

Datenbanken und Bioinformatik-Werkzeuge für die Glykobiologie

Claus-W. von der Lieth, Andreas Bohne-Lang, Klaus K. Lohmann
Deutsches Krebsforschungszentrum – Zentrale Spektroskopie (B090), Heidelberg

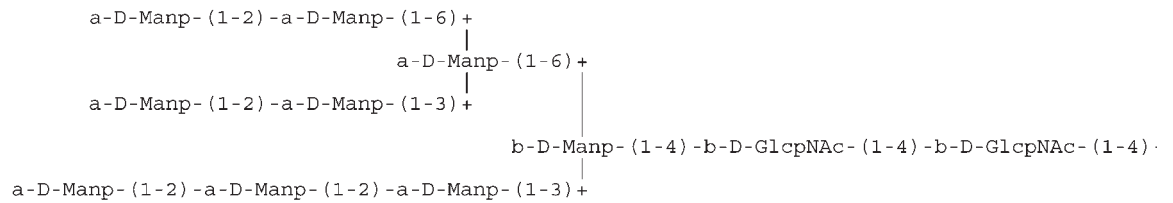


Abb. 1a: Vollständige Beschreibung eines N-Glykans im CarbBank-Format. Anstelle der griechischen Buchstaben α und β zur Bezeichnung der Konfiguration am anomeren C-Atom werden *a* und *b* verwendet.

Das menschliche Genom scheint nicht mehr als 30.000 bis 40.000 Proteine zu kodieren. Eine der überraschenden Erkenntnisse, die seine Sequenzierung zu Tage gefördert hat, ist die im Vergleich mit anderen Organismen relativ geringe Anzahl von Genen. Es ist daher eine wichtige Herausforderung zu untersuchen, wie ko- und posttranslationale Modifikationen die Eigenschaften und Funktionen von Proteinen beeinflussen und moderieren können. Glykosylierung, die Anheftung von Zuckerstrukturen an Residuen auf der Proteinoberfläche, ist eine der am häufigsten vorkommenden Modifikationen, die strukturell sehr vielfältiger Art sein können^[1]. Gemessen an der Gesamtmasse eines Proteins können Glykosylierungen einen recht großen Anteil ausmachen. Ausgehend von gut untersuchten Glykoproteinen, die in Proteindatenbanken enthalten sind, kann abgeschätzt werden, dass wahrscheinlich mehr als die Hälfte aller Proteine glykosyliert ist. Die Zuckerstrukturen sind entweder über ein Stickstoffatom der Aminosäure Asparagin (N-Glykane) oder über das Sauerstoffatom von Serin oder Threonin (O-Glykane) mit dem Protein verknüpft. Während für die N-Glykane eine Erkennungssequenz auf Proteinebene existiert (Asn-X-Ser/Thr, wobei X ungleich Prolin ist), kann praktisch jedes auf der Proteinoberfläche exponierte Serin oder Threonin glykosyliert werden. Die Größe der verknüpften Zuckerstrukturen reicht von kleinen, linearen Di- und Trisacchariden (O-Glykane) bis zu komplexen, mehrfach verzweigten Strukturen (N-Glykane) mit bis zu vierzig monomeren Einheiten.

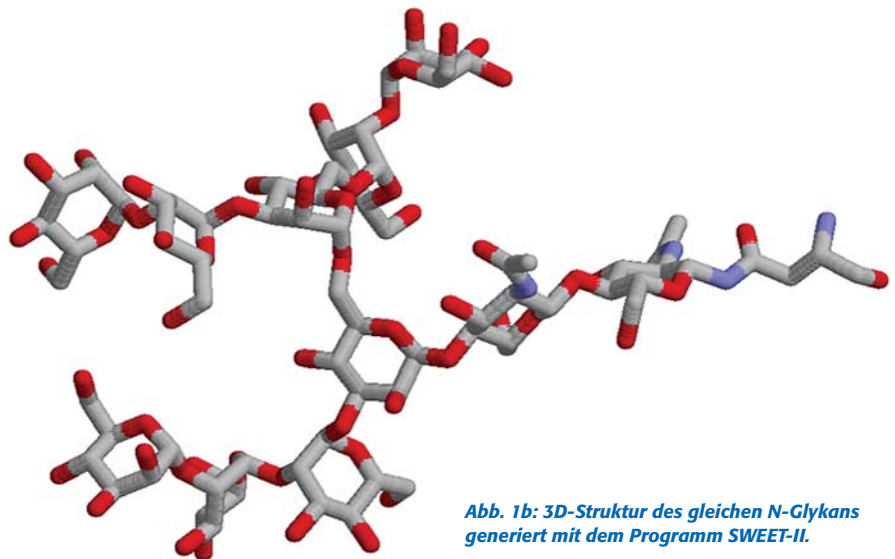


Abb. 1b: 3D-Struktur des gleichen N-Glykans generiert mit dem Programm SWEET-II.

Glykosylierungsmuster als empfindliche Marker

► Die spezifischen Glykanstrukturen, die sich an einer bestimmten Glykosylierungsstelle des Proteins ausbilden, werden nur indirekt durch das Genom bestimmt. Diese Arbeit verrichten Enzyme, die unter dem Begriff Glykosyltransferasen (GT) zusammengefasst werden. Sie bauen in einer sehr spezifischen Synthese – ein Enzym katalysiert nur die Bildung einer speziellen glykosidischen Verknüpfung – N- und O-Glykane auf. Verglichen mit der Biosynthese von Proteinen gibt es also einen zusätzlichen Schritt der Verschlüsselung. Die in einer Zelle vorhandenen GT stellen die Anweisungen für eine molekulare Maschinerie dar, durch welche das Repertoire von Glykan-

strukturen bestimmt wird, das einer Zelle zugänglich ist. Es sind mittlerweile einige hundert GT bekannt.

Obwohl die gleichen Glykosylierungswerkzeuge allen Kopien eines Proteins zur Verfügung stehen, die einen bestimmten sekretorischen Pfad einer Zelle durchlaufen, sind viele Proteine nicht einheitlich mit bestimmten Glykanen dekoriert, sondern zeigen ein Muster, das für eine spezifische Glykosylierungsstelle charakteristisch ist^[2]. Aus diesem Grunde ist das Glykosylierungsmuster ein sehr empfindlicher Marker für Veränderungen in den Zellen.

Das ‚glycan profiling‘ – also die Bestimmung aller Glykan-Strukturen, die in einer Zelle oder einem Organ vorkommen (Glykan-Repertoire) – kann unter anderem dazu dienen, krankes von normalem Gewebe zu

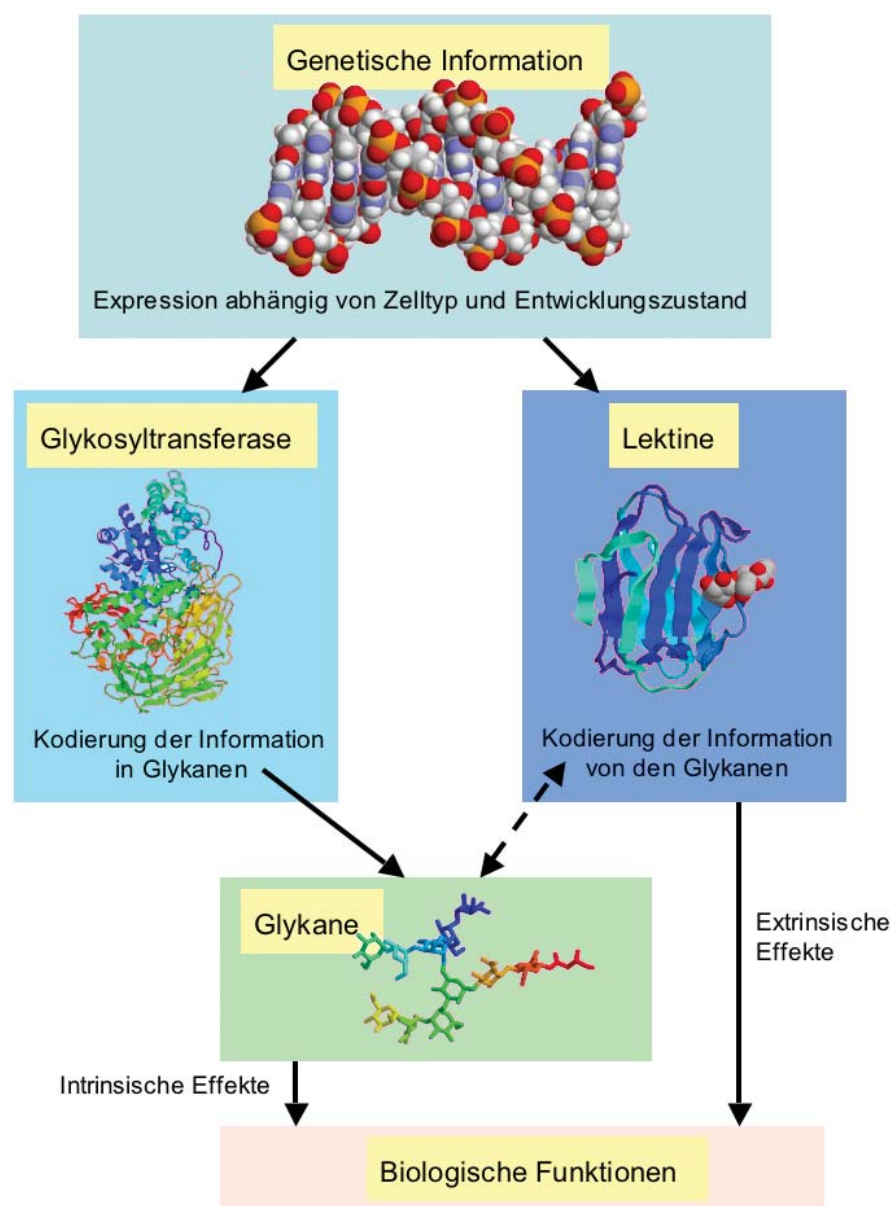


Abb. 2: Kodierung und Dekodierung der biologischen Information von Glykan-Strukturen. Der Aufbau von Glykanstrukturen wird nur indirekt durch das Genom bestimmt. Diese Arbeit verrichten Glykosyltransferasen. Sie bauen in einer sehr spezifischen Synthese – ein Enzym katalysiert nur die Bildung einer speziellen glykosidischen Verknüpfung – N- und O-Glykane auf. Verglichen mit der Biosynthese von Proteinen gibt es also einen zusätzlichen Schritt der Verschlüsselung. Beeinflussen die Glykane direkt die Eigenschaften der Proteine wie z.B. ihre Löslichkeit und Stabilität so wird dies als intrinsische Funktion der Glykane bezeichnet. Als extrinsische Effekte werden Funktionen von Glykanen bezeichnet, die aus einer spezifischen Wechselwirkung (z.B. molekulare Erkennungen und Signalübertragungen) mit den Zucker-erkennenden Proteinen (Lektine) resultieren.

unterscheiden. So zeigen sich im krankhaften Gewebe bei der rheumatoiden Arthritis und dem Rinderwahnsinn deutliche Unterschiede im Glykosylierungsmuster bestimmter Proteine^[3], die darauf hindeuten, dass zelluläre oder genetische Veränderungen die Aktivität der GT beeinflussen^[4]. Eine weitere wichtige Technik, mit der man sowohl die Synthese von Glykanen als auch ihre Funktion untersuchen kann, sind Experimente mit ‚Knock-out‘ (KO)-Mäusen, bei denen die Gene für bestimmte GT entfernt wurden. Die genaue Charakterisierung

der strukturellen Unterschiede in den gefundenen Glykanmustern von normalen und KO-Mäusen ermöglichen detaillierte Einblicke in die Rolle der GT einer bestimmten Zelle eines Organs.

Analytische Strategien

Weltweit wird in verschiedenen Laboratorien intensiv daran gearbeitet, geeignete analytische Strategien auszutesten, die eine schnelle Kartierung des Glykan-Repertoires eines Proteins ermöglichen. Ähnlich wie im

Bereich der Proteomik ist die Massenspektrometrie (MS)^[5] in Kombination mit verschiedenen enzymatischen Abbaureaktionen besonders gut geeignet für ein schnelles und automatisches Screening von Glykomen. Andererseits ist mittels verschiedener NMR-Techniken^[6] eine vollständige Strukturbestimmung – Art und Konfiguration der monomeren Einheiten, Art der glykosidischen Verknüpfung, Konformation der Glykane – zugänglich. Allerdings liegt die benötigte Menge an reiner Substanz um Größenordnungen höher als bei den MS-Verfahren, sodass NMR-Techniken weniger gut geeignet sind für Hochdurchsatzverfahren. Generell ist die Bestimmung der Struktur von Glykanen aufgrund der unterschiedlichen Verknüpfungsmöglichkeiten der monomeren Einheiten (siehe *Kasten strukturelle Vielfalt von Zuckern*) wesentlich aufwändiger als die Analyse von DNA- und Proteinsequenzen.

Angesichts der zu erwartenden Datenmengen, die im Rahmen sich abzeichnender Glykomiik-Projekte anfallen werden, ist es dringend notwendig, dass – analog zu den Anstrengungen im Bereich der Genomik und Proteomik – auch entsprechende bioinformatische Werkzeuge und Datenbanken für die Glykobiologie entwickelt werden. Aufgrund der Unterschiede in ihren strukturellen Merkmalen – lineare, immer gleichartige Verknüpfung der monomeren Bausteine einerseits und verzweigte Strukturen mit unterschiedlichen Verknüpfungen andererseits – lassen sich die zur Beschreibung von Ähnlichkeiten und Eigenschaften von DNA- und Proteinsequenzen entwickelten Algorithmen nicht einfach übertragen.

Bioinformatik für die Glykobiologie

Eine zentrale, allgemein zugängliche Anlaufstelle zum Abrufen von Kohlenhydratstrukturen, wie etwa das *European Bioinformatics Institute (EBI)* oder das US-amerikanische National Center for Biotechnology Information (NCBI) für Gen- und Proteinsequenzen, gibt es bisher nicht. Der *Tabelle 1* ist zu entnehmen, dass etliche im Internet verfügbare Datenbanken und Anwendungen existieren. Diese werden aber unter verschiedenen thematischen Aspekten und mit unterschiedlichen Intentionen aufgebaut und betreut. Während einige der angebotenen Dienste lediglich über das Netz verfügbar gemachte Ergebnisse von Diplom- und/oder Doktorarbeiten darstellen und nicht sicher gestellt ist, dass diese Dienste auch in Zukunft gepflegt werden, zeigt sich bei anderen Anwendungen, dass kontinuierlich neue Daten eingegeben werden und der Funktionsumfang der angebotenen

Strukturelle Vielfalt von Zuckern

Wie alle biologischen Makromoleküle setzen sich auch die Kohlenhydrate aus einer begrenzten Anzahl von monomeren Untereinheiten – den Monosacchariden, die zu meist als ringförmige Verbindungen vorliegen – zusammen. Diese werden durch eine Dreibuchstaben-Kodierung – z.B. Gal für Galactose, Glc für Glucose – beschrieben. Anders jedoch als bei den Proteinen, bei denen die Aminosäuren immer über die gleiche Art miteinander verknüpft sind und ein kettenförmiges Molekül bilden, können Glykane über verschiedene Atome verbunden sein. Deshalb ist auch die Angabe der verknüpften Atome für jede glykosidische Verbindung notwendig. Es kön-

nen sich verzweigte oder gar ringförmige Strukturen ausbilden. Bisher wurden – wenn man die Pflanzenwelt mit einbezieht – mehr als 100 verschiedene monomere Untereinheiten gefunden. Als zusätzliche Information für eine vollständige Beschreibung wird die Konfiguration am anomeren C-Atom benötigt, die mit α oder β gekennzeichnet wird. Die große strukturelle Vielfalt von Glykanen verglichen mit den DNA- und Proteinsequenzen, liefert eine einleuchtende Erklärung dafür, dass die Analytik von Zuckern aufwendiger ist, und ihre digitale Beschreibung zusätzliche Terme berücksichtigen muss.

Im Gegensatz zu den Gen- und Proteindatenbanken, bei denen man oft die interessierenden Sequenzen selbst eingeben muss, bekommt der Benutzer die LINUCS-Notation nie zu Gesicht. Die Eingabe der Glykan-Strukturen erfolgt in der für den Benutzer gewohnten Schreibweise (siehe *Abb. 1a*). Spezielle Routinen, die wir auch allen Entwicklern von Datenbanken zur Verfügung stellen, sorgen für eine Umsetzung in die interne, lineare Darstellung. Somit haben wir die inhaltlichen Voraussetzungen geschaffen, dass unsere Daten einfach mit anderen Datenbeständen verknüpft werden können.

Werkzeuge für die Massenspektrometrie

Ähnlich wie bei der Identifizierung von Proteinen in der Proteomik zeichnet sich ab, dass die MS auch in Glykomik-Projekten die analytische Methode der Wahl zur schnellen Bestimmung des Repertoires von Glykanen

Werkzeuge sich fortlaufend erweitert. Zusätzlich gibt es auch Datenbanken, die von kommerziellen Anbietern erstellt werden. Diese sind jedoch nur gegen Bezahlung zugänglich. Etablierte Prozeduren mittels derer Wissenschaftler die von ihnen erzeugten Primärdaten in einer zentralen Datenbank ablegen können – wie dies z.B. bei Gen- und Proteinsequenzen oder bei räumlichen Proteinstrukturen der Fall ist – gibt es gegenwärtig nicht.

Die Beschreibung von Gen- und Proteinsequenzen als Abfolge von Buchstaben stellt eine eindeutige Beschreibung ihrer Struktur dar, die sich gut digital verarbeiten lässt. Mittels der Sequenzinformation lassen sich Datenbestände effizient verknüpfen. Leider gibt es keine vergleichbare, vereinheitlichte digitale Beschreibung der Struktur von Glykanen, sodass eine schnelle Verknüpfung der weltweit vorhandenen Datensammlungen schwierig bis unmöglich ist. Gerade das einfache Abrufen aller zu einem Protein verfügbaren Daten durch die automatische Abfrage von Datenbanken mit unterschiedlichen Themenschwerpunkten ist eine der wesentlichen Voraussetzungen für die hohe Akzeptanz von Bioinformatik-Ansätzen in der molekular-biologischen und medizinischen Forschung.

Wir haben uns daher bemüht, auch für Kohlenhydrate eine eindeutige, lineare Beschreibung zu entwickeln, die sich gut digital verarbeiten lässt. Alle von der Zentralen Spektroskopie im DKFZ entwickelten Datenbanken und Anwendungen sind mittlerweile über die LINUCS-Notation (**L**inear **N**otation for **U**nique description of **C**arbohydrate **S**equences)^[7] miteinander verknüpft.

Für Glykobiologie im Web verfügbare Anwendungen und Datenbanken

Name	Inhalte	URL
Kohlenhydrat-relevante Informationen in Protein Datenbanken		
CAZy	Kohlenhydrat-aktive Enzyme	afmb.cnrs-mrs.fr/CAZY/
Lectines	3D Strukturen von Lektinen	www.cermav.cnrs.fr/lectines
CTDL	tierische Lektine	ctdl.glycob.ox.ac.uk/
PDB2LINUCS	Glykoproteine in der PDB	www.dkfz.de/spec/pdb2linucs/
Vorhersage von Glykosierungsstellen von Proteinen		
NetNGlyc	N-Glykosylierung	www.cbs.dtu.dk/services/NetNGlyc/
NetOGlyc	O-Glykosylierung	www.cbs.dtu.dk/services/NetOGlyc/
YinOYang	Glyko-, Phosphorylierung	www.cbs.dtu.dk/services/YinOYang/
big-PIPredictor	GPI-Anker Vorhersage	mendel.imp.univie.ac.at/sat/gpi/gpi_server.html
DGPI	GPI-Anker Vorhersage	129.194.185.165/dgpi/index_en.html
Werkzeuge zur Glykan-Strukturbestimmung		
Glycofragment	Masse von Glykan-Fragmenten	www.dkfz.de/spec/projekte/fragments/
GlycoSearchMS	MS-Spektrenvergleich	www.dkfz.de/spec/sweetdb/
GlycoMod	Glykanstruktur aus Molpeak	www.expsy.org/tools/glycomod/
GlycoMass	Masse von Kompositionen	www.expsy.org/tools/glycomod/glycanmass.html
GlyPeps	Glykoprotein Erkennung	www.dkfz.de/spec/glypeps/
CASPER	¹ H, ¹³ C-NMR Abschätzung	www.casper.organ.su.se/casper/
SugarBase	¹ H, ¹³ C-NMR Suche	boc.chem.uu.nl/sugabase/sugabase.html
NMR-Serach	¹ H, ¹³ C-NMR Suche	www.dkfz.de/spec/sweetdb/
Weitere nützliche Werkzeuge		
LINUCS	lineare Codierung	www.dkfz.de/spec/linucs/
LiGraph	Graphische Darstellung	www.dkfz.de/spec/ligraph/
IUPAC	Nomenklatur	www.chem.qmw.ac.uk/iupac/2carb/
Räumliche Strukturen		
SWEET-II	Generierung 3D-Strukt.	www.dkfz.de/spec/sweet2/
Disaccharides	Konformationskarten	www.cermav.cnrs.fr/cgi-bin/dj/di.cgi
GlycoMaps DB	Konformationskarten	www.dkfz.de/spec/glycomaps/
Dynamic Molecules	Molekular-Dynamik Simulationen	www.md-simulations.de
Kohlenhydrat-Datenbanken		
Name	Anbieter	URL
SWEET-DB	DKFZ-Heidelberg	www.dkfz.de/spec/sweetdb/
Carbohydrate DB	Consortium for Functional Glycomics	web.mit.edu/glycomics/carb/carbdb.shtml
GlycoSuite	Proteome Systems Ltd	www.glycosuite.com/
Glycomic Database	GlycoMinds	www.glycominds.com/GlycoInfo.asp.

ist. Der automatische Vergleich der MS-Spektren von enzymatisch verdauten Proteinen mit den theoretisch berechneten Massenspektren der entsprechenden Proteinfragmente erlaubt in vielen Fällen eine eindeutige Identifizierung der untersuchten Proteine. Da bisher vergleichbare Verfahren und Algorithmen für den Glykomik-Bereich nicht verfügbar waren, haben wir damit begonnen, die entsprechenden Werkzeuge zu entwickeln. Das Programm **GlycoFragment**^[8] berechnet die Massen aller theoretisch möglichen Bruchstücke eines Glykans. Bei der Berechnung dieser Bruchstücke kann auf empirische Regeln aufgebaut werden, welche Bindungen bevorzugt gebrochen werden. Ausgehend von den Einträgen, die in der **SWEET-DB**^[9] enthalten sind, haben wir für den MS-Vergleich etwa 10.000 theoretische Spektren von N- und O-Glykanen sowie Glykolipiden generiert. Vergleicht man nun – ähnlich wie bei den Proteinen – das gemessene MS-Spektrum mit allen theoretisch berechneten Spektren, so kann eine eindeutige Identifizierung von schon bekannten Strukturen erfolgen. Über die Anzahl der Treffer innerhalb einer gewissen Fehlertoleranz und der Abweichung von berechneten und gemessenen Massen, berechnet das Programm **GlycoSearchMS** einen Bewertungsfaktor, der ein Maß für die Güte der Übereinstimmung ist. Je höher dieser Bewertungsfaktor ist, umso höher ist die Wahrscheinlichkeit, dass die gefundenen Substanzen auch tatsächlich korrekt identifiziert sind.

MS-Methoden sind nicht für eine vollständige Strukturauflösung von Glykanen geeignet, da sie z.B. nicht erlauben, zwischen Stereoisomeren wie Galactose und Glucose zu unterscheiden, da diese identische Massen haben. Eine eindeutige Identifizierung erlauben verschiedene chromatographische Trennverfahren, die jedoch alle den Nachteil haben, dass sie mit bekannten Oligosacchariden kalibriert und validiert werden müssen. Daher ist die NMR-Spektroskopie die wichtigste physikalische Methode zur Bestimmung der Konformation eines Glykans. Sie erlaubt eine vollständige Strukturbestimmung der Glykane, benötigt aber eine größere Menge an zu untersuchender Substanz, erfordert eine größere Reinheit und längere Messzeiten.

Werkzeuge für NMR

Auch für die NMR-Spektren sind Auswertungsverfahren zur automatischen Erkennung von Glykanen in der Entwicklung. Allerdings sind die chemischen Verschiebungen von Verbindungen, die zum Vergleich herangezogen werden, keine physikalischen Größen, die sich mit hinreichender Genau-

igkeit einfach aus der chemischen Struktur ableiten lassen. Nach wie vor stellen Datenbanken, die experimentell gemessene NMR-Spektren von Referenzsubstanzen enthalten, die verlässlichste Grundlage für entsprechende Vergleichs- und Abschätzungsalgorithmen dar. Leider ist es bisher nicht gelungen, eine Datenbank aufzubauen, in der alle interessierten Wissenschaftler ihre NMR-Spektren ablegen können und die von der wissenschaftlichen Gemeinschaft allgemein akzeptiert wird. Da die wissenschaftlichen Zeitschriften immer mehr dazu neigen, die Spektren nicht mehr in den Publikationen zu dokumentieren, ist die Gefahr hoch, dass diese mühsam gewonnenen Daten für die Wissenschaft verloren gehen.

Unsere Datenbank **SWEET-DB** enthält die ¹H- und ¹³C-NMR Spektren von etwa 1000 Oligosacchariden, die aus der Literatur extrahiert wurden. Verschiedene Suchstrategien erlauben das gezielte Auffinden von NMR-Spektren aufgrund struktureller Vorgaben – z.B. bestimmten Teilstrukturen – oder die Zuordnung von Strukturen zu bestimmten Spektren. Auch wenn somit ein praktikabler Weg aufgezeichnet ist, wie NMR-Spektren für die automatische Analyse von Glykanen genutzt werden können, so reicht die momentan verfügbare Menge an Referenzdaten nicht aus, um den Anforderungen der Glykomik-Projekte zu genügen. Aus diesem Grunde haben wir Werkzeuge entwickelt, die es NMR-Spektroskopikern ermöglichen, ihre Daten zunächst in einer lokal installierten Datenbank einzugeben, um sie dann nach Veröffentlichung der Allgemeinheit zur Verfügung zu stellen. Unsere Hoffnung ist, dass ein solches Konzept greift und akzeptiert wird und somit dem Verlust von wissenschaftlichen Daten entgegen wirkt.

Spitzentechnologie „Glycomiks“

Die Entwicklung von Datenbanken und Bioinformatik-Werkzeugen für die Glykobiologie befindet sich noch in Kinderschuhen^[10]. Es ist jedoch offensichtlich, dass die jetzt entstehenden Glykomik-Projekte eine starke Nachfrage erzeugen werden. Bei dem größten bisher laufenden Projekt – dem amerikanischen *Consortium for Functional Glykomik* – ist die Entwicklung einer Bioinformatik-Infrastruktur ein ganz zentraler Punkt. Momentan wird intensiv daran gearbeitet, drei neue Datenbanken für Glykosyltransferasen, Lektine und Kohlenhydrate aufzubauen. In einer Studie der Zeitschrift *Technology Review* des *Massachusetts Institute of Technology*, die im Februar 2003 veröffentlicht wurde, wird die ‚Glykomik‘ als eine der zehn Spitzentechnologien genannt, „die die Zukunft verändern wird“. Zieht man die übliche Übertreibung in solchen Äuße-

rungen ab, so ist dennoch klar, dass es gute wissenschaftliche Gründe gibt, der Bedeutung von Zuckern als biologisch wichtigen Akteuren des Organismus einen größeren Raum einzuräumen.

Literatur

- [1] **Taylor M.E., Drickamer, K.** (2003): Introduction to Glycobiology, Oxford University Press
- [2] **Ritchie G.E., Moffatt B.E., Sim R.B., Morgan B.P., Dwek R.A., Rudd P.M.** Glycosylation and the complement system. *Chem Rev.* 2002 Feb; 102(2):305–20–19.
- [3] **Rudd P.M., Merry A.H., Wormald M.R., Dwek R.A.** (2002): Glycosylation and prion protein. *Curr Opin Struct Biol.* 12(5):578–86
- [4] **Peracaula R., Royle L., Tabares G., Mallorqui-Fernandez G., Barrabes S., Harvey D.J., Dwek R.A., Rudd P.M., De Llorens R.** (2003): Links Glycosylation of human pancreatic ribonuclease: differences between normal and tumor states. *Glycobiology.* 13(4):227–44
- [5] **Dell A., Morris H.R.** (2001): Glycoprotein structure determination by mass spectrometry. *Science* 291(5512):2351–6
- [6] **Manzi A.E., Norgard-Sumnicht K., Argade S., Marth J.D., van Halbeek H., Varki A.** (2000): Exploring the glycan repertoire of genetically modified mice by isolation and profiling of the major glycan classes and nano-NMR analysis of glycan mixtures. *Glycobiology.* 10(7):669–89
- [7] **Bohne A., Lang E., Forster T., von der Lieth C.-W.** (2001): LINUCS: LInear Notation for Unique Description of Carbohydrate Sequences. *Carbohydrate Research* 336 1–11
- [8] **Lohmann K.K., von der Lieth C.-W.** (2003): Glyco-Fragment: a Web Tool to Support the Interpretation of Mass Spectra of Complex Carbohydrates. *Proteomics, in press*
- [9] **Loss A., Bunsmann P., Bohne A., Loss A., Schwarzer E., Lang E., von der Lieth C.-W.** (2002): SWEET-DB: an attempt to create annotated data collections for carbohydrates. *Nucleic Acids Res.* 1; 30(1):405–8
- [10] **Marchal I., Golfier G., Dugas O., Majed M.** (2003): Bioinformatics in glycobiology. *Biochimie.* 85(1–2):75–81

Korrespondenzadresse:

Dr. Claus-W. von der Lieth
Deutsches Krebsforschungszentrum (DKFZ)
Zentrale Spektroskopie – B090
Im Neuenheimer Feld 280
D-69120 Heidelberg
Tel.: 06221-424541
Fax: 06221-424554
w.vonderlieth@dkfz.de
a.bohne@dkfz.de
k.lohmann@dkfz.de
www.dkfz.de/spec/