

1 **Syntactic Sugars: Crafting a Regular Expression Framework for Glycan Structures**

2 Alexander R. Bennett¹, Daniel Bojar^{2,*}

3 ¹Department of Medical Biochemistry, Institute of Biomedicine, University of Gothenburg, 41390
4 Gothenburg, Sweden.

5 ²Department of Chemistry and Molecular Biology, University of Gothenburg, 41390 Gothenburg,
6 Sweden. Wallenberg Centre for Molecular and Translational Medicine, University of Gothenburg,
7 41390 Gothenburg, Sweden.

8 *Corresponding author

9 **Abstract**

10 **Summary**

11 Structural analysis of glycans pose significant challenges in glycobiology due to their complex
12 sequences. Research questions such as analyzing the sequence content of the α 1-6 branch in *N*-
13 glycans, are biologically meaningful yet can be hard to automate. Here, we introduce a regular
14 expression system, designed for glycans, feature-complete, and closely aligned with regular
15 expression formatting. We use this to annotate glycan motifs of arbitrary complexity, perform
16 differential expression analysis on designated sequence stretches, or elucidate branch-specific
17 binding specificities of lectins in an automated manner. We are confident that glycan regular
18 expressions will empower computational analyses of these sequences.

19 **Availability and implementation**

20 Our regular expression framework for glycans is implemented in Python and is incorporated
21 into the open-source glycowork package (version 1.1+). Code and documentation are available
22 at <https://github.com/BojarLab/glycowork/blob/master/glycowork/motif/regex.py>.

23 **Contact:** daniel.bojar@gu.se

1

2 **Main**

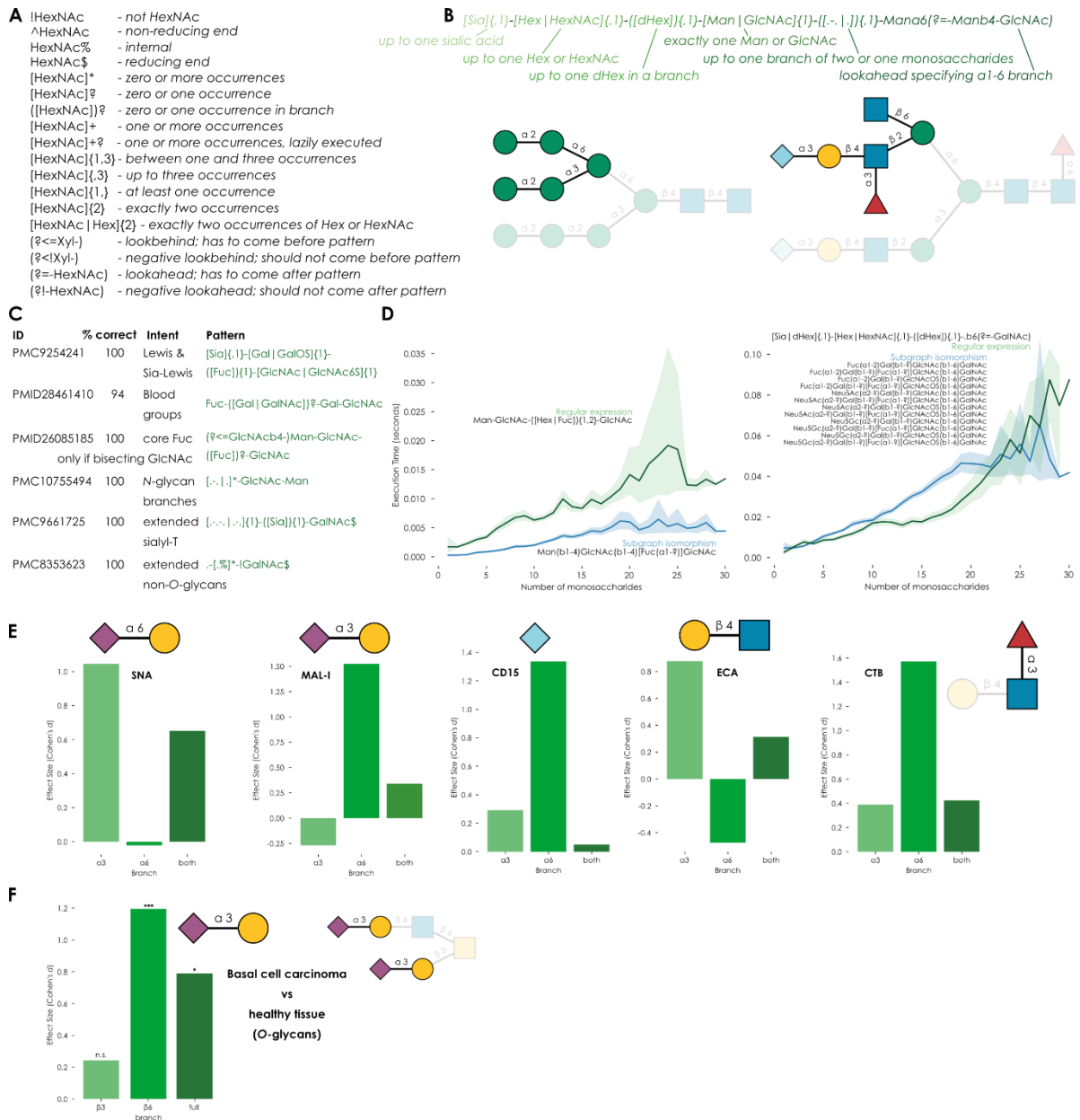
3 Glycans are complex carbohydrates, key in many biological processes (Varki, 2017).
4 Monosaccharides, linked together in specific patterns, are attached to proteins and lipids,
5 influencing their function and properties. Diversity in glycans arises from different
6 monosaccharide monomers and the ways in which these can be linked, leading to a vast array
7 of possible structures even with a limited set of building blocks (Bojar *et al.*, 2021).

8 The biological importance of glycans lies in motifs, distinct sequence patterns within the larger
9 structure. For example, the sialyl Lewis^x motif, Neu5Ac α 2-3Gal β 1-4(Fuc α 1-3)GlcNAc, is
10 crucial for cell adhesion and is implicated in cancer metastasis (Jin and Wang, 2020), while
11 high-mannose motifs are recognized by proteins from the immune system (Dommett *et al.*,
12 2006), and blood group antigens, such as the ABO blood groups (Stanley *et al.*, 2022), are
13 determined by specific glycan motifs on red blood cells.

14 Understanding glycans, and their motifs, is vital for unraveling their roles in health and disease
15 (Reily *et al.*, 2019), and in developing targeted therapeutic strategies. While motifs drive
16 biological functions, their position within the glycan is crucial for functionality. The simplest
17 example of this is distinguishing terminal and internal forms of the same motif, which already
18 affects biological properties, as many proteins or lectins only recognize one form (Bojar *et al.*,
19 2022). More distal context effects, such as whether a motif is presented on the α 1-3 or the,
20 often not outstretched (Fogarty *et al.*, 2020), α 1-6 branch in *N*-glycans, also can have drastic
21 effects, as most lectins have an arm preference for their binding motif (Li *et al.*, 2020).
22 Examples include MAL-I preferring Neu5Ac α 2-3 on the α 1-6 branch, whereas SNA strongly
23 favored Neu5Ac α 2-6 on the α 1-3 branch. Especially this latter category of motifs, for which
24 the larger sequence context is key, are hard to capture with existing methods.

1 Here, we present a regular expression (RegEx) system for glycans to search for, and extract,
2 sequence patterns of arbitrary complexity. Regular expressions are used in computer science
3 for pattern matching within text (Friedl, 2006), providing a concise and flexible means to find
4 specific sequences of characters within a string with a search pattern. Originally developed for
5 use in theoretical computer science, RegEx is now widely used, from simple string matching
6 to complex text processing and data extraction, using a series of special symbols to define
7 patterns, which can include specific words, numbers, or more complex text structures.

8 The formulaic (monosaccharides/linkages as building blocks) yet nonlinear nature of glycans
9 requires an adaptation of RegEx for optimal usage. Our RegEx framework closely mimics that
10 developed for text and supports ambiguity, modifiers, quantifiers, greedy and lazy execution,
11 as well as lookahead and lookbehind. A summary is shown in Fig. 1A, with further
12 explanations within glycowork (<https://bojarlab.github.io/glycowork/motif.html#regex>), since
13 our RegEx system is contained within our open-source Python package glycowork (Thomès *et*
14 *al.*, 2021) (v1.1+) as the new *glycowork.motif.regex* module, ensuring long-term maintenance.
15 This allows easy access to probe glycan sequences with RegEx patterns via the *get_match*
16 function. RegEx also enhances existing glycowork functionalities, such as specifying RegEx
17 patterns within motif highlighting capabilities of GlycoDraw (Lundstrøm, Urban, Thomès, *et*
18 *al.*, 2023) (Fig. 1B) or adding patterns to motif annotation and network highlighting functions.



1

2 **Figure 1. Developing a regular expression system for glycans.** A) Modifiers and quantifiers for the glycan
 3 RegEx system. B) An example query designed to extract the α 1-6 branch of *N*-glycans, illustrated with GlycoDraw
 4 via the “highlight_motif” functionality. C) Average accuracy of glycan RegEx queries across multiple datasets,
 5 checked manually. D) Speed comparison of glycan RegEx query versus normal motif search for core fucosylation
 6 or β 1-6 branch in *O*-glycans, against all human glycans up to 30 monosaccharides within glycowork. E) Sequence
 7 contexts of the α 1-3 branch, α 1-6 branch, and the entire glycan were extracted from the asymmetric *N*-glycan
 8 array data (Li *et al.*, 2020) and were analyzed with the *get_pvals_motifs* function of glycowork to establish branch-
 9 dependent enrichment patterns for various lectins. F) Differential glycomics expression of mono- and disaccharide

1 motifs from the β 1-3 and β 1-6 branch of basal cell carcinoma (Möginger *et al.*, 2018) were analyzed with the
2 *get_differential_expression* function of glycowork to identify how branch-specific motifs drive overall patterns.

3 Our RegEx system first chunks the provided pattern into homogeneous modules (Fig. 1B).
4 Segments without modifiers or quantifiers (but supporting structural ambiguity) are
5 immediately used for subgraph isomorphism operations to detect their location(s) within the
6 glycan graph as prepared by glycowork. Complete wildcards can be specified with
7 “Monosaccharide” or “.”, and linkages can be specified if desired (e.g., “Man α 6” or “Gal β 3/4”)
8 Segments with modifiers or quantifiers are processed into dictionaries of type (glycan
9 substructure: allowed counts), followed by subgraph isomorphisms of the substructures.
10 Glycan branches are indicated by parentheses. Lookahead and lookbehind operations include
11 the specified sequence in matching operations but exclude their node indices from extracted
12 matches. All this is followed by an iterative procedure of tracing a path through the matches of
13 the individual segments, such that all requirements of the regular expression are fulfilled (e.g.,
14 minimum/maximum counts) and the match represents a connected portion of the glycan graph.
15 Our RegEx system by default aims for a greedy match, akin to normal RegEx, yet lazy
16 execution is supported by a “?” modifier (e.g., “+?” instead of “+”).

17 We designed our glycan RegEx system to be functionally analogous to standard RegEx
18 formulation in Python, though we hasten to add that standard RegEx could not easily be used
19 to interrogate glycans in a similar manner, as it would only operate on the string level and fall
20 prey to ambiguities of glycan string nomenclatures, which we avoid with graph operations. An
21 example query (Fig. 1B) illustrates the potential of the herein developed method. Applied to a
22 variety of real-world datasets, followed by manual inspection of correct extraction, yields an
23 average accuracy of 99% (Fig. 1C), for the analyzed datasets.

24 Simpler glycan RegEx operations are computationally more expensive than motif searches
25 within glycowork (Fig. 1D), yet more complex queries, requiring many conventional motif

1 searches, gain in computational efficiency in comparison. Glycan RegEx allows us to extract
2 larger sequence contexts and use them in downstream analyses, such as the investigation
3 whether sequence patterns of the α 1-3/ α 1-6 branch in *N*-glycans differ across conditions. This
4 presents a non-trivial endeavor as the structural heterogeneity of *N*-glycans does not currently
5 enable a generalizable and easy way to extract α 1-6 branches across many sequences. We
6 applied our RegEx strategy to glycan array data from asymmetric glycans, to probe branch-
7 specific binding of well-known lectins (Fig. 1E). This confirmed strong preferences of SNA
8 for Neu5Ac α 2-6 on the α 1-3 branch and MAL-I recognizing Neu5Ac α 2-3 on the α 1-6 branch.
9 We further find enhancement of ECA binding by LacNAc repeats only on the α 1-3 branch, and
10 strong preferences for Neu5Gc binding of the α CD15 antibody, only on the α 1-6 branch.

11 To illustrate that this extends beyond glycan array analysis, we engaged in differential
12 glycomics expression analysis (Lundstrøm, Urban, and Bojar, 2023) to find out whether
13 features of specific *O*-glycan branches are differentially expressed in a skin cancer dataset
14 (Möginger *et al.*, 2018) (Fig. 1F). We then contrasted this with the current approach of
15 analyzing entire glycan sequences for differential expression and conclude that the increase in
16 α 2-3 linked sialylation in cancer, seen in the overall dataset, is entirely driven by the β 1-6
17 branch in this dataset, and hence exclusive to core 2/4/6 structures.

18 In summary, we envision that a glycan-specific RegEx system will facilitate analyzing glycans
19 and their changes in disease. As this method allows for the analysis of unprecedented and
20 arguably arbitrary motif complexity, we expect the discovery of new patterns in glycan
21 dysregulation. We are also convinced that this approach further strengthens the usefulness of
22 the glycowork ecosystem and will, in general, empower glycoinformatics workflows.

23 **Author contributions**

1 Conceptualization: D.B., Funding Acquisition: D.B., Resources: D.B., Software: A.R.B., D.B.,
2 Supervision: D.B., Visualization: A.R.B., D.B., Writing—Original Draft Preparation: D.B.,
3 Writing—Review & Editing: A.R.B., D.B.

4 **Funding**

5 This work was supported by a Branco Weiss Fellowship – Society in Science awarded to D.B.;
6 by the Knut and Alice Wallenberg Foundation; and the University of Gothenburg, Sweden.

7 *Conflict of Interest:* none declared.

8 **Data availability**

9 Code and documentation are available via glycowork and
10 <https://github.com/BojarLab/glycowork/blob/master/glycowork/motif/regex.py>.

11 **References**

- 12 Bojar,D. *et al.* (2022) A Useful Guide to Lectin Binding: Machine-Learning Directed
13 Annotation of 57 Unique Lectin Specificities. *ACS Chem. Biol.*, [acschembio.1c00689](#).
14 Bojar,D. *et al.* (2021) Deep-Learning Resources for Studying Glycan-Mediated Host-
15 Microbe Interactions. *Cell Host & Microbe*, **29**, 132-144.e3.
16 Dommett,R.M. *et al.* (2006) Mannose-binding lectin in innate immunity: past, present and
17 future. *Tissue Antigens*, **68**, 193–209.
18 Fogarty,C.A. *et al.* (2020) How and why plants and human N-glycans are different: Insight
19 from molecular dynamics into the “glycoblocks” architecture of complex
20 carbohydrates. *Beilstein J. Org. Chem.*, **16**, 2046–2056.
21 Friedl,J.E.F. (2006) Mastering regular expressions 3rd ed. O’Reilly, Sebastapol, CA.

- 1 Jin,F. and Wang,F. (2020) The physiological and pathological roles and applications of sialyl
2 Lewis x, a common carbohydrate ligand of the three selectins. *Glycoconj J*, **37**, 277–
3 291.
- 4 Li,L. *et al.* (2020) Microarray analyses of closely related glycoforms reveal different
5 accessibilities of glycan determinants on N-glycan branches. *Glycobiology*, **30**, 334–
6 345.
- 7 Lundstrøm,J., Urban,J., and Bojar,D. (2023) Decoding glycomics with a suite of methods for
8 differential expression analysis. *Cell Reports Methods*, 100652.
- 9 Lundstrøm,J., Urban,J., Thomès,L., *et al.* (2023) GlycoDraw: a python implementation for
10 generating high-quality glycan figures. *Glycobiology*, cwad063.
- 11 Möglinger,U. *et al.* (2018) Alterations of the Human Skin N- and O-Glycome in Basal Cell
12 Carcinoma and Squamous Cell Carcinoma. *Front. Oncol.*, **8**, 70.
- 13 Reily,C. *et al.* (2019) Glycosylation in health and disease. *Nat Rev Nephrol*, **15**, 346–366.
- 14 Stanley,P. *et al.* (2022) Structures Common to Different Glycans. In, Varki,A. *et al.* (eds),
15 *Essentials of Glycobiology*. Cold Spring Harbor Laboratory Press, Cold Spring
16 Harbor (NY).
- 17 Thomès,L. *et al.* (2021) Glycowork: A Python package for glycan data science and machine
18 learning. *Glycobiology*, cwab067.
- 19 Varki,A. (2017) Biological roles of glycans. *Glycobiology*, **27**, 3–49.
- 20