

O-GLYCBASE version 4.0: a revised database of O-glycosylated proteins

Ramneek Gupta*, Hanne Birch, Kristoffer Rapacki, Søren Brunak and Jan E. Hansen

Center for Biological Sequence Analysis, The Technical University of Denmark, DK-2800 Lyngby, Denmark

Received October 13, 1998; Revised October 16, 1998; Accepted October 22, 1998

ABSTRACT

O-GLYCBASE is a database of glycoproteins with O-linked glycosylation sites. Entries with at least one experimentally verified O-glycosylation site have been compiled from protein sequence databases and literature. Each entry contains information about the glycan involved, the species, sequence, a literature reference and http-linked cross-references to other databases. Version 4.0 contains 179 protein entries, an approximate 15% increase over the last version. Sequence logos representing the acceptor specificity patterns for GalNAc, GlcNAc, mannosyl and xylosyl transferases are shown. The O-GLYCBASE database is available through the WWW at <http://www.cbs.dtu.dk/databases/OGLYCBASE/>

INTRODUCTION

Glycosylation is an important post-translational modification of proteins and affects more than half of all the proteins in a cell (1). O-linked glycosylation refers to the event in which a carbohydrate is covalently linked to the hydroxyl group of serine or threonine, and is one of the three main types of post-translational modifications involving carbohydrates. This modification influences a number of properties of proteins (2,3) including proteolytic resistance, solubility, immunological properties, ligand binding and is also involved in recognition during sperm-egg binding. Glycosylation is important for the biotechnology industry (4,5) since altered glycosylation has been shown to impair protein function (6-8).

Mapping the exact sites of modification is important, both in understanding the glycosylation process per se and in inferring interactions with other processes such as phosphorylation (9). Additionally, it is important while designing constructs for recombinant glycoproteins to ensure that acceptor sites are placed within the correct sequence motif. While current experimental methods for mapping glycosylation sites (such as Edman degradation and mass spectrometry) are cumbersome, they are presently the only precise way of accurate determination. However, with the large amount of experimental data being made available, this data can be employed for constructing prediction methods for acceptor specificities. Thus, an added-value dataset

Table 1.

O-linked sugar	Abbreviation	Number of sites in O-GlycBase 4.0
N-Acetylgalactosamine	GalNAc	680 (194 Ser, 486 Thr)
N-Acetylglucosamine	GlcNAc	88
Mannose	Man	158
Xylose	Xyl	16
Glucose	Glc	11
Fucose	Fuc	7
Others/Unspecified	-	31

of known and verified O-glycosylation sites would facilitate this research.

NEW FEATURES IN VERSION 4.0

Since version 3.0, 88 new O-glycosylation sites have been included in the database. A few earlier entries have been revised from recent literature and a few slight inconsistencies were corrected. Version 4.0 includes Medline references for the literature quoted, and database cross-references have been http-linked to sources available on the web.

O-GLYCOSYLATION TYPES AND SEQUENCE MOTIFS

The database includes entries with six types of sugars O-linked to the hydroxy amino acid (Table 1).

Only a minor fraction of serine and threonine residues are modified by glycosylation and no clear consensus sequence exists for acceptor sites in most types of O-glycosylation. While certain rules and acceptor motifs have been proposed (10-16), there are no definite rules which distinguish glycosylated hydroxy amino acid residues from non-glycosylated residues. Sequence logos displaying the fuzzy sequence context for mucin-type (GalNAc), GlcNAc, mannose and xylose glycosylation are shown in Figure 1.

DATA SOURCES

Proteins with a carbohydrate assigned either to a serine or a threonine residue were extracted from SWISS-PROT (17) (release 36) and PIR (18) (release 50) databases as well as directly from published reports. Only a few entries were submitted directly to the database. All proteins included in our database have at least one experimentally verified O-glycosylation site and consist largely of glycoproteins expressed *in vivo*.

*To whom correspondence should be addressed. Tel: +45 4525 2472; Fax: +45 4593 1585; Email: ramneek@cbs.dtu.dk

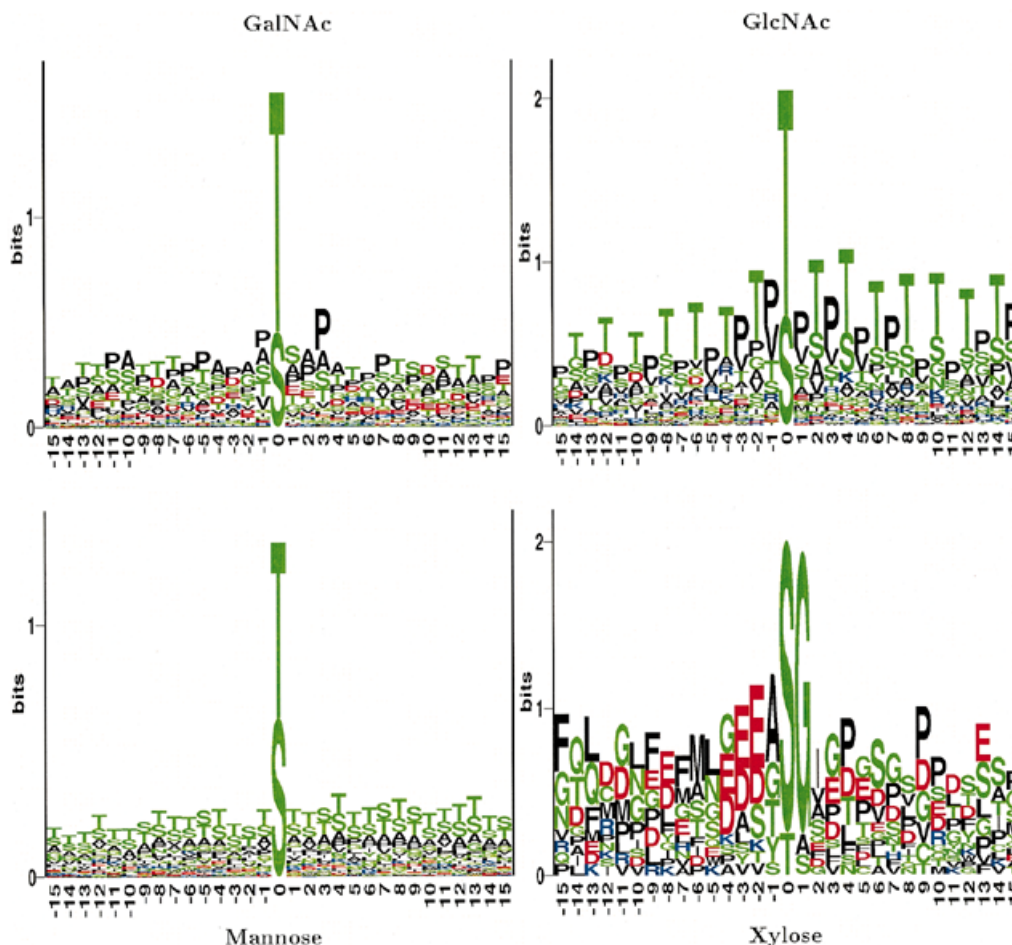


Figure 1. Acceptor site specificities are shown in the form of Shannon information content (21) represented as sequence logos (22). Logos are shown for GalNAc, GlcNAc, mannose and xylose-linked glycosylations and incorporate all sites available in OGLYCBASE 4.0 (Table 1). Sites were aligned with the glycosylated residue at position 0. The height of each column reflects the bias from random of particular residues surrounding the glycosylation site. The size of each residue letter reflects the frequency of the residue on that position. The height of the central serine/threonine has been rescaled to magnify the context of other positions, and is thus non-informative. Neutral and polar amino acids are shown in green, acidic and basic charged residues in red and blue, respectively, and the neutral and hydrophobic in black. The high frequency of serine and threonine in the vicinity of the glycosylated position 0 is largely indicative of the close clustering of O-glycosylation sites in general (11,12,20,23). All four acceptor patterns display a low frequency of charged amino acids at position -1 as originally found for the GalNAc transferase (24). GalNAc and GlcNAc acceptor sites are influenced to a large extent by proline vicinity (23–26). Further, the clustering of GlcNAc acceptor sites indicates an even positioning of acceptor sites and an odd positioning of proline/valine sites (Gupta *et al.*, in preparation). For mannose modified sites, the most frequent hydrophobic residue is alanine on all positions. Xylose acceptor sites agree with the characteristic ‘SGXG’ motif suggested for glycosaminoglycan attachment (16).

DESCRIPTION

Version 4.0 (7 Oct, 1998) of O-GLYCBASE contains 179 glycoprotein entries with 991 experimentally determined O-glycosylation sites. The entries consist of 11 fields (Fig. 2) including database and literature references, the O-linked glycan, sequence and a list of O-linked sites. N-glycosylated sites have also been marked on the included proteins where information was readily available. The database is extensively cross-referenced to sequence and structure databases including SWISS-PROT, PROSITE, TREMBL, PIR, PDB, EMBL, HSSP, GenPept and GenBank. The database is available in a plain-text as well as an html format with http links to Medline abstracts (maintained at PubMed, NCBI) and database sources on the internet. The database will be continually updated with the availability of more information. For a more comprehensive description of the database, see Hansen *et al.* (19).

PREDICTION OF O-GLYCOSYLATION

We have earlier reported a neural network based prediction method of mucin type O-glycosylation on mammalian proteins, NetOGlyc 2.0 (20). The method is available at the URL: <http://www.cbs.dtu.dk/services/NetOGlyc/> both as a mail and web server. We have recently also made available a prediction method for α -linked *N*-acetylglucosamine (GlcNAc) acceptor sites on *Dictyostelium discoideum* proteins at <http://www.cbs.dtu.dk/services/DictyOGlyc/> (Gupta *et al.*, in preparation). Further prediction servers for other types of O-glycosylation will be developed when a statistically sufficient number of acceptor sites is available.

Predictions indicate that more than one-third of all SWISS-PROT entries may be O-glycosylated, which is in sharp contrast to the very small number of entries annotated as O-glycosylated.

