

A Probabilistic Model for Mining Labeled Ordered Trees: Capturing Patterns in Carbohydrate Sugar Chains

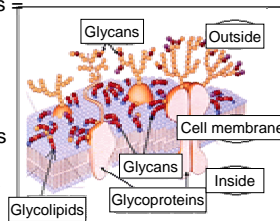
Hiroshi Mamitsuka
Bioinformatics Center
Kyoto University

Outline

- What is a Carbohydrate Sugar Chain (Glycan)?
- Glycan = Labeled Ordered Tree
- Databases on Glycans
- Probabilistic Models for Labeled Ordered Trees and their empirical experimental results
 - Probabilistic Sibling Dependent Tree Markov Model (PSTMM)
 - Profile PSTMM
 - Ordered Tree Markov Model
- Concluding Remarks

Glycans

- Carbohydrate Sugar Chains = **Glycans**
- Third major class of biomolecules next to DNA and proteins
- Often found on cell surfaces
- Crucial to the development and function of multicellular organisms
- However, many still unknowns in glycobiology



Glycans: Third Major Class

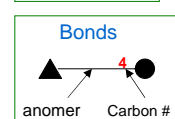
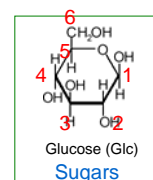
- DNA: Genome - Genomics
- Protein: Proteome - Proteomics
- Glycan: Glycome - Glycomics
 - The collective identity of the entirety of carbohydrates in an organism
 - The collective identity of the entirety of carbohydrates in a cell

Glycan Structure

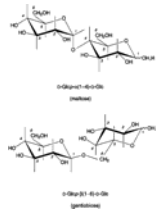
- DNA: String of four letters
 - four kinds of nucleotides (A,G,C,T)
- Protein: String of twenty letters
 - twenty types of amino acids
- Glycan: Tree structure of *monosaccharides* (sugars) and *linkages*
 - More than 10,000 structures known
 - ~twelve main types of monosaccharides (i.e., Glucose [Glc], N-acetylglucosamine [GlcNAc], Mannose [Man])
 - 10-15 classes (i.e., N-Glycans, O-Glycans, GPI anchors, etc.)

Building Blocks of Glycans

- | | |
|-----------|---------------------------|
| ● Galp | Galactose |
| ■ GalpNAc | N-acetylgalactosamine |
| ● Glcp | Glucose |
| ■ GlcpNAc | N-acetylglucosamine |
| ● Manp | Mannose |
| ▲ Fucp | Fucose |
| ☆ Xylp | Xylose |
| ◆ NeupAc | N-acetylneuraminic acid |
| ◇ NeupGc | N-glycolylneuraminic acid |
| ◆ KDN | Ketodeoxynonulosonic acid |



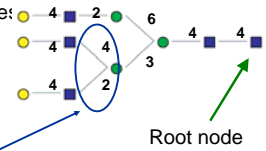
What Do Glycans Look Like?



- Variations in conformations between just two sugars (monosaccharides)
- Typical glycan structures contain 10-15 sugars!

What Do Glycans Look Like?

- IUPAC 2D Representation
- Tree structures of monosaccharides and linkage:
- Nodes = sugars/monosaccharides
- Edges = bonds/linkages
- Features:
 1. Rooted tree
 2. Monosaccharides = Labels
 3. Ordered children



Glycan Structure

- Glycan is
 - Rooted Tree:
 - Tree with root
 - (Rooted) Labeled Tree:
 - Tree whose nodes have labels attached
 - Monosaccharide names
 - (Rooted) Labeled Ordered Tree
 - Labeled tree whose children are ordered

Outline


- What is a Glycan?
- Glycan = Labeled Ordered Tree
- **Databases on Glycans**
- Probabilistic Models for Labeled Ordered Trees and their empirical experimental results
 - Probabilistic Sibling Dependent Markov Model (PSTMM)
 - Profile PSTMM
 - Ordered Tree Markov Model
- Concluding Remarks

General Database Systems for Glycans

- CarbBank
- SWEET-DB / glycosciences.de
- KEGG GLYCAN
- Consortium for Functional Glycomics
- EuroCarbDB
- Commercial databases:
 - GlycoSuite (Proteome Systems, Ltd.)
 - Glycomics DB (Glycominds, Ltd.)

CarbBank

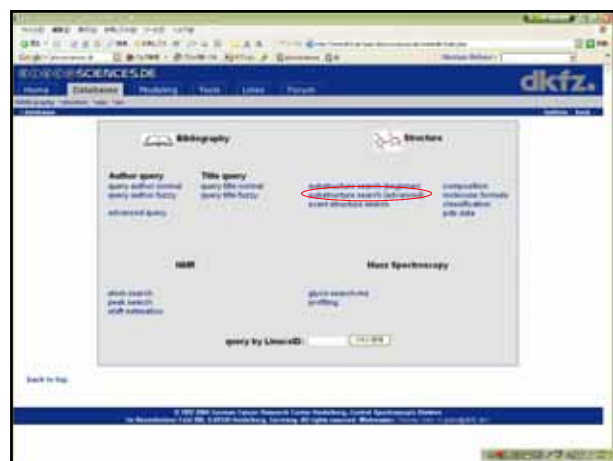
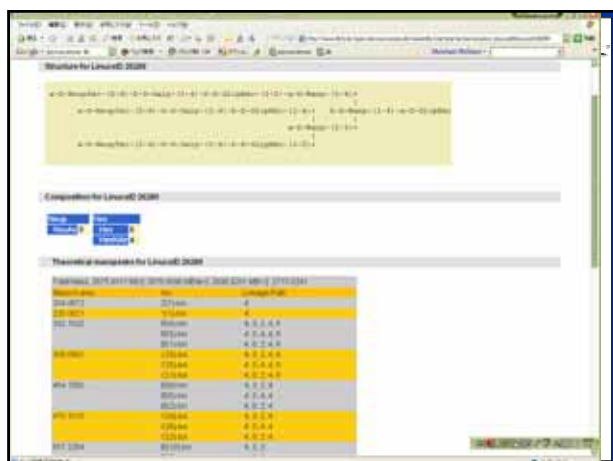
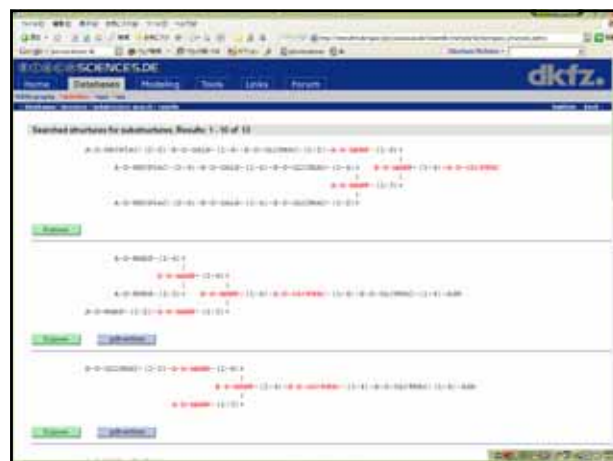
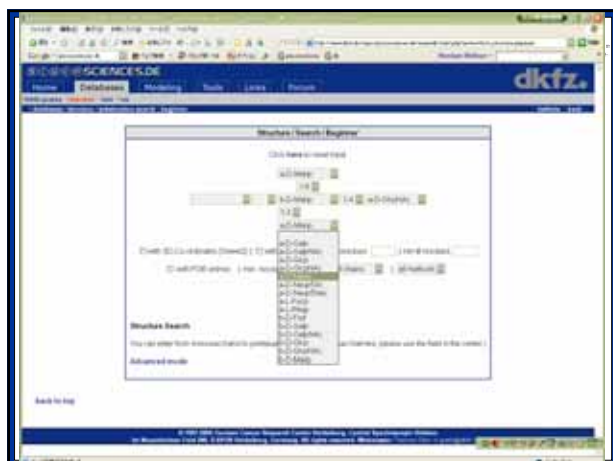
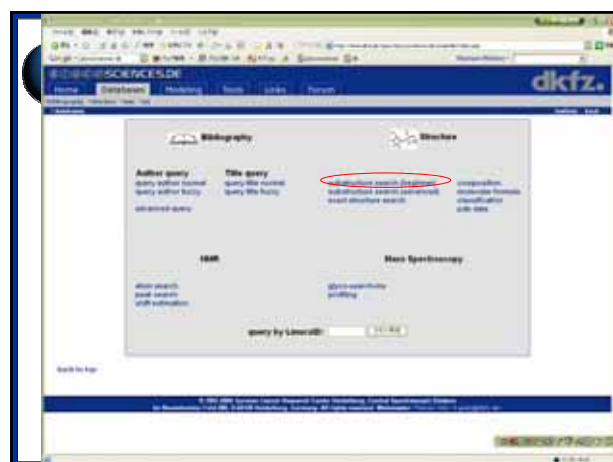
- Developed by Complex Carbohydrate Research Center, University of Georgia
- Community database of carbohydrates
- Project ended due to lack of funding in 1996
- Continued in Japan until around 2000



SWEET-DB

Bioinformatics Center
Institute for Chemical Research, Kyoto University
京都大学 化学研究所 バイオインフォマティクス研究センター

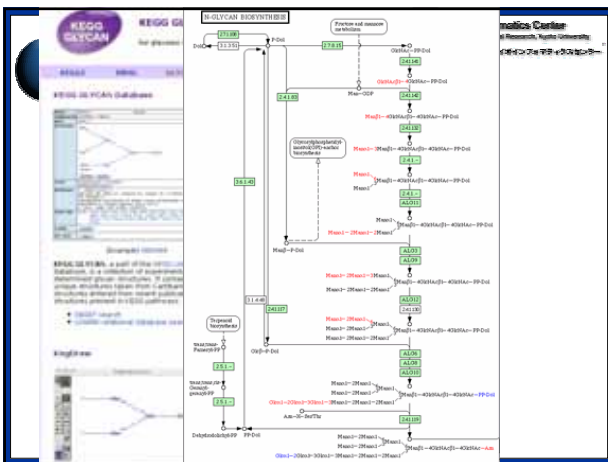
- A part of Glycoscience.de
- <http://www.dkfz-heidelberg.de/spec/sweetdb/>
- Combines CarbBank and Sugabase using a common web-based interface
- Provides searching by bibliography, structure, NMR and MS, as well as by LINUCS ID





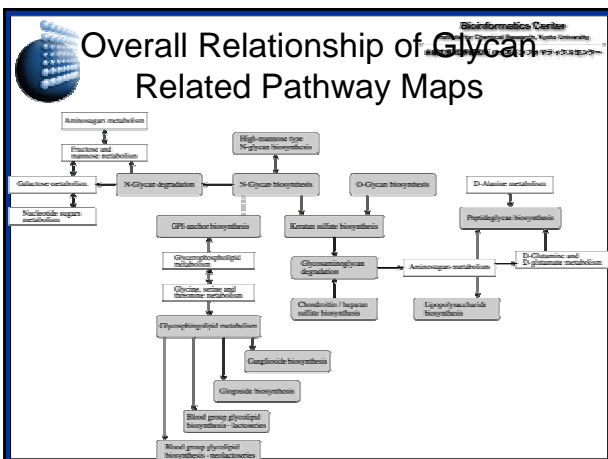
KEGG GLYCAN

- <http://www.genome.jp/kegg/glycan/>
- Based on CarbBank as well as input from scientists
- All data is linked with KEGG's other resources: GENES, PATHWAY, KO (KEGG Ontology) and literary databases



Glycan Biosynthesis and Metabolism Pathways

N-Glycan biosynthesis	Glycosylphosphatidylinositol(GPI)-anchor biosynthesis
High-mannose type N-glycan biosynthesis	Glycosphingolipid metabolism
N-Glycan degradation	Blood group glycolipid biosynthesis - lactoseries
O-Glycan biosynthesis	Blood group glycolipid biosynthesis - neo-lactoseries
Chondroitin / heparan sulfate biosynthesis	Globoside metabolism
Keratan sulfate biosynthesis	Ganglioside biosynthesis
Glycosaminoglycan degradation	Glycan structures - biosynthesis 1
Lipopolysaccharide biosynthesis	Glycan structures - biosynthesis 2
Peptidoglycan biosynthesis	Glycan structures - degradation



KEGG Glycan Search

Enter query glycan: (in one of the three forms)

Glycan ID: (Example: G00021) View structure

KCF File Name:

KCF File Text:

Select target database:

☒ KEGG GLYCAN ☐ CarbBank

Select program:

☒ Gapped (Approximate match) ☐ Ungapped (Exact match)

Select option:

☐ Global search ☐ Local search

KEGG Glycan Search

Query: G00078

Structure:

Database: ☒ KEGG GLYCAN ☐ CarbBank

Program: ☒ Gapped (Approximate match) ☐ Ungapped (Exact match)

Option: ☒ Global search ☐ Local search

-> Show advanced options

Glycan Data Search Result

Entry	Structure	Name	Composition	Class
G00078		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00079		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00080		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00081		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00082		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00083		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00084		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00085		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00086		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00087		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00088		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00089		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00090		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00091		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00092		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00093		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00094		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00095		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00096		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00097		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00098		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00099		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00100		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide

KEGG GLYCAN: G04450

Entry: G04450

Composition: (Gal)6 (GlcNAc)3 (Lac)2 (Cer)1

Mass: 1712.6 (Da)

Structure:

Class: Disaccharide

Other IDs: G04450

Links: All links

KCF data: [button]

Glycan Data Search Result

Entry	Structure	Name	Composition	Class
G00078		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00079		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00080		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00081		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00082		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00083		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00084		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00085		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00086		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00087		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00088		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00089		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00090		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00091		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00092		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00093		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00094		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00095		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00096		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00097		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00098		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00099		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide
G00100		GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	GalNAc1-4GlcNAc1-3GalNAc1-4GlcNAc1-1Cer	Disaccharide

Similarity-Score: 700

Query:

Entry: G04450

KEGG Glycan Search

Enter query glycan: (in one of the three forms)

Glycan ID: **G00078** (Example: G00021) [View structure](#)

KCF File Name:

KCF File Text:

Select target database: ☒ KEGG GLYCAN ☒ CarbBank

Select program: ☒ Gapped (Approximate match) ☐ Ungapped (Exact match)

Select option: ☒ Global search ☐ Local search

CarbBank Data Search Result

Number of entries in a page: 20

Page 1 of 42 Items: 1 - 20 of 833

No.	Entry	Similarity score	Title
1	CCSD-1313	800	On neutral fucoylglycosyls having long branch carbohydrate chains: H-active glycosphingolipids of human erythrocyte membranes
2	CCSD-1512	800	Isolation and characterization of gangliosides with a new sialosyl linkage and core structures. G ₂ Gangliosides of human erythrocyte membranes
3	CCSD-9241	800	A multiplicity of erythrocyte glycolipids of the neolacto series revealed by immuno-thin-layer chromatography with monoclonal anti-I and anti-x antibodies
4	CCSD-12348	800	Glycosylated antigens with blood-group I and x specificities from human adult and umbilical cord erythrocytes
5	CCSD-30390	800	Branched monosialo gangliosides of the lacto-series isolated from bovine erythrocytes: characterization of a novel ganglioside, Neu5Ac2Sia6GalNAc6S
6	CCSD-33213	800	Human neonatal gangliosides. Characterization of a novel I type ganglioside with the fucose alpha 2-6Gal structure
7	CCSD-33815	750	Blood group antigens on human erythrocytes: distribution, structure and possible functions
8	CCSD-1725	750	Characterization of an epitope (determinant) structure in a developmentally regulated glycolipid antigen defined by a cold agglutinin Fx. Recognition of alpha-1,6Galactose and alpha-4,6Galactose groups in a branched structure
9	CCSD-1730	750	On neutral fucoylglycosyls having long branch carbohydrate chains: H-active glycosphingolipids of human erythrocyte membranes
10	CCSD-1733	750	Characterization of blood-group I active gangliosides: structural requirements for I- and x-specificities
11	CCSD-1738	750	Isolation and characterization of an I-active Ceramide Hexacosahexide from Rabbit Erythrocyte Membrane
12	CCSD-1748	750	Human placenta gangliosides
13	CCSD-4580	750	Glycosphingolipids in cellular interaction, differentiation, and oncogenesis
14	CCSD-4583	750	Glycosphingolipids in cellular interaction, differentiation, and oncogenesis
15	CCSD-9089	750	A multiplicity of erythrocyte glycolipids of the neolacto series revealed by immuno-thin-layer chromatography with monoclonal anti-I and anti-x antibodies
16	CCSD-15890	750	Interaction of Myxobolus pneumoniae with erythrocyte glycolipids of I and x antigen types
17	CCSD-20753	750	Monoclonal antibodies directed to tumor-associated gangliosides and fucoylglycosides, method for production thereof, and use in passive immunization and diagnosis
18	CCSD-37111	750	Free specificity of a monoclonal anti-fetocellular cell antibody for glycolipids with terminal N-acetyl-D-glucosamine structure
19	CCSD-45257	750	Glycosylation in autoimmunity
20	CCSD-1728	650	Structural identification of two ten-sugar branched chain glycosphingolipids of blood-group H type present in epithelial cells of rat small intestine

Consortium for Functional Glycomics

- Consortium home page: <http://www.functionalglycomics.org/>
- Consortium of major universities and research institutes worldwide
- Aim: to provide a central resource for glycomics research
- Also provides requested resources to promote participating investigators' research
 - Glycan arrays and data
 - Mass spectra analysis...
- CFG glycan database web page: <http://www.functionalglycomics.org/glycomics/molecule/jsp/carbohydrate/carbMoleculeHome.jsp>

Glycan Database

Search by: Sub-Structure | MALDI | Composition | Linear Nomenclature | Multiple Criteria

Updates

- First full version of glycan structures database
- Contains nearly 7000 entries
- Each entry contains structural and chemical information as well as related references
- Different search interfaces are provided via the menu above
- The database will be regularly updated with newly synthesized or discovered glycans

Search for glycan

- Sub-structure
- Composition
- Linear nomenclature
- Use multiple search criteria

Search of glycan structures

- N- and O-linked glycans from Cellbank
- Glycosylated cell, seed database
- N- and O-linked glycans identified in tissues and cells analyzed by the Analytical Glycomics Database Core (AGDC)
- Glycans submitted to the Carbohydrate Structure Core (CSC) and available as a reference

Glycan nomenclature

- Glycans are displayed in various formats for ease of use
- The Consortium nomenclature for representing glycans can be found here

© 2002-2005 Consortium for Functional Glycomics. All rights reserved.

Glycan Search

Sub-Structure Search

Create a new structure to search for

Create structure starting with the template

Create structure starting with the template

[To find glycan structures from the database containing specific sub-structures.]

© 2002-2005 Consortium for Functional Glycomics. All rights reserved.

Glycan Search

Sub-Structure Search

Create a new structure to search for

Create structure starting with the template

Create structure starting with the template

[To find glycan structures from the database containing specific sub-structures.]

© 2002-2005 Consortium for Functional Glycomics. All rights reserved.

Glycan Search

Sub-Structure Search

Create a new structure to search for

Create structure starting with the template

Create structure starting with the template

[To find glycan structures from the database containing specific sub-structures.]

© 2002-2005 Consortium for Functional Glycomics. All rights reserved.

Motivation (from Biological Side)

- Glycans = Labeled Ordered Trees
- Many unknowns in glycobiology
 - High uncertainty and noisy
- The leaves of glycans are important in recognition by various pathogens
- Differences in these patterns affect biological functions
- Pattern mining method robust against noise required

Motivation (from Informatics Side)

- Labeled Ordered Trees: Semi-structured (or Unstructured) data
 - Other examples found in web and text mining, e.g. XML
- Mining semi-structured data, like graphs and/or trees, becoming important in machine learning and data mining
- Frequent pattern mining and kernels already developed recently
- New mining approaches robust against noise required

Probabilistic Modeling

- Statistical machine learning
- Represent uncertainty
- Robust against noise in data
- Efficient learning schemes already known
- Modeling labeled ordered trees not developed yet

Mining Labeled Ordered Trees Based-on Probabilistic Modeling

- To “learn” patterns from the tree structures of glycans by estimating the probability parameters of our model

Three Problems

- Must be solved to be used in real-world applications
 1. **Computing likelihood**: computing how likely a given example can be generated from a model
 2. **Learning**: estimating probability parameters of a model from given data
 3. **Parsing**: finding the most likely state transition on an example given a model

Three Problems

1. Computing the likelihood of a set of trees:
 - To determine which data are considered to belong to same class as training (learned) data
2. Estimating the parameters:
 - To “learn” patterns found in given data
3. Finding the most likely state transition:
 - To retrieve the learned patterns
 - To apply to multiple alignments

Bioinformatics Center
Institute for Chemical Research, Kyoto University
京都大学 化学研究所 バイオインフォマティクス研究センター

Three Problems

- Learning (estimating probability parameters) is the most important, since ...
- Computing the likelihood is a part of learning
- Parsing can be done by modifying the likelihood computation

Bioinformatics Center
Institute for Chemical Research, Kyoto University
京都大学 化学研究所 バイオインフォマティクス研究センター

Three Problems

- Must be solved to be used in real-world applications
- 1. **Computing likelihood**: computing how likely a given example can be generated from a model
- 2. **Learning**: estimating probability parameters of a model from given data
- 3. **Parsing**: finding the most likely state transition of a tree given a model

Bioinformatics Center
Institute for Chemical Research, Kyoto University
京都大学 化学研究所 バイオインフォマティクス研究センター

Hidden Markov Model (HMM)

- Model: set of states connected by directed edges
- Parameters:
 - State transition probability: a_{ij}
 - Label output probability (at state): $b_j(k)$
- Operation: series of Markov transition on states, generating sequence of labels (string) with likelihood

○ : state
← : state transition
↑ : Label output

Bioinformatics Center
Institute for Chemical Research, Kyoto University
京都大学 化学研究所 バイオインフォマティクス研究センター

Three Problems

- Hidden Markov model case:
 1. **Computing likelihood**: computing auxiliary probabilities: Forward or Backward probabilities
 2. **Learning**: Maximizing the likelihood by Baum-Welch (Forward-Backward) algorithm, an EM (Expectation-Maximization) algorithm
 3. **Parsing**: Viterbi algorithm

Bioinformatics Center
Institute for Chemical Research, Kyoto University
京都大学 化学研究所 バイオインフォマティクス研究センター

Forward Probability, HMM Case

- The probability that the current state is j and string $[1..t]$ is already generated: $\alpha_\sigma[t, j]$
- Dynamic programming over t available
- Updating formula: $\alpha_\sigma[t, i] = \sum_j a_{ij} b_j(\sigma_t) \alpha_\sigma[t-1, j]$

Bioinformatics Center
Institute for Chemical Research, Kyoto University
京都大学 化学研究所 バイオインフォマティクス研究センター

Computing Likelihood, HMM Case

- Compute forward probabilities over whole a given string
- Use final forward probabilities
- Likelihood: $\sum_i \alpha_\sigma[T, i]$

Learning, HMM Case

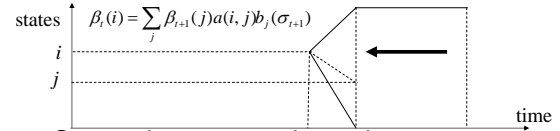
- Baum-Welch (Expectation-Maximization) algorithm
 1. Compute forward probabilities
 2. Compute backward probabilities: $\beta_\sigma[t, j]$
 - Probability that the current state is j and string $[t..T]$ is already generated
 3. Compute expectation value of transition ij :

$$E_{p_\sigma}[\#((i, j), \sigma)] \propto \sum \alpha_\sigma[t, i] a_{ij} b_j(\sigma_{t+1}) \beta_\sigma[t+1, j]$$
 4. Update transition probability a_{ij} using expectation values:

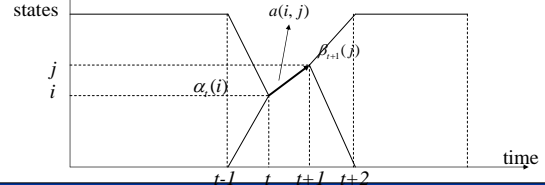
$$\hat{a}_{ij} = \frac{E_{p_\sigma}[\#((i, j), \sigma)]}{\sum E_{p_\sigma}[\#((i, j), \sigma)]}$$
- The above steps iterated until convergence

Supplement for HMM Learning

Computing backward probabilities



Computing expectation values



Learning, HMM Case Summary

- Maximizing the Likelihood
- Baum-Welch: EM algorithm
- Repeat the following two steps alternately until some stopping condition satisfied
 - E-step:
 1. Compute forward and backward
 2. Compute expectation values
 - M-step:
 1. Update transition probabilities

Back to Probabilistic Models for Trees

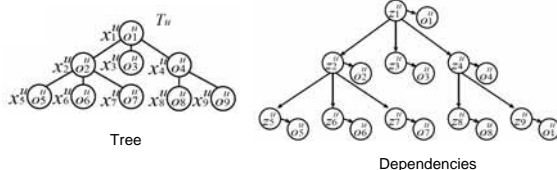
- Notations
 - Tree: T
 - Node: x
 - State type: s
 - State of node: z
 - Label: o

Related work:

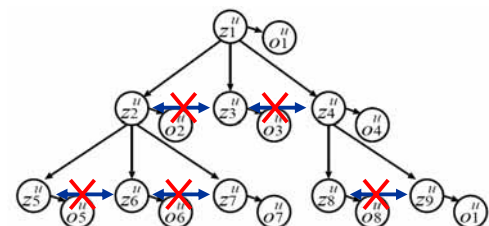
Hidden Tree Markov Model (HTMM)

[Deligenti et al., 2003]

- Probabilistic model for labeled trees, not for labeled ordered trees
- State depends on that of the parent only



HTMM Cannot Capture Sibling Dependencies!



Outline

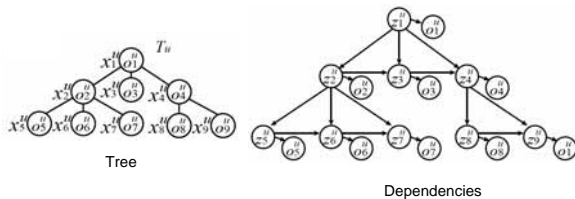
- What is a Glycan?
- Glycan = Labeled Ordered Tree
- Databases on Glycans
- Probabilistic models for Labeled Ordered Trees and their empirical experimental results
 - Probabilistic Sibling Dependent Markov Model (PSTMM)
 - Profile PSTMM
 - Ordered Tree Markov Model
- Concluding Remarks

Proposed model

- Probabilistic Sibling dependent Tree Markov Model (PSTMM)
 - Modeling sibling dependency as well as parent-child dependency
 - Extension of hidden Markov model (HMM) and hidden tree Markov model (HTMM)

PSTMM

- State depends on those of both the parent and the immediately elder sibling



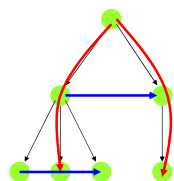
Define PSTMM Parameters

- Three probability parameters
 - Initial state probability: $\pi[s_l] (= P(z_l^u = s_l))$
 - Probability that the state of the root is s_l
 - State transition probability:

$$a[\{s_q, s_l\}, s_m] (= P(z_j^u = s_m \mid z_p^u = s_q, z_i^u = s_l))$$
 - Probability that the state of j is s_m given that the state of the immediately elder sibling is s_l and the state of the parent is s_q
 - Label output probability: $b[s_l, \sigma_h]$
 - Probability that the state s_l outputs σ_h

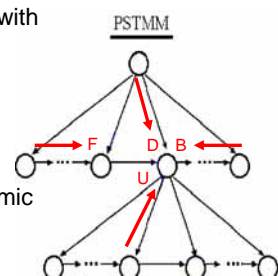
Auxiliary Probabilities for PSTMM

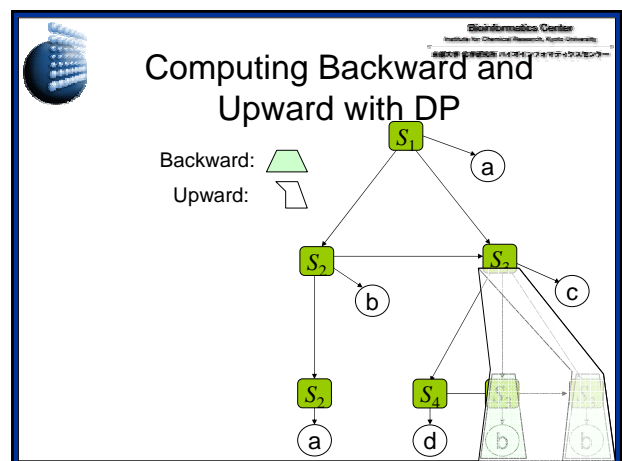
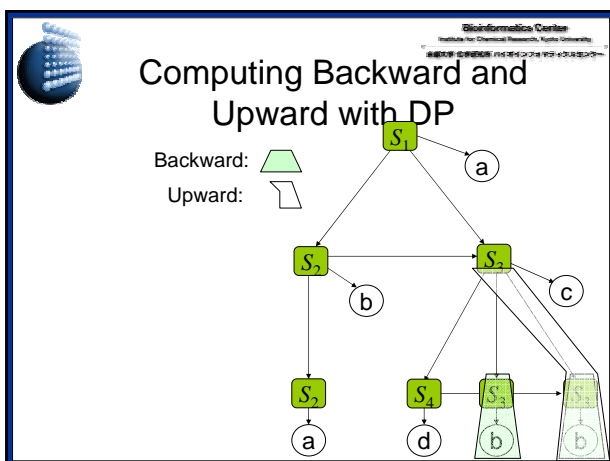
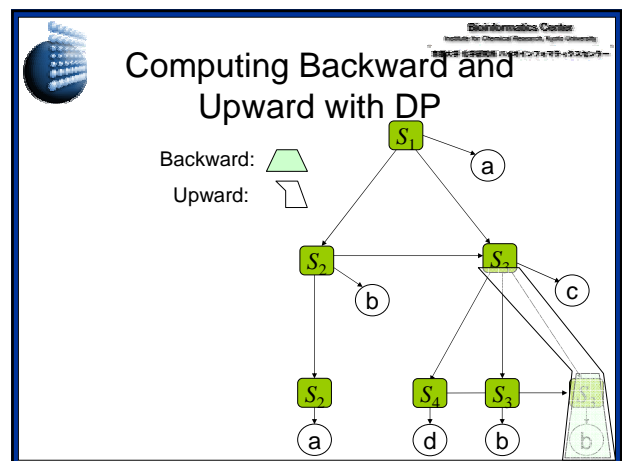
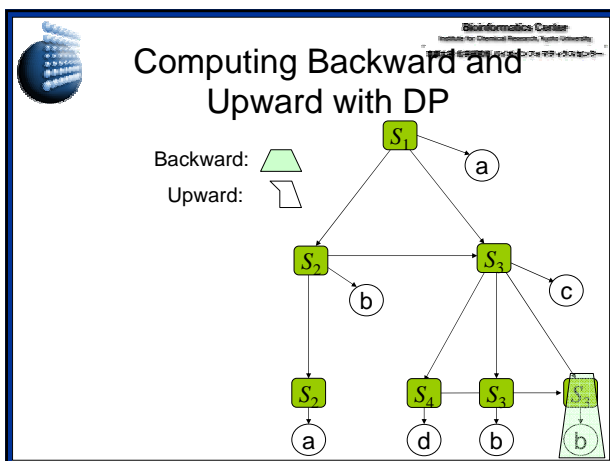
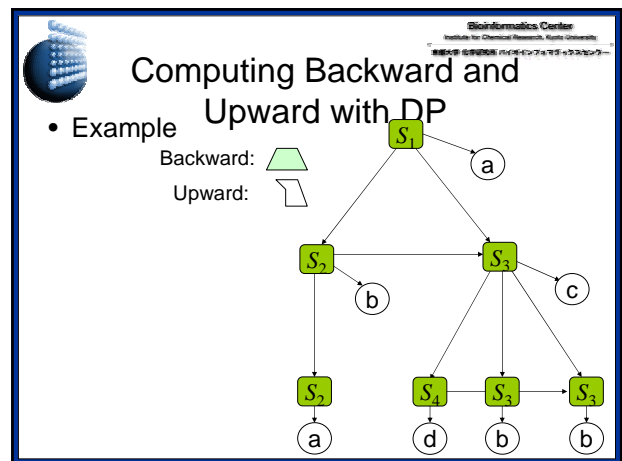
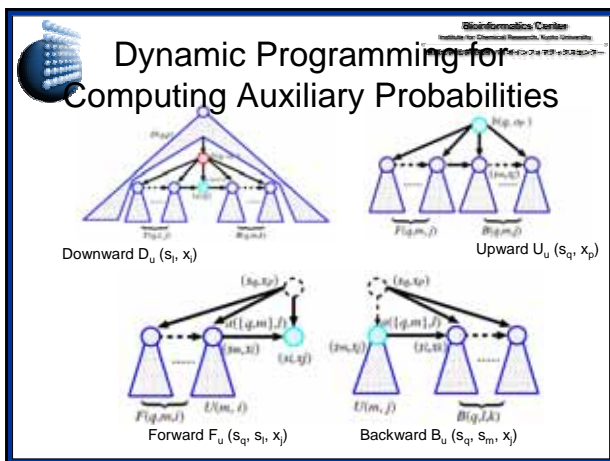
- Extension of HMM for Labeled Ordered Trees
 - Labeled Ordered Tree
 - Two directions:
 - Parent-Child directions
 - Sibling directions
 - Sibling: Forward and Backward
 - Parent-Child: ?
- ↓
- Downward and Upward!

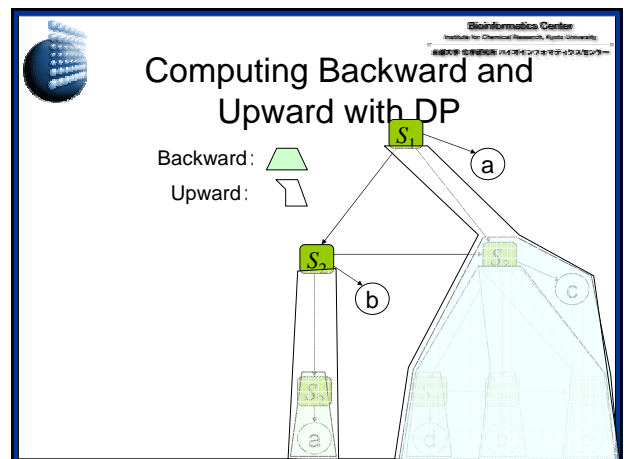
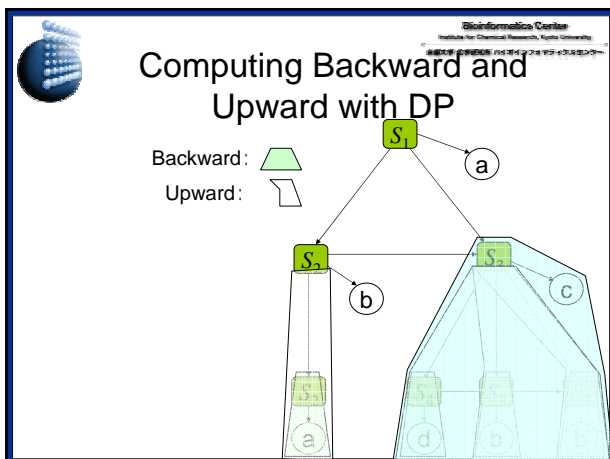
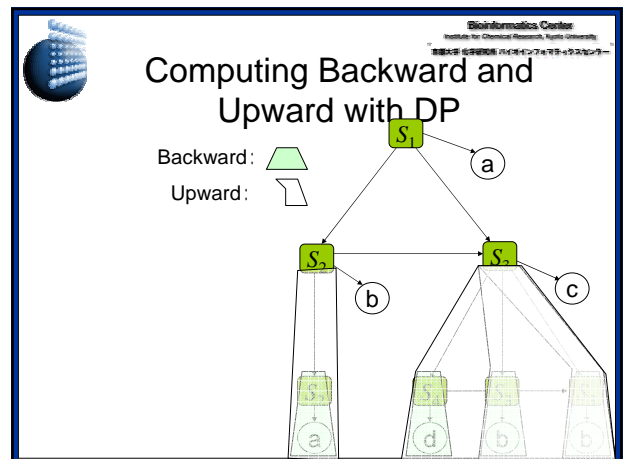
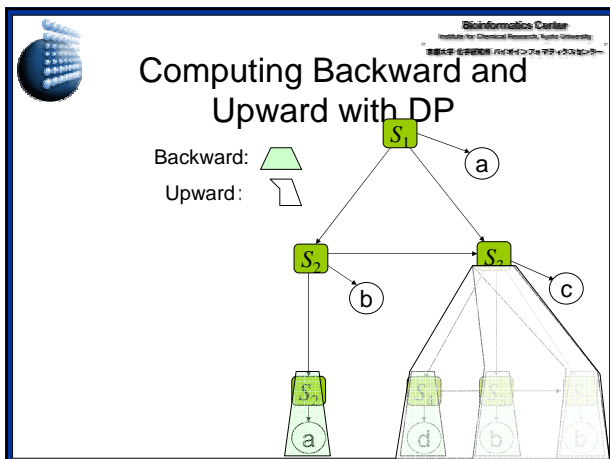
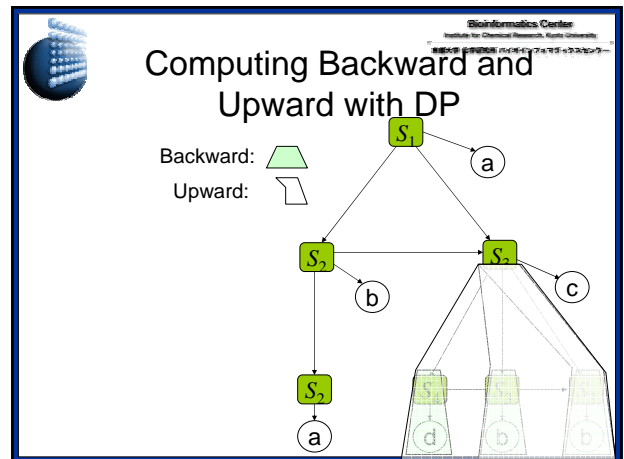
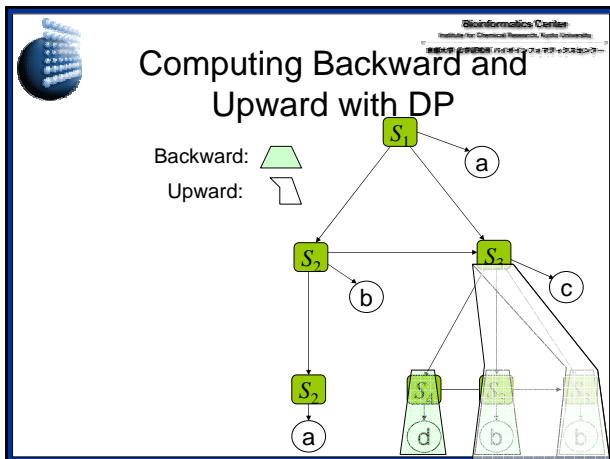


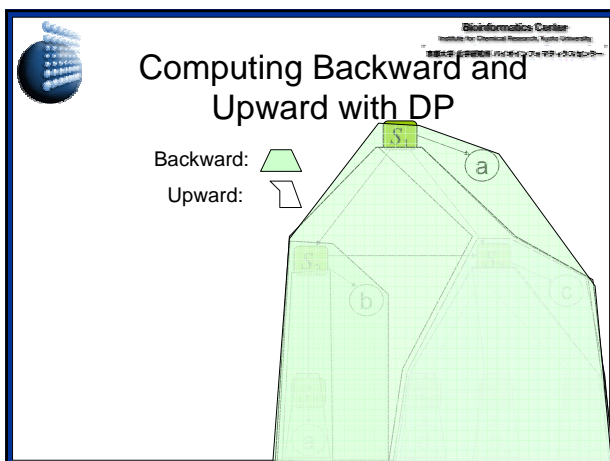
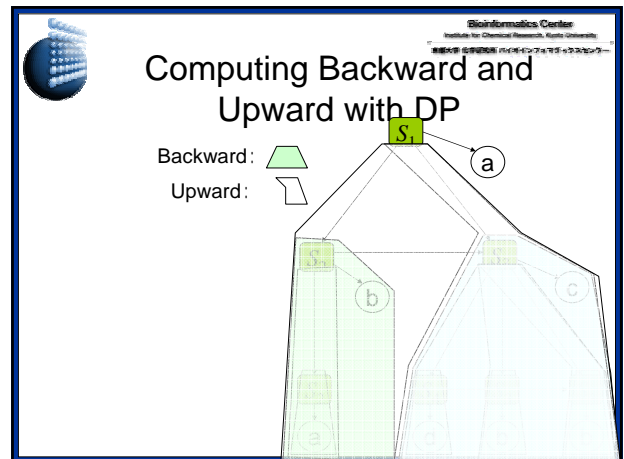
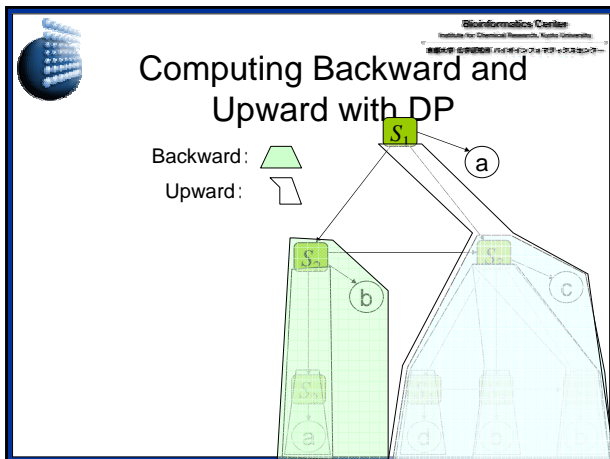
Four Auxiliary Parameters

- Define the following four probabilities for a tree u with nodes x_1, \dots, x_n
 - Forward $F_u(s_q, s_p, x_j)$
 - Backward $B_u(s_q, s_m, x_j)$
 - Upward $U_u(s_q, x_p)$
 - Downward $D_u(s_p, x_i)$
- Each computed by dynamic programming (DP)

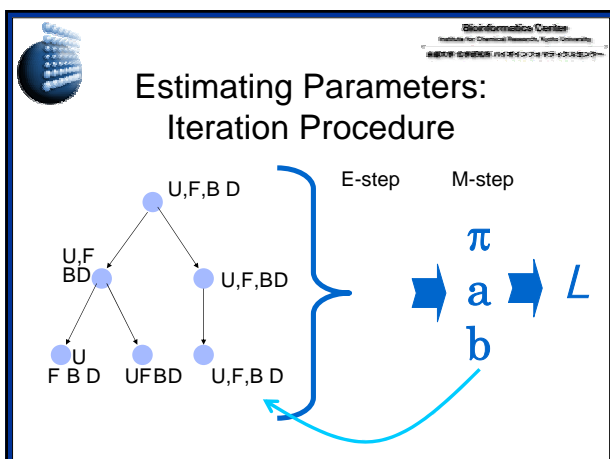








- Bioinformatics Center
Institute for Chemical Research, Kyoto University
- ### Estimating Parameters
- Maximum Likelihood and EM (Expectation-Maximization) as shown for HMMs
 - E-step computes three expectation values in E-step:
 - (s_m) = initial state expectation value
 - $\langle \{s_q, s_m\}, s_l \rangle$ = state transition expectation value
 - $(s_l, h) =$ label output expectation value
 - M-step updates our probability parameters $[s_l]$, $a[\{s_q, s_l\}, s_m]$, and $b[s_l, h]$ using these expectation values
 - Repeat E-M until likelihood is maximized



- Bioinformatics Center
Institute for Chemical Research, Kyoto University
- ### Computing the Likelihood
- The likelihood of tree T_u given a set of parameters can thus be found at the root node x_1 as:
 - $L(T_u;) = [s_l] U_u(s_l, x_1)$
 - The likelihood for a set of trees $T = \{T_u, \dots, T_n\}$ given a set of parameters can thus be computed as a product of the likelihood of each tree:
 - $L(T;) = \prod_u [s_l] U_u(s_l, x_1)$
 - $= \prod_u L(T_u;)$

Most Likely State Transition

- HMMs:
 - Viterbi algorithm
 - Used for multiple sequence alignment
- PSTMMs:
 - Viterbi algorithm
 - For Multiple tree alignment

The Most Likely State Transition

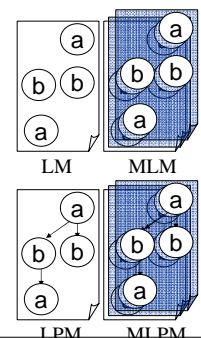
- Based on maximum likelihood parameters, we can calculate:
 - $\pi(s_q, s_m, x_i)$ = maximum state transition probability
 - $\mu(s_q, x_p)$ = maximum output label probability
 - $\pi(s_q, s_m, x_i)$ and $\mu(s_q, x_p)$, which retrieve the actual states corresponding to these maximum values (argmax of)

Empirical Experiments

- Experimental Setting
 - Synthetic
 - Performance comparison
 - Capture sibling-dependent patterns
 - Discriminate between those that do and do not contain these patterns
 - Use five-fold cross validation
 - Real data: glycans
 - Performance comparison
 - Analyzing patterns found

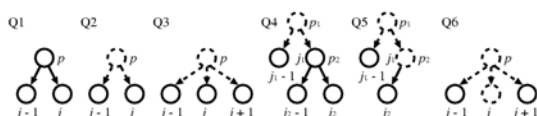
Synthetic Data Experiment

- Compared to five models that do not incorporate any dependencies among children.
 - HTMM
 - Label model (LM)
 - Mixture of label model (MLM)
 - Label pair model (LPM)
 - Mixture of label pair model (MLPM)



Synthetic Data Experiment

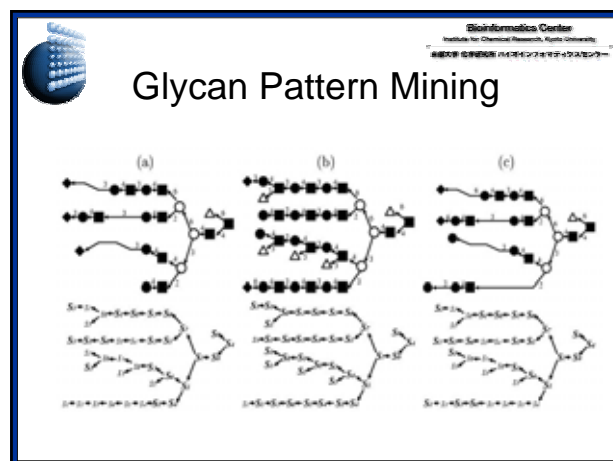
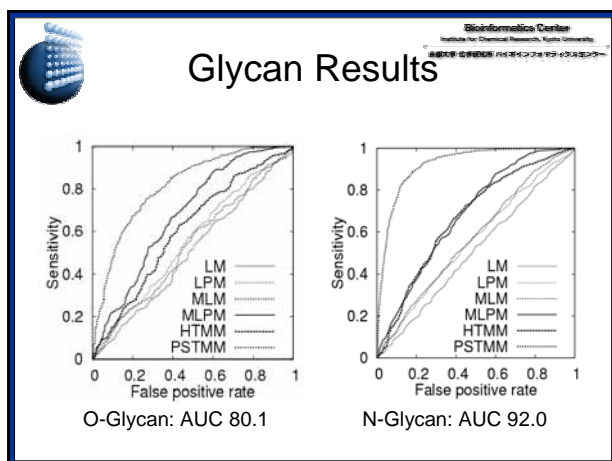
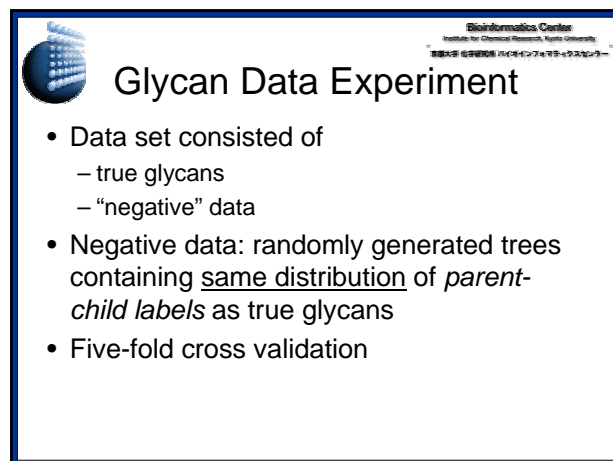
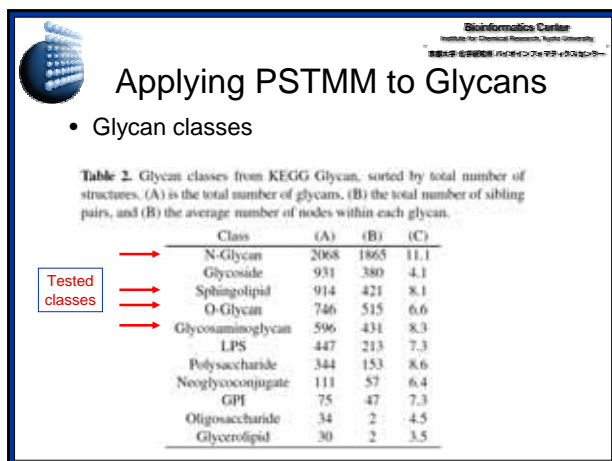
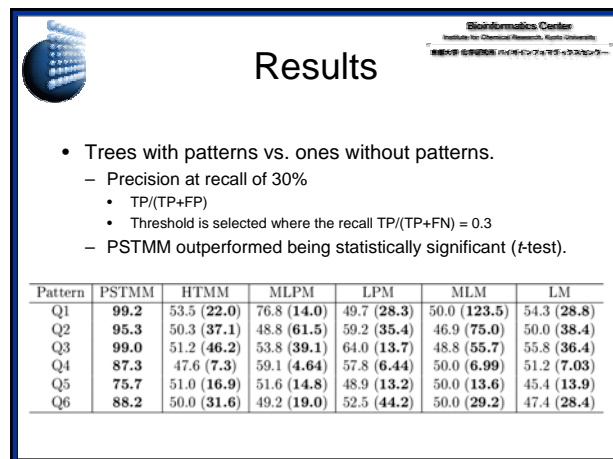
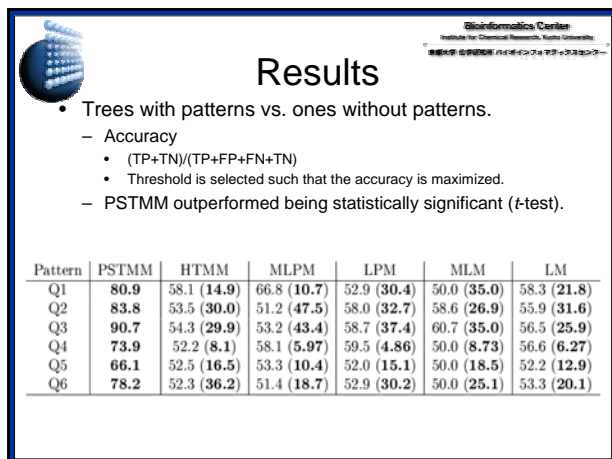
- Positive data set consisted of trees of equivalent size, embedded with various types of patterns:
- Negative data set (which is used only for test) consisted of trees of equivalent size, keeping the same distribution of parent-child pairs as that in positive data set

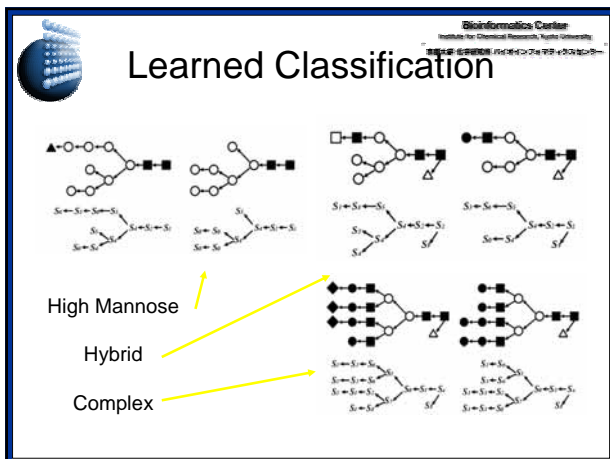
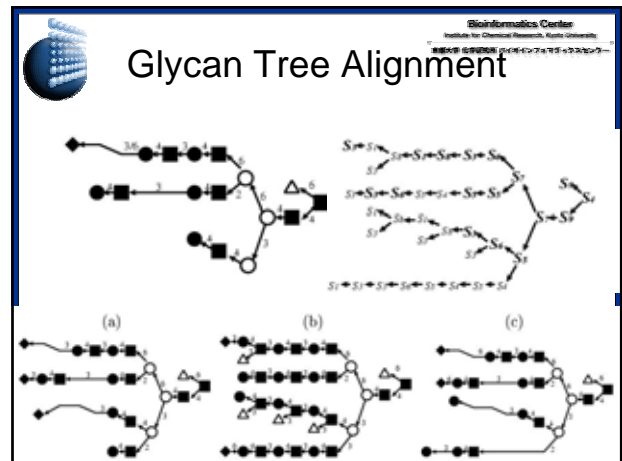
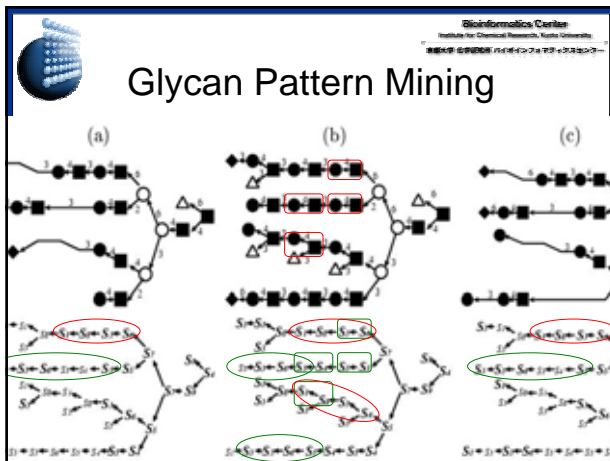


Results

- Trees with patterns vs. ones without patterns.
 - AUC (Area under the ROC curve)
 - Equivalent to Mann-Whitney-Wilcoxon test and Gini index
 - PSTMM outperformed being statistically significant (t -test).

Pattern	PSTMM	HTMM	MLPM	LPM	MLM	LM
Q1	87.6	56.9 (13.0)	71.8 (12.5)	50.7 (40.0)	49.9 (64.0)	57.6 (40.9)
Q2	89.1	51.4 (24.4)	48.5 (32.8)	58.7 (27.6)	52.6 (39.0)	53.9 (33.9)
Q3	96.1	53.0 (55.0)	51.2 (75.2)	58.4 (36.8)	55.6 (37.8)	56.3 (31.2)
Q4	80.0	48.9 (10.3)	58.9 (6.2)	60.3 (6.5)	49.9 (9.7)	54.2 (8.9)
Q5	70.3	50.5 (13.1)	51.8 (16.1)	50.0 (18.3)	49.9 (21.7)	47.1 (16.2)
Q6	84.5	50.3 (48.4)	49.1 (21.5)	51.5 (39.1)	49.9 (29.6)	49.8 (27.4)





Computational Complexity of PSTMM

- Equivalent to context free grammars for strings: maximal practical bound

	Time	Space
U, ϕ_U, τ_U	$O(T \cdot S ^2 \cdot V)$	$O(S \cdot V)$
B, ϕ_B, τ_B	$O(T \cdot S ^3 \cdot V)$	$O(S ^2 \cdot V)$
F	$O(T \cdot S ^3 \cdot V)$	$O(S ^2 \cdot V)$
D	$O(T \cdot S ^3 \cdot V)$	$O(S \cdot V)$
$\mu(a)$	$O(T \cdot S ^3 \cdot V \cdot C)$	$O(S ^3)$
$\mu(b)$	$O(T \cdot S \cdot V)$	$O(S \cdot \Sigma)$
$\mu(\pi)$	$O(T \cdot S \cdot V)$	$O(S)$
\hat{a}	$O(T \cdot S ^3)$	$O(S ^3)$
\hat{b}	$O(T \cdot S \cdot V)$	$O(S \cdot \Sigma)$
$\hat{\pi}$	$O(T \cdot S)$	$O(S)$

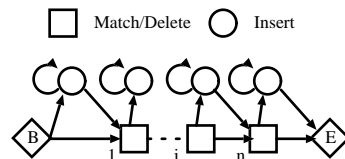
- ## Drawbacks of PSTMM
- Computational complexity is maximal practical bound
 - Overfitting problems
 - Difficult to retrieve patterns from learned states

- ## New Models of Labeled Ordered Trees
- Profile PSTMM
 - Incorporate match, insert and delete states
 - Utilize new state transitions: Down and Right
 - Ordered Tree Markov Model (OTMM)
 - Reduce dependencies on parents: State depends on that at the immediately elder sibling only, except that the eldest siblings which depend on their parents

Outline

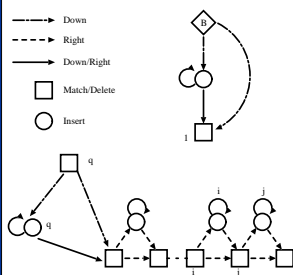
- What is a Glycan?
- Glycan = Labeled Ordered Tree
- Databases on Glycans
- Probabilistic Models for Labeled Ordered Trees and their empirical experimental results
 - Probabilistic Sibling Dependent Markov Model (PSTMM)
 - **Profile PSTMM**
 - Ordered Tree Markov Model
- Concluding Remarks

ProfileHMMs



- Match and Delete states can be found together at same positions.
- Insert states loop back to themselves or to the next position.
- Begin state at start and End state at end.

ProfilePSTMM State Model



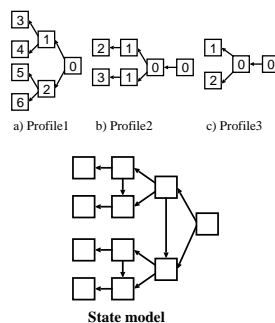
- Two types of transitions:
 - parent-child
 - between siblings
- Begin state transitions down to root node
 - Also represents End state
- Positions (1, ..., i, ..., n) are fixed
- Each position has C(i) children positions

Probability Parameters

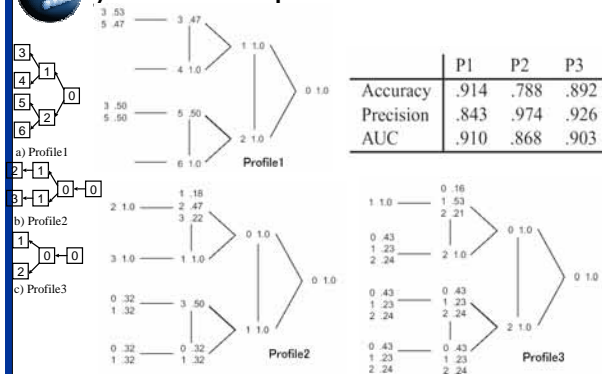
- Similar Forward, Backward, Upward and Downward parameters
- State positions are fixed:
 - Tree is traversed together with state model
 - No need to traverse every combination of states (as when they are free)
 - Much more efficient
- Profiles could be retrieved directly from match states' label output probabilities

Synthetic Experiment

- Three types of profiles
- Fifty trees containing a profile
- Fifty negative data trees with same parent-child label distribution
- Fixed state model

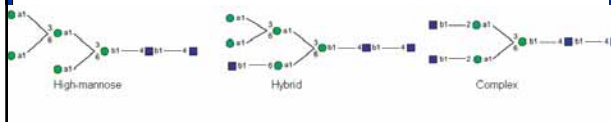


Synthetic Experiment Results

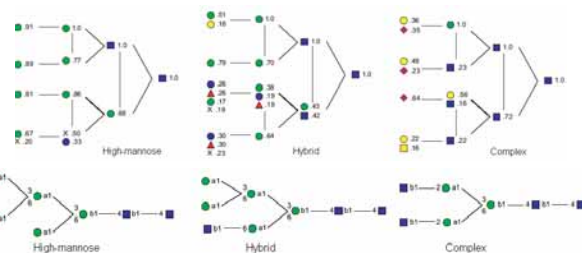


Glycan Class Experiment

- Test the three sub-classes of N-Glycans:
 - High-mannose, Hybrid, Complex type
- 50 examples for each subtype (positive examples)
- 50 negative examples generated from distribution of positive examples



Glycan Class Experiment

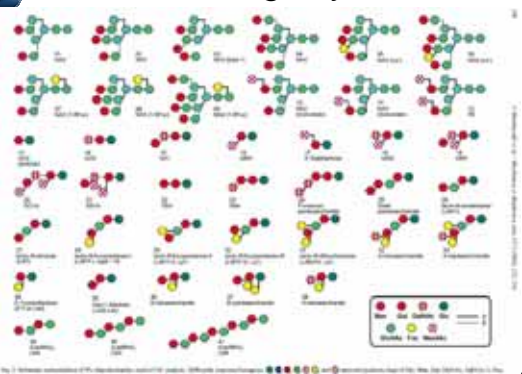


	High-mannose	Hybrid	Complex
Accuracy	.978	.982	.970
Precision	.882	.904	.882
AUC	.959	.966	.954

Galectin Experiment

- Galectins are glycan-binding proteins
- Recognize galactose at leaves
- Details still not completely understood
- Binding affinity for specific glycans tested for major galectins in Hirabayashi et al. 2002.
- Positive data: 30 weighted glycans
- Negative data: 30 glycans with same parent-child label distributions

Galectin Binding Glycan Data

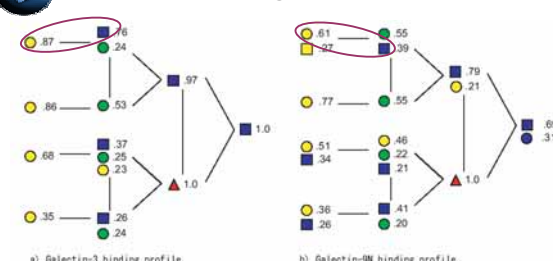


Galectin Binding Affinity Data

Table 3. Binding affinities and weights for Galectin-3 and Galectin-9N. Affinity values are normalized and inverted from the original data by Hirabayashi [16] such that higher values indicate higher affinity. Abbreviations: NA3: triantennary N-Glycan; fuc. NA3: core-fucosylated NA3; NA4: tetraantennary N-Glycan; fuc. NA4: core-fucosylated NA4; penta.: penta-saccharide; A-hexa: A-hexasaccharide; LN3: LAcNAc; LN5: (LacNAc)₅.

	Gal-3 affinity (weight)	Gal-9N affinity (weight)
NA3	1.28205 (1)	2.6316 (2)
fuc. NA3	1.21951 (1)	2.2222 (2)
NA3 type1	1.08696 (1)	1.6949 (0)
NA4	1.44928 (1)	5.5556 (5)
fuc. NA4	1.40845 (1)	4.3478 (4)
Galili penta.	1.47059 (1)	0.2273 (0)
Forssman penta.	0.16129 (0)	11.111 (11)
A-hexa	1.5873 (1)	3.8462 (3)
LN3	2.85714 (2)	1.2346 (0)
LN5	5.26316 (5)	8.3333 (8)

Galectin Binding Glycan Profiles



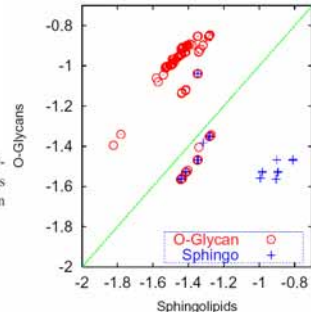
	Gal-3	Gal-9N
Acc	.847	.91
Prec	1.0	.918
AUC	.93	.931

Testing Profile Differentiation

- Train two profiles on different data sets
- Test both on both data sets to assess differentiation ability

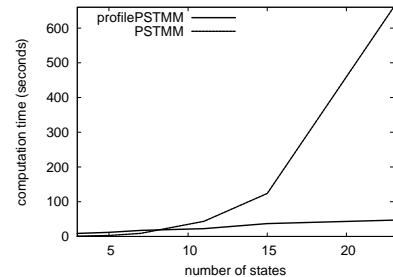
Table 5. 2x2 Contingency matrix for O-Glycan vs. Sphingolipid (Sphingo) class differentiation. The resulting discrimination rate is 923/968=95.4%.

	O-Glycan model	Sphingo model
O-Glycans	445	44
Sphingo	1	478



Improvement in Computation Time

- Efficiency increased by $O(|S|)$

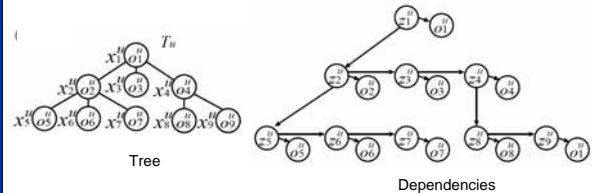


Outline

- What is a Glycan?
- Glycan = Labeled Ordered Tree
- Databases on Glycans
- Probabilistic Models for Labeled Ordered Trees and their empirical experimental results
 - Probabilistic Sibling Dependent Markov Model (PSTMM)
 - Profile PSTMM
 - **Ordered Tree Markov Model**
- Concluding Remarks

Ordered Tree Markov Model (OTMM)

- State depends on that of the immediately elder sibling, except the eldest siblings which depend on the parents



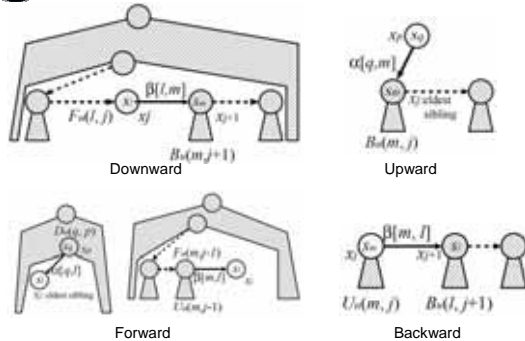
Define OTMM Parameters

- Three probability parameters
 - Initial state probability: $\pi[s_i](=P(z_i^u = s_i))$
 - Probability that the state of the root is s_i
 - State transition probability (two cases):
 - $a[s_q, s_m](=P(z_j^u = s_m | z_i^u = s_q))$
 - Probability that the state of j is s_m given that the state of the parent is s_q (eldest siblings)
 - $a[s_j, s_m](=P(z_j^u = s_m | z_i^u = s_i))$
 - Probability that the state of j is s_m given that the state of the immediately elder sibling is s_i (otherwise)
 - Label output probability: $b[s_i, \sigma_h]$
 - Probability that the state s_i outputs σ_h

Ordered Tree Markov Model (OTMM)

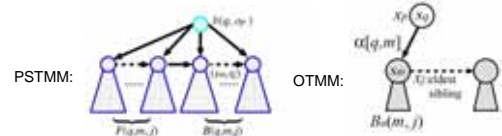
- Same learning scheme applied, i.e. EM algorithm
- Need four auxiliary probabilities again: Forward, Backward, Upward and Downward
- However, significant modification required for dynamic programming updating, since a state in OTMM does not depend on that of a parent, except the eldest siblings

DP for Computing Four Auxiliary Probabilities for OTMM



Comparison with PSTMM

- Relatively similar two auxiliary probabilities
 - Upward: child to parent



- Backward: younger sibling to elder

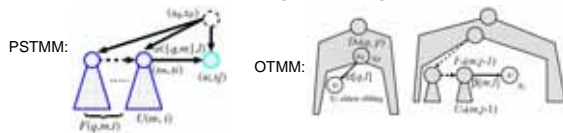


Comparison with PSTMM

- Significantly different auxiliary probabilities
 - Downward: parent to child



- Forward: older sibling to younger



EM update

- The same as that of PSTMM
- E-step computes three expectation values:
 - (s_m) = initial state expectation value
 - (s_m, s_l) = state transition expectation value
 - (s_l, h) = label output expectation value
- M-step updates our probability parameters $[s]$, $a[s]$, s_m , and $b[s]$, h using these expectation values
- Repeat E-M until some stopping condition satisfied

Computational Complexity of OTMM

- Comparison with PSTMM and HTMM
 - Efficiency always increased by $O(|S|)$

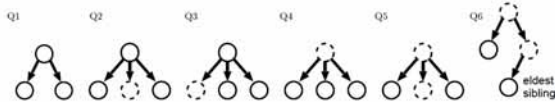
	Time
OTMM	$O(T \cdot S ^2 \cdot V)$
HTMM	$O(T \cdot S ^2 \cdot V)$
PSTMM	$O(T \cdot S ^3 \cdot V \cdot C)$
	Space
OTMM	$\max\{O(S \cdot V), O(S ^2), O(S \cdot \Sigma)\}$
HTMM	$\max\{O(S \cdot V), O(S ^2), O(S \cdot \Sigma)\}$
PSTMM	$\max\{O(S ^2 \cdot V), O(S ^3), O(S ^2 \cdot \Sigma)\}$

Experimental Setting

- Synthetic
 - Performance comparison with PSTMM
 - Capture sibling-dependent patterns
 - Discriminate between those that do and do not contain these patterns
 - Evaluate predictive accuracy and computation time using five-fold cross validation
- Real data: glycans
 - Performance comparison
 - Analyzing patterns found

Synthetic Data Experiment

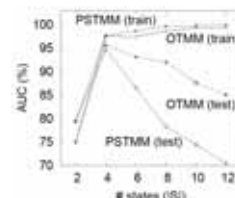
- Data set consisted of trees of equivalent size, embedded with various types of patterns:



- Parameters:
 - K: generated patterns
 - |T|: #training trees
 - |S|: #states

Synthetic Data Results

- Trees with patterns vs. ones without patterns.
 - Fixed Q1, |T|=100 and K=1,
 - Changed |S|=2,...,12
 - AUC (Area under the ROC curve) for both training and test
 - OTMM avoided overfitting to the data found in PSTMM.



Synthetic Data Results

- Trees with patterns vs. ones without patterns.
 - Fixed Q1 and K=3,
 - Changed |S|=2,...,12 and |T|=100,...,600
 - Computation time in training
 - AUC (Area under the ROC curve) for test
 - OTMM avoided overfitting, keeping much less computation time

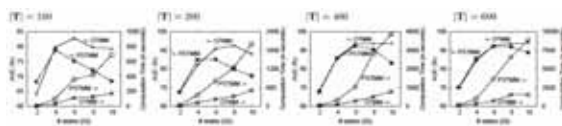


Figure 11: AUC and computation time for K=3.

Synthetic Data Results

- Trees with patterns vs. ones without patterns.
 - Fixed Q1 and |T|=200
 - Changed |S|=2,...,12 and K=1,...,4
 - Computation time in training
 - AUC (Area under the ROC curve) for test
 - OTMM avoided overfitting, keeping much less computation time

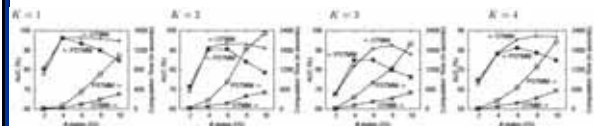


Figure 12: AUC and computation time for |T|=200.

Synthetic Data Results

- Trees with patterns vs. ones without patterns.
 - Fixed K=2, |T|=400 and |S|=6 where overfitting avoided for both OTMM and PSTMM
 - Computation time in training
 - AUC (Area under the ROC curve)
 - OTMM reduced the computation time drastically, keeping the same predictive performance.

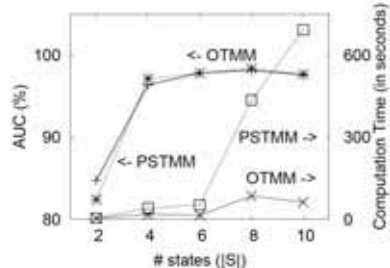
	Model	Q1	Q2	Q3	Q4	Q5	Q6
AUC (%)	OTMM	93.4	87.2	88.6	96.6	81.8	82.0
	PSTMM	92.3	89.9	91.8	95.0	79.9	75.2
Comp. Time (in seconds)	OTMM	438.7 (0.304)	583.8 (0.269)	608.8 (0.309)	379.5 (0.193)	829.8 (0.239)	581.4 (0.179)
	PSTMM	2145.6 (1.0)	2173.8 (1.0)	1970.8 (1.0)	1961.9 (1.0)	3475.4 (1.0)	3257.1 (1.0)

Glycan Data Experiment Performance Comparison

- Data set consisted of
 - Positives: N-glycans
 - Negatives: O-glycans
- Used cross-validation in the same manner as synthetic data
- Used parameter settings achieved the best performance in synthetic data

Glycan Data Experiment Performance Comparison Results

- OTMM more efficient computationally, keeping the same predictive performance



Glycan Data Experiment Glycan Pattern Mining 1

- N-Glycan

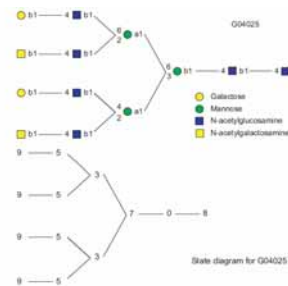


Figure 14: Example of the states learned using OTMM for a specific glycan structure.

Glycan Data Experiment Glycan Pattern Mining 2

- Three subclasses in N-glycan

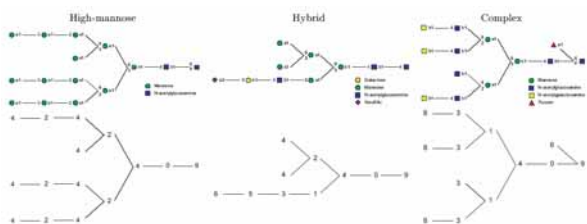


Figure 15: (top) The actual glycans, and (bottom) the most likely state paths.

Summary

- Proposed a family of probabilistic models for labeled ordered trees and their efficient learning scheme
 - OTMM reduced the complexity reasonably, avoiding overfitting and keeping the predictive performance
- Structure-based analysis: First step of glycome informatics
- There indeed seem to exist sibling-dependent relationships in glycans!
- Statistical analysis of glycans seem appropriate considering the noisiness of the data
- Important to link with other information
 - Functional annotations of genes and proteins that interact and bind with glycans

Acknowledgement

- Kiyoko F. Aoki-Kinoshita
 - Nobuhisa Ueda
 - Kosuke Hashimoto
 - Minoru Kanehisa
- (Bioinformatics Center, Kyoto University)

References

- N. Ueda, K. F. Aoki and H. Mamitsuka, SDM 2004
- K. F. Aoki et al, SIGMOD Record 2004
- K. F. Aoki et al, ISMB 2004
- N. Ueda et al, TKDE 2005