



UNIVERSIDADE FEDERAL DO RIO GRANDE DO NORTE
CENTRO DE CIÊNCIAS EXATAS E DA TERRA
DEPARTAMENTO DE INFORMÁTICA E MATEMÁTICA APLICADA
DISCIPLINA DE LINGUAGENS FORMAIS E AUTÔMATOS



Análise do uso de RegEx na biblioteca Glycowork

Andriel Vinicius de Medeiros Fernandes,
Gabrielle de Vasconcelos Borja,
Jeremias Pinheiro de Araújo Andrade,
Lucas Vinicius Dantas de Medeiros,
María Paz Marcato,
Ramon Cândido Jales de Barros

Natal-RN
Outubro, 2025

Nome completo do autor

Análise de uso de RegEx na biblioteca Glycowork

Pesquisa elaborada e apresentada para a disciplina DIM0606 - Linguagens Formais e Autômatos, ofertada pelo Departamento de Informática e Matemática Aplicada da Universidade Federal do Rio Grande do Norte e ministrada no semestre 2025.2 pelo Prof. Dr. Valdicleis da Silva Costa, como requisito parcial para a obtenção de nota para a 1ª unidade.

DIMAP – DEPARTAMENTO DE INFORMÁTICA E MATEMÁTICA APLICADA
CCET – CENTRO DE CIÊNCIAS EXATAS E DA TERRA
UFRN – UNIVERSIDADE FEDERAL DO RIO GRANDE DO NORTE

Natal-RN

Outubro, 2025

Análise do uso de RegEx na biblioteca Glycowork

Autores: Andriel Vinicius de Medeiros Fernandes,
Gabrielle de Vasconcelos Borja,
Jeremias Pinheiro de Araújo Andrade,
Lucas Vinicius Dantas de Medeiros,
María Paz Marcato,
Ramon Cândido Jales de Barros
Professor: Valdigleis

RESUMO

O resumo deve apresentar de forma concisa os pontos relevantes de um texto, fornecendo uma visão rápida e clara do conteúdo e das conclusões do trabalho. O texto, redigido na forma impessoal do verbo, é constituído de uma seqüência de frases concisas e objetivas e não de uma simples enumeração de tópicos, não ultrapassando 500 palavras, seguido, logo abaixo, das palavras representativas do conteúdo do trabalho, isto é, palavras-chave e/ou descritores. Por fim, deve-se evitar, na redação do resumo, o uso de parágrafos (em geral resumos são escritos em parágrafo único), bem como de fórmulas, equações, diagramas e símbolos, optando-se, quando necessário, pela transcrição na forma extensa, além de não incluir citações bibliográficas.

Palavras-chave: Palavra-chave 1, Palavra-chave 2, Palavra-chave 3.

Lista de figuras

Lista de tabelas

Lista de abreviaturas e siglas

Sumário

1 Referencial Teórico

O trabalho feito na Universidade de Gotemburgo, na Suécia, tem como bases duas áreas distintas, a biologia e a ciência da computação. Na área biológica, o foco é sobre os glicanos, biopolímeros com inúmeras aplicações e funções biológicas. Já na computação, foi estudado o uso de expressões regulares (RegEx), padrões formados por sequências de caracteres utilizados para busca, análise e manipulação de textos. Dessa forma, os pesquisadores aplicaram expressões regulares para formalizar uma estrutura de busca voltada à identificação de padrões específicos dos glicanos.

1.1 Glicanos: Complexidade e Importância Biológica

Para compreender melhor a abordagem biológica do estudo, é importante explorar os glicanos, suas características e funções, bem como outros conceitos relacionados abordados no trabalho.

Definição 1.1.1 (Definição de Glicanos). *Glicanos são polissacarídeos estruturais, longas cadeias formadas por unidades de açúcar (monossacarídeos) ligadas entre si por ligações glicosídicas, presentes abundantemente na Terra como componentes importantes de estruturas como glicoproteínas, glicolipídeos e proteoglicanos, além de fazerem parte de paredes celulares de fungos e leveduras.*

Essa diversidade de funções ocorre porque, em sua composição, existem subestruturas chamadas de motivos de glicanos, uma sequência ou arranjo particular de açúcares que funciona como um sinal de reconhecimento molecular, ou seja, a parte que carrega o significado mais importante e que é reconhecida por outras moléculas. São eles que trazem a importância biológica dos glicanos e, por isso, é de extrema importância que eles sejam entendidos pela ciência.

Contudo, a grande diversidade e versatilidade traz à estrutura interesse contínuo de pesquisas de diversas áreas, como a imunologia, biotecnologia e parasitologia. Porém,

como existem inúmeros motivos de glicanos, catalogá-los e reconhecê-los é um desafio.

1.2 RegEx na Ciência da Computação

Na ciência da computação, RegEx são uma sequência de caracteres que define um padrão de busca. Elas representam uma linguagem formal que pode ser utilizada para identificar, extrair e manipular subconjuntos de texto com base em regras e padrões específicos.

Essa notação envolve uma combinação de cadeias de símbolos do alfabeto da linguagem. Por exemplo, o operador main (+) é utilizado para denotar união, o asterisco (*) para o fecho estrela, indica zero ou mais ocorrências do caractere anterior, e o ponto (.) para concatenação. Essa sintaxe permite a criação de padrões que vão desde buscas simples, como encontrar uma palavra, até a validação de estruturas complexas em um texto.

1.3 A Aplicação Inovadora de RegEx em Glicanos

Diante do desafio na análise dos glicanos, o estudo propôs a aplicação de expressões regulares como ferramenta para buscar e extrair seus padrões estruturais de forma precisa e flexível. Essa técnica foi adaptada no trabalho para lidar com a estrutura ramificada e não linear dos glicanos, permitindo uma análise muito mais eficiente e escalável. O artigo demonstra que, embora as aplicações tradicionais de RegEx sejam conhecidas, sua versatilidade permite o uso em contextos menos convencionais, como a bioinformática.

2 Capítulo 2

Neste capítulo, descrevemos o funcionamento do sistema **Glycan RegEx** implementado no módulo `glycowork.motif.regex`, suas funções principais, o papel das expressões regulares e as vantagens do uso desse sistema.

O **Glycan RegEx**, introduzido por Bennett e Bojar (2024) no pacote `glycowork`, representa um avanço significativo na análise computacional de glicanos ao acrescentar o uso de RegEx para identificação e extração de motivos. Tal sistema foi proposto como uma adaptação das tradicionais RegEx de ciência da computação, permitindo a detecção de motivos na estrutura não linear dos glicanos, possibilitando buscas precisas e otimizadas nestes elementos.

2.1 Estrutura e funções do módulo

O **Glycan RegEx** se baseia na tradução de padrões RegEx para operações de isomorfismo de subgrafos dentro das estruturas moleculares de glicanos. Quando o usuário fornece um padrão o sistema decompõe esse padrão em unidades menores, chamadas de módulos homogêneos, correspondentes a monossacarídeos e ligações individuais, e cada módulo é processado para identificar as possíveis correspondências no grafo do glicano, permitindo localizar subestruturas equivalentes de forma independente da forma textual em que o glicano foi representado. Isso é essencial, já que o módulo aceita diversos formatos de entrada, como GlycoCT, Oxford e demais notações.

Estas RegEx glicosídicas têm funcionamento muito similar às clássicas, com suporte a modificadores, quantificadores e operadores de busca contextual, como *lookahead* e *lookbehind*. Isso permite representar características como a quantidade de ocorrências de um módulo, ligações opcionais, dentre outras; o sistema também aceita os curingas "." e `Monosaccharide`, que representam estruturas genéricas para a busca. Em relação às novidades, o sistema permite especificar subgrafos específicos a partir da representação por

parênteses.

TODO: especificar as diferenças (o chat respondeu!)

O funcionamento interno do **Glycan RegEx** ocorre, a princípio, com a segmentação da RegEx em módulos homogêneos: cada módulo é classificado como simples quando não há a presença de modificadores ou quantificadores ou complexo quando há; módulos simples são armazenados no formato de **string**, enquanto os complexos são salvos em dicionários. Em seguida, cada módulo é convertido em um grafo e por meio de operações de isomorfismo de subgrafos é detectado onde cada subgrafo se encaixa no glicano completo. Durante tal processo de detecção ocorre também a aplicação dos modificadores e quantificadores dos módulos complexos, o que permite descartar de forma otimizada subgrafos que não atinjam aos requisitos definidos, como número específico de ocorrências do padrão.

Ao fim, com o processamento de todos os módulos, ocorre a construção do caminho através do glicano que une todas as correspondências parciais e validadas pelos requisitos definidos na expressão completa reconstruindo, dessa forma, o motif exato a ser buscado. Essa abordagem possibilita representar ligações específicas, como $\alpha 1-3$ ou $\beta 1-6$, bem como ambiguidade estrutural e ramificações expressas por meio de parênteses. Dessa forma, a notação **RegEx** é adaptada à topologia molecular, permitindo uma modelagem altamente expressiva e precisa das estruturas glicosídicas.

Todo este processo é executado por diversas funções internas da biblioteca, dentre as quais se destacam:

- **preprocess_pattern()**: Particiona a RegEx fornecida em uma lista de padrões menores (módulos);
- **process_complex_pattern()**: Verifica por meio de isomorfismo de subgrafos se determinado padrão está presente no glicano;
- **match_it_up()**: Para cada módulo, formata este em uma cadeia de caracteres ou dicionário, executa a função **process_complex_pattern()** e retorna os módulos que representam os subgrafos encontrados;
- **trace_path()**: Conecta cada subgrafo encontrado na estrutura do glicano e retorna o caminho completo;
- **get_match()**: Extrai trechos da estrutura do glicano que correspondem ao motif buscado pelo usuário, concentrando todo o processo por chamadas às demais funções.

2.2 Vantagens da RegEx no contexto do glycowork

Como comentado, a inovação trazida pelo Glycan RegEx ao framework glycowork consiste na união do método tradicional de busca por motifs com isomorfismo de subgrafo e o uso de um sistema de expressões regulares específico para glicanos.

Sem o uso das expressões regulares, há a necessidade de uma base estática de motifs já conhecidos. Essa base pode ser observada no arquivo `common_names.json` do próprio framework. Além disso, torna-se necessário o uso de múltiplas combinações de padrões para detectar motifs complexos e específicos do glicano, bem como lidar com a dificuldade de capturar contextos mais amplos, como a ocorrência de um motif presente apenas em um ramo do glicano. Todos esses fatores culminam em um uso bastante rígido e, por vezes, difícil de representar de forma computacional.

Com a introdução do sistema Glycan RegEx é possível escrever padrões de buscas de motifs em um formato mais declarativo e eficaz, com uso de operadores clássicos que permitem negação, repetição e demais operações. Dessa forma, o algoritmo de busca de motifs tem a capacidade de gerar diversos subgrafos intermediários (ao contrário do mecanismo anterior), encontrando a combinação ideal de motifs de forma refinada e mais simples para o programador.

Esta camada adicional que o sistema traz pode gerar um *overhead* em estruturas mais simples do glicano quando comparado com a busca tradicional do glycowork por motifs. Entretanto, quando aplicado em glicanos mais complexos ou maiores, tal mecanismo se sobressai tanto em sua eficácia, podendo ser mais simples escrever uma única expressão regular que atenda aos critérios, quanto em sua eficiência.

3 Capítulo 3

4 Considerações finais

As considerações finais formam a parte final (fechamento) do texto, sendo dito de forma resumida (1) o que foi desenvolvido no presente trabalho e quais os resultados do mesmo, (2) o que se pôde concluir após o desenvolvimento bem como as principais contribuições do trabalho, e (3) perspectivas para o desenvolvimento de trabalhos futuros, como listado nos exemplos de seção abaixo. O texto referente às considerações finais do autor deve salientar a extensão e os resultados da contribuição do trabalho e os argumentos utilizados estar baseados em dados comprovados e fundamentados nos resultados e na discussão do texto, contendo deduções lógicas correspondentes aos objetivos do trabalho, propostos inicialmente.

4.1 Principais contribuições

Texto.

4.2 Limitações

Texto.

4.3 Trabalhos futuros

Texto.