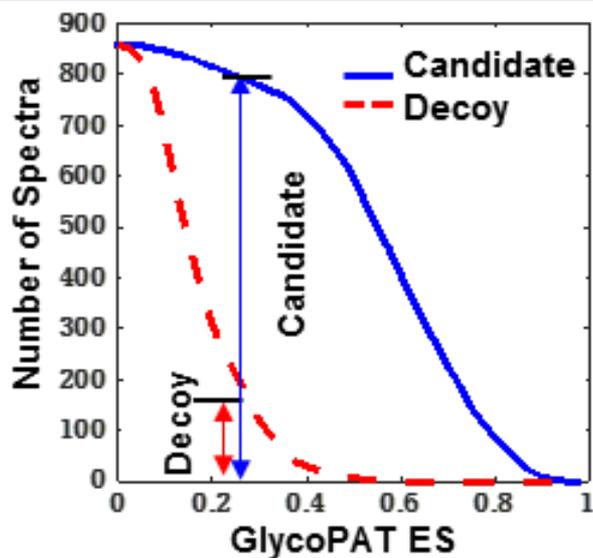
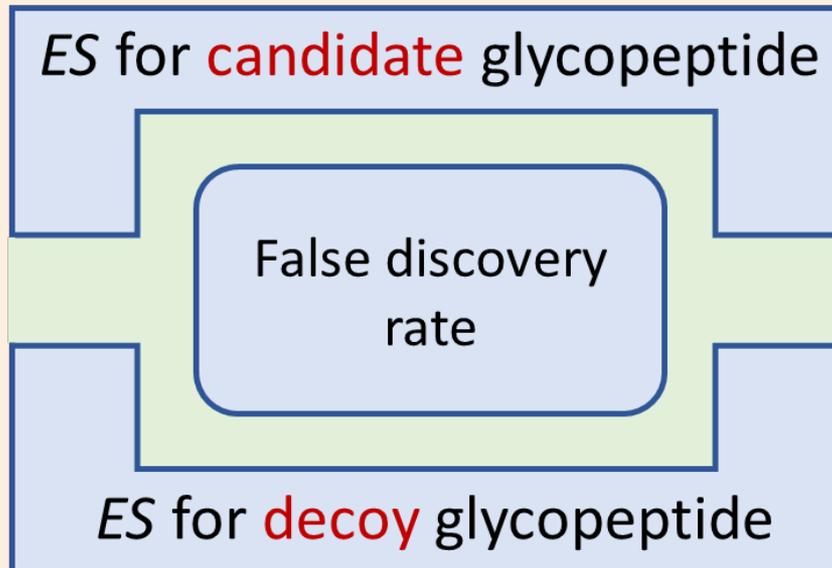
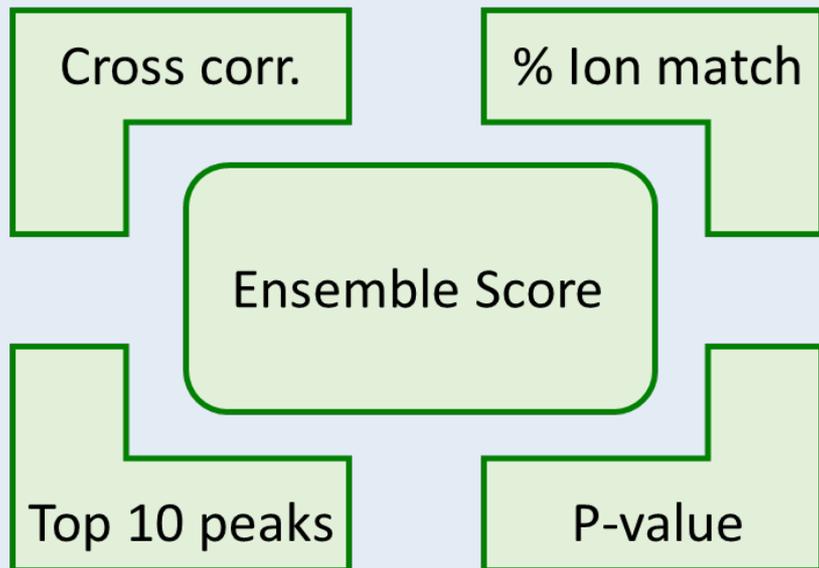


# Methods

## GlycoPAT: GUI



# GlycoPAT: ensemble score ( $ES$ ) and false discovery rate ( $FDR$ )

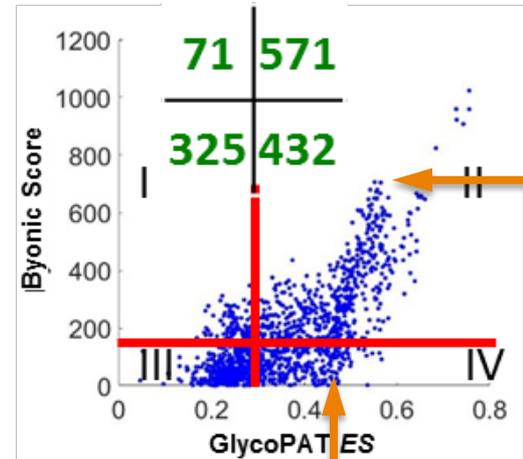


# GlycoPAT: scoring of HCD spectra

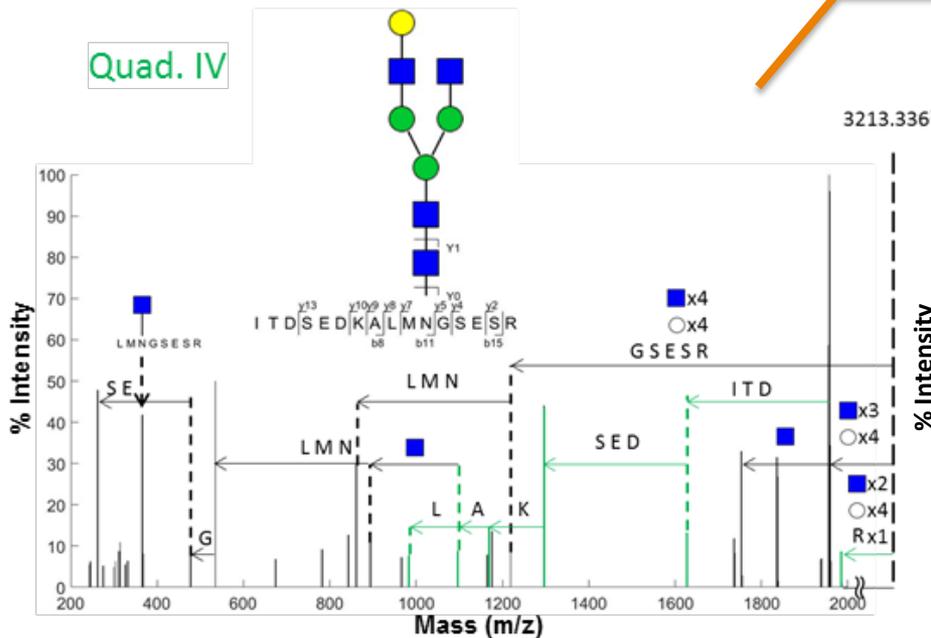
## Basigen | HCD

*J Proteome Res.* 15(10):3904-3915, 2016

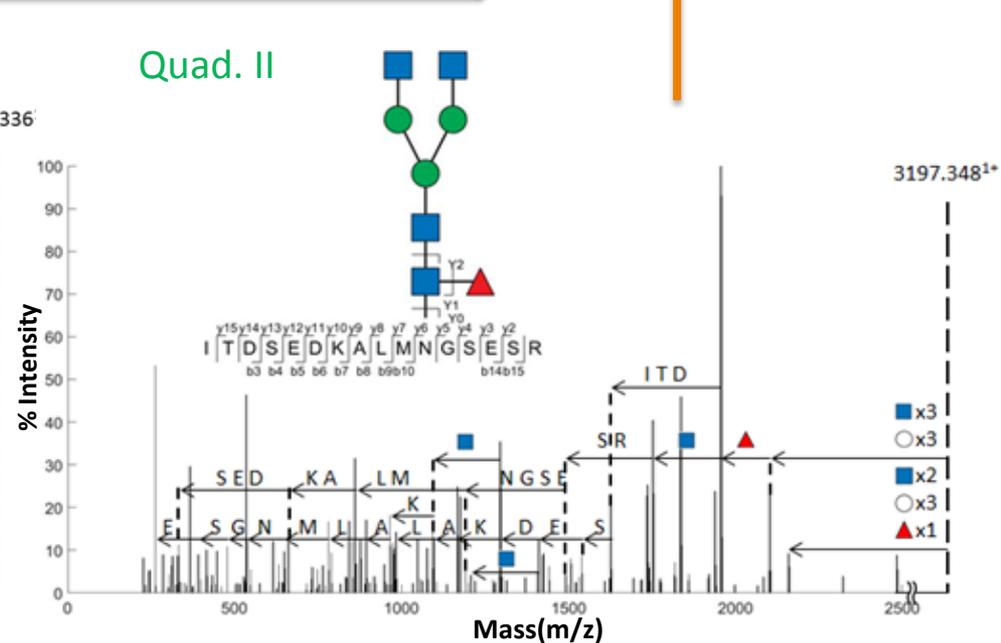
*\* Simultaneous breakage of peptide backbone and glycan structures*



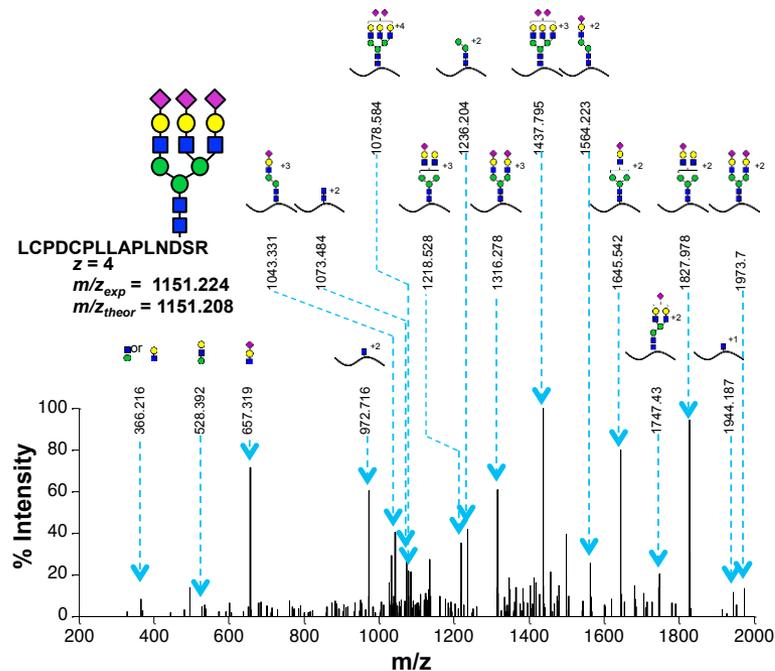
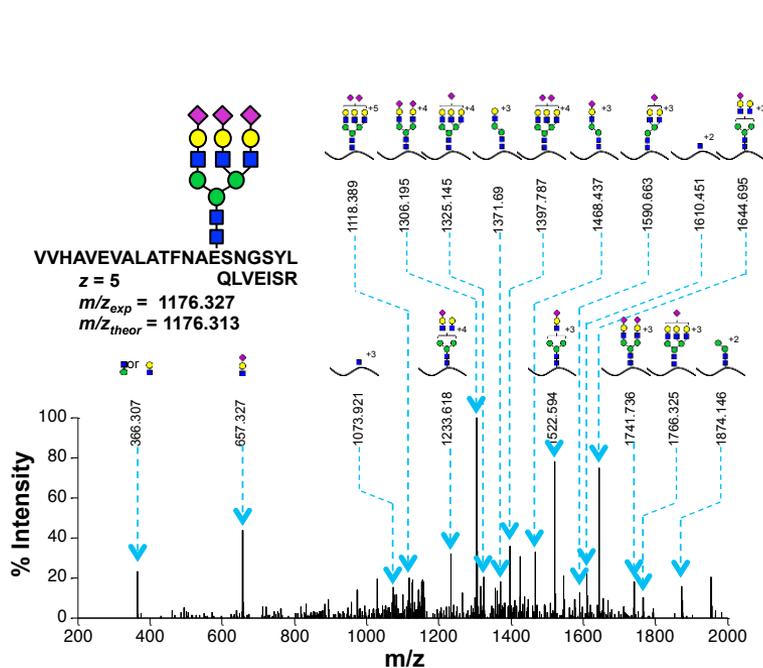
Quad. IV



Quad. II

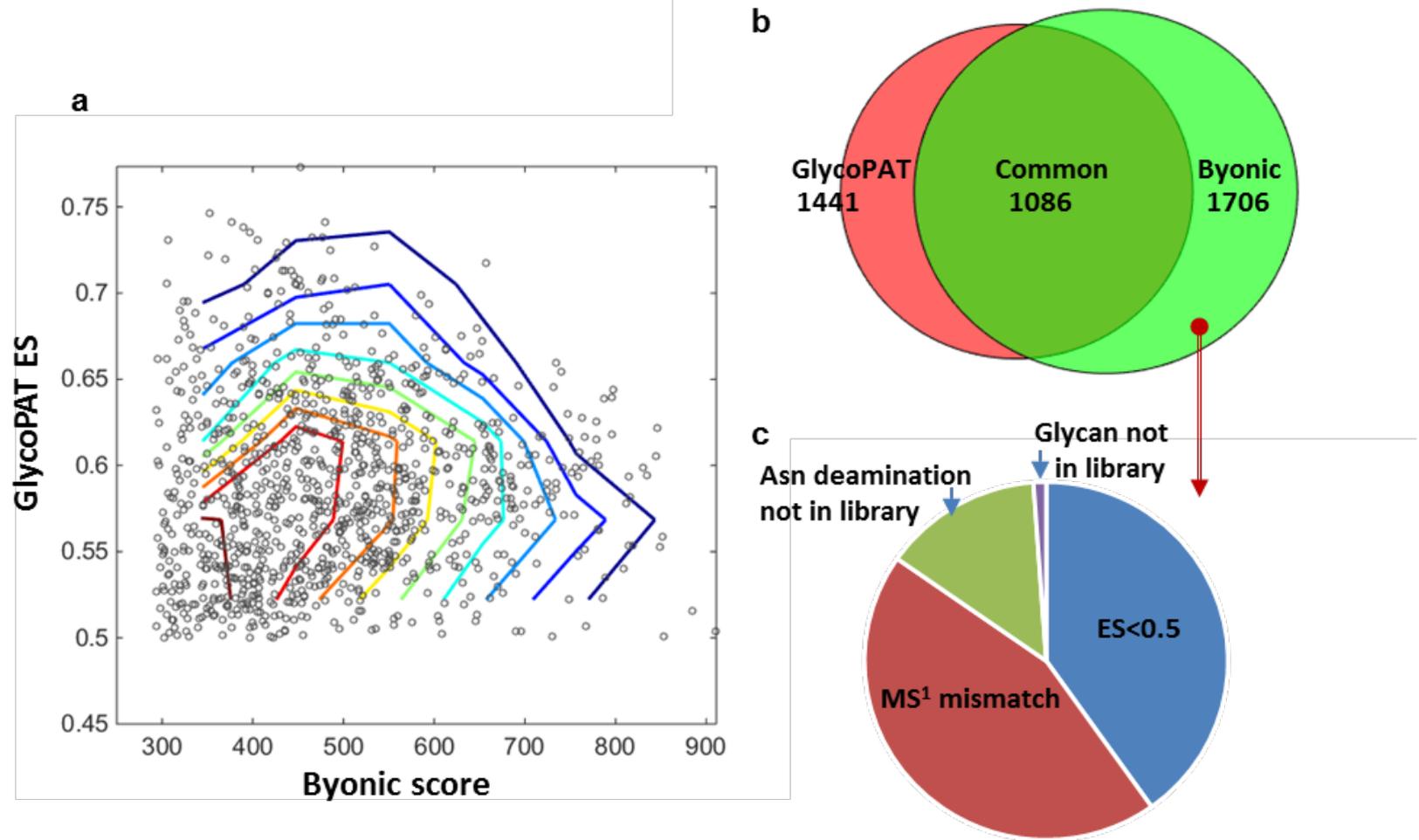


# GlycoPAT: scoring of CID spectra



*\* Analysis of ladder-like breakdown of glycans*

# Analysis of glycopeptides in whole prostate cancer lysates

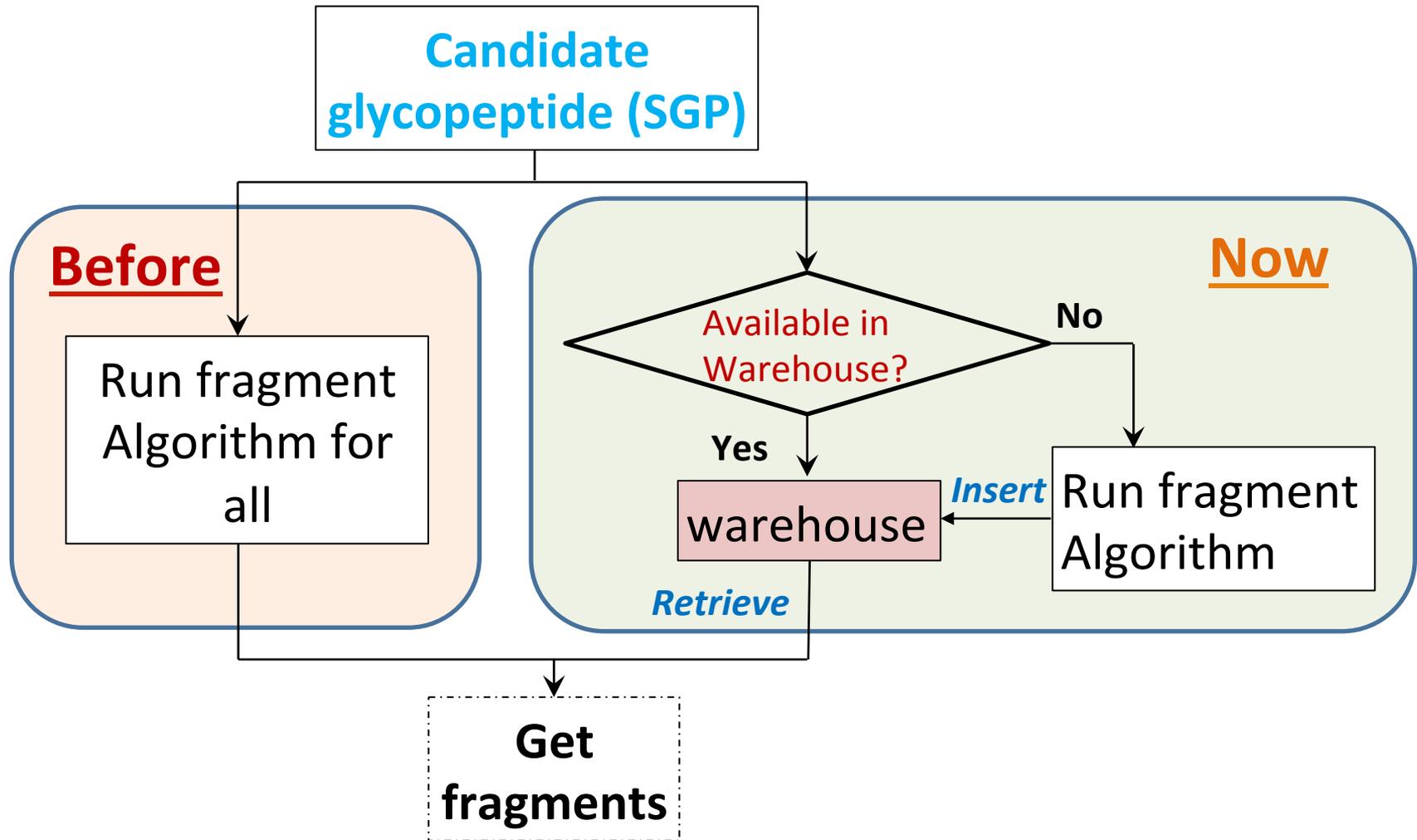


# Work in progress

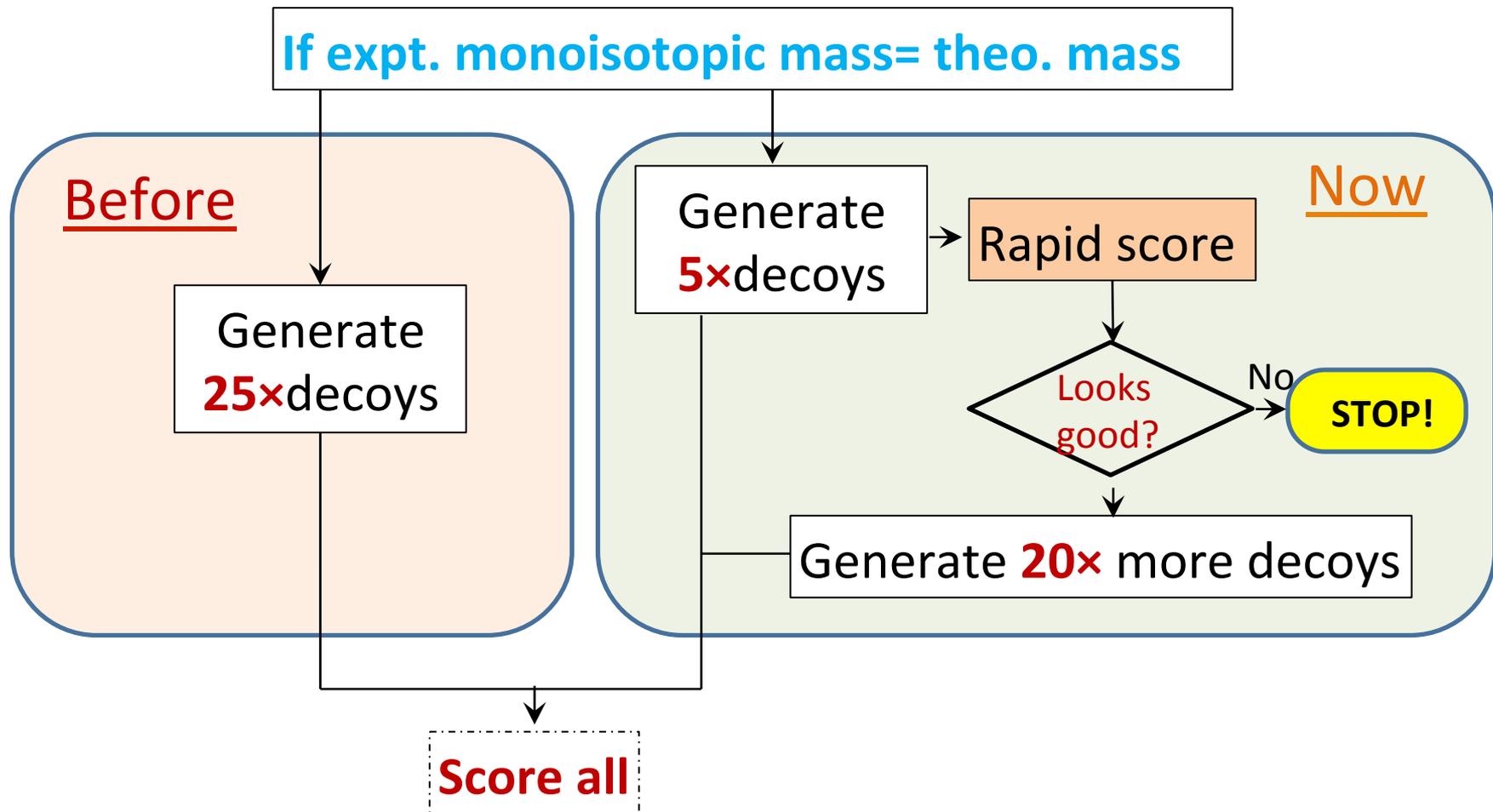
- Improving calculation speed:

***FOCUS!! Don't do everything for everyone!***

# Improve speed: *The fragment warehouse*



# Improved speed: *Selective scoring*



# Work in progress

- Improving calculation speed:

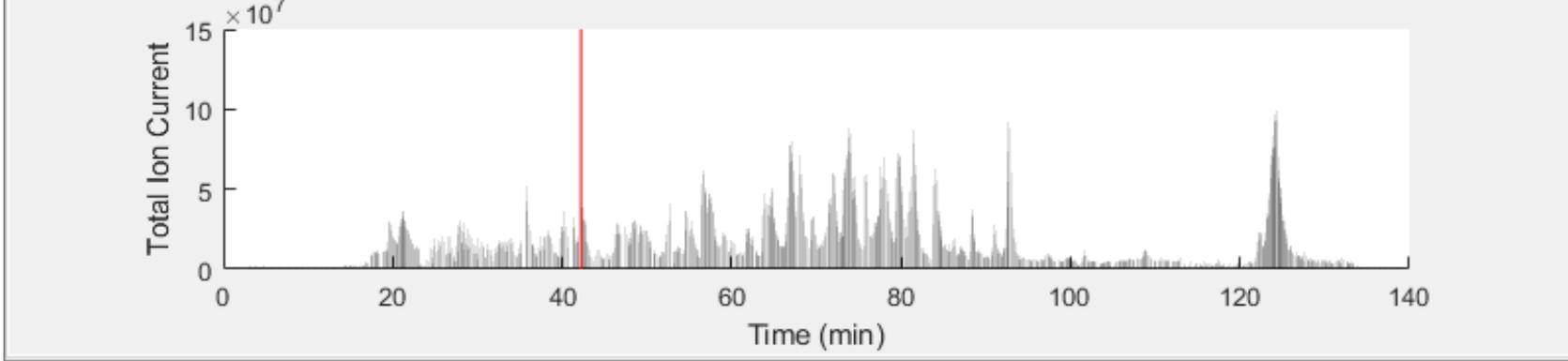
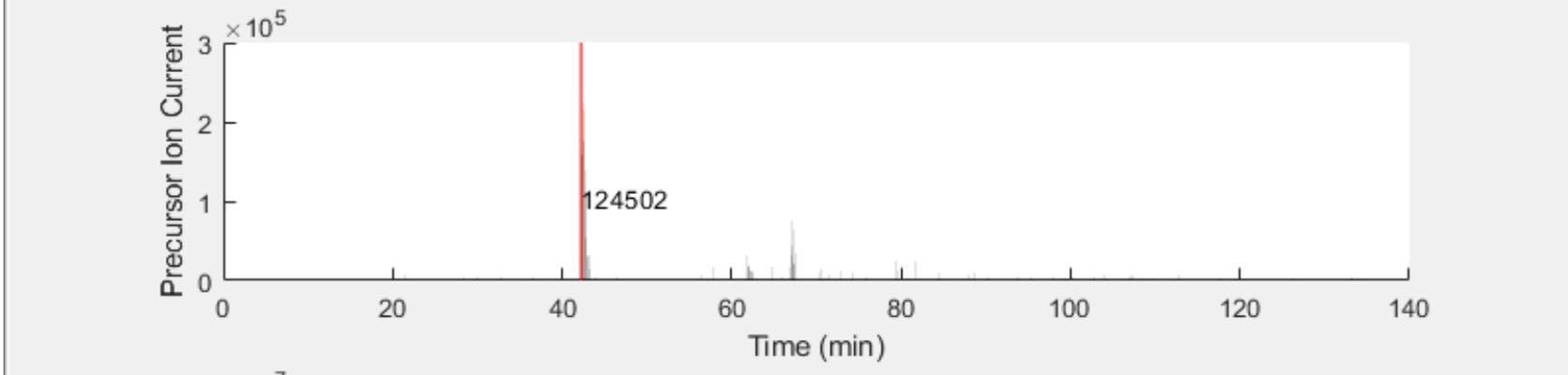
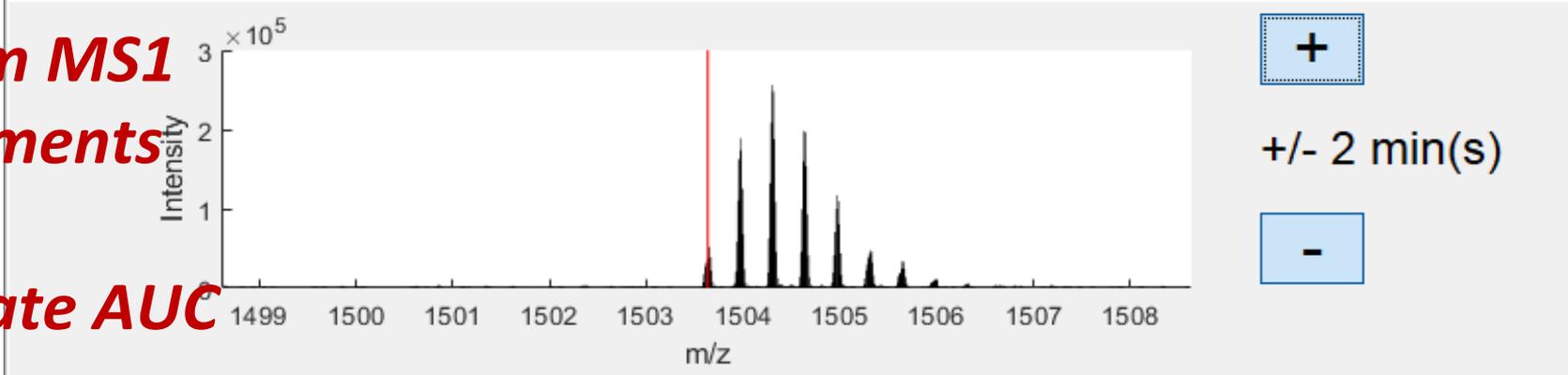
***FOCUS!! Don't do everything for everyone!***

- Streamlining result visualization:

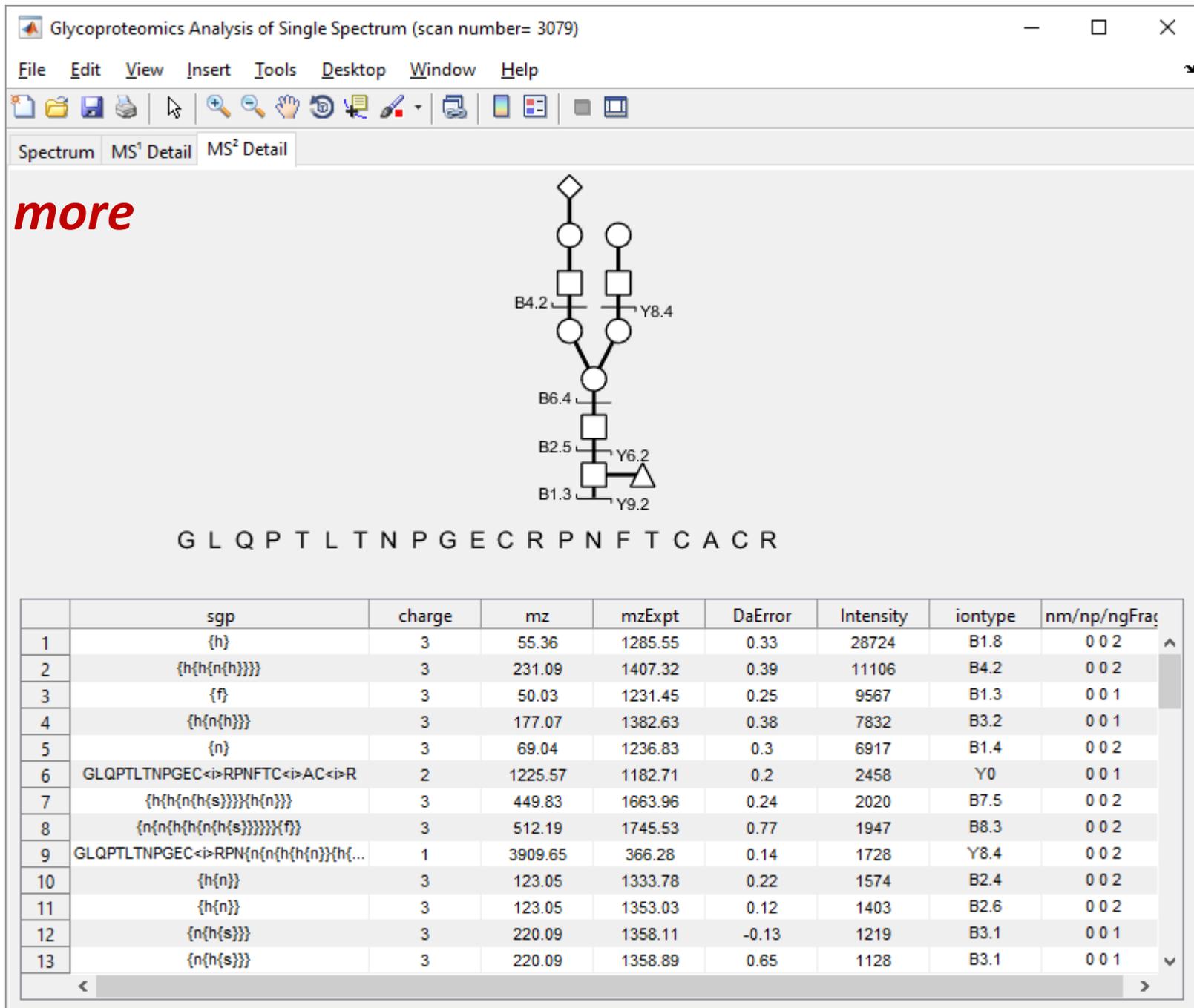
***SIMPLIFY!! Tell the story in pictures!***



**Confirm MS<sup>1</sup> assignments & calculate AUC**



**Provide more details**



# Conclusions

- **GNAT-Web**: Build reaction networks efficiently
  - Use this for *in silico* deterministic and stochastic simulations
  - Display of experimental data sets
- **DrawGlycan-SNFG**: Easy and Robust
- **GlycoPAT**: MS data analysis toolbox
  - Improve computational time
  - Integrate analysis from different fragmentation modes
  - Test in more biomedical applications

# Acknowledgements

## Lab members:

Anju Kelkar, Ph.D.

Virginia del Solar Fernandez, Ph.D.

## *Graduate students*

Ted Groth

Xinheng Yu

Arezoo Momeni

Changjie Zhang

Yuqi Zhu

Gabbie Pawlowski

Kai Cheng

Yusen Zhou

## Collaborators:

Alan Friedman and Jun Qu, Buffalo

Anne Dell, Stuart Haslam  
Imperial College

## Funding support:

**NHLBI** Systems Biology Collaborations



**NIGMS:** General Medicine



NYSTEM

NEW YORK STATE STEM CELL SCIENCE



American  
**Heart**  
Association®

