

## **INFORMATION TO USERS**

The most advanced technology has been used to photograph and reproduce this manuscript from the microfilm master. UMI films the text directly from the original or copy submitted. Thus, some thesis and dissertation copies are in typewriter face, while others may be from any type of computer printer.

**The quality of this reproduction is dependent upon the quality of the copy submitted.** Broken or indistinct print, colored or poor quality illustrations and photographs, print bleedthrough, substandard margins, and improper alignment can adversely affect reproduction.

In the unlikely event that the author did not send UMI a complete manuscript and there are missing pages, these will be noted. Also, if unauthorized copyright material had to be removed, a note will indicate the deletion.

Oversize materials (e.g., maps, drawings, charts) are reproduced by sectioning the original, beginning at the upper left-hand corner and continuing from left to right in equal sections with small overlaps. Each original is also photographed in one exposure and is included in reduced form at the back of the book.

Photographs included in the original manuscript have been reproduced xerographically in this copy. Higher quality 6" x 9" black and white photographic prints are available for any photographs or illustrations appearing in this copy for an additional charge. Contact UMI directly to order.

# **U·M·I**

University Microfilms International  
A Bell & Howell Information Company  
300 North Zeeb Road, Ann Arbor, MI 48106-1346 USA  
313/761-4700 800/521-0600



Order Number 9108187

**Biochemical and molecular genetic studies of human  
alpha-N-acetylgalactosaminidase and Schindler disease**

Wang, Anne May, Ph.D.

City University of New York, 1990

Copyright ©1990 by Wang, Anne May. All rights reserved.

**U·M·I**  
300 N. Zeeb Rd.  
Ann Arbor, MI 48106



**BIOCHEMICAL AND MOLECULAR GENETIC STUDIES OF  
HUMAN ALPHA-N-ACETYLGALACTOSAMINIDASE  
AND SCHINDLER DISEASE**

by

**ANNE MAY WANG**


A dissertation submitted to the Graduate Faculty in Biomedical Sciences  
in partial fulfillment of the requirements for the degree of  
Doctor of Philosophy, The City University of New York.

1990

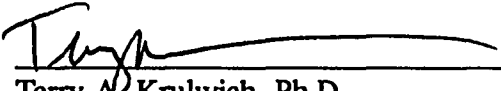
**Copyright 1990**  
**ANNE MAY WANG**  
**All Rights Reserved**

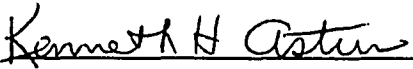
This manuscript has been read and accepted for the Graduate Faculty in Biomedical Sciences in satisfaction of the dissertation requirement for the Degree of Doctor of Philosophy, The City University of New York.

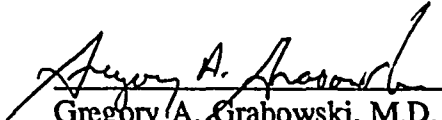
5/25/90  
\_\_\_\_\_  
Date


  
\_\_\_\_\_  
Robert J. Desnick, Ph.D., M.D.  
Chair of Examining Committee

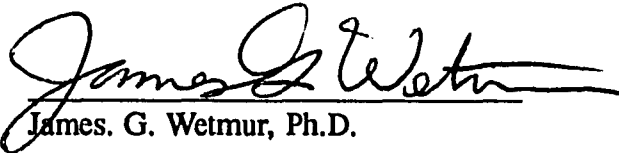
8/22/90  
\_\_\_\_\_  
Date

  
\_\_\_\_\_  
Terry A. Krulwich, Ph.D.  
Executive Officer

  
\_\_\_\_\_  
Kenneth H. Astrin, Ph.D.

  
\_\_\_\_\_  
Gregory A. Grabowski, M.D.

  
\_\_\_\_\_  
Detlev Schindler, M.D.

  
\_\_\_\_\_  
James G. Wetmur, Ph.D.

Supervisory Committee

**ABSTRACT****BIOCHEMICAL AND MOLECULAR GENETIC STUDIES OF  
HUMAN ALPHA-N-ACETYL GALACTOSAMINIDASE AND SCHINDLER DISEASE**

by

**ANNE M. WANG**

Advisor: Robert J. Desnick, Ph.D., M.D.

Human  $\alpha$ -N-acetylgalactosaminidase ( $\alpha$ -GalNAc; E.C. 3.2.1.49), the lysosomal glycohydrolase that cleaves  $\alpha$ -N-acetylgalactosaminyl moieties from glycoconjugates, is encoded by a gene localized to chromosome 22q13→ter. The deficient activity of this enzyme results in Schindler disease, an autosomal recessive disorder characterized by the increased urinary excretion of glycopeptides and oligosaccharides containing  $\alpha$ -N-acetylgalactosaminyl moieties. Two clinical phenotypes of Schindler disease have been described. The Type I disease is similar to the infantile form of neuroaxonal dystrophy, whereas the Type II disorder is characterized by disseminated angiokeratoma corporis diffusum with no neurological involvement. In these studies, a full-length cDNA encoding human  $\alpha$ -GalNAc was isolated and characterized. The full-length 2158 bp  $\alpha$ -GalNAc sequence encoded 411 amino acids and predicted a signal peptide sequence of 17 amino acids as well as six N-glycosylation sites. The functional integrity of the cDNA was demonstrated by transient expression in COS-1 cells. Northern hybridization analysis of mRNA revealed two transcripts of about 3.6 and 2.2 kb. The  $\alpha$ -GalNAc cDNA had 46.9% amino acid homology with the human  $\alpha$ -galactosidase A ( $\alpha$ -Gal A) cDNA suggesting the evolutionary relatedness of these genes. The isolation of a full-length  $\alpha$ -GalNAc cDNA facilitated the characterization of the molecular lesions in the first identified affected homozygotes with Schindler disease. The specific molecular lesions were determined by the

polymerase chain reaction (PCR) amplification and sequencing of reverse-transcribed  $\alpha$ -GalNAc transcripts from affected homozygotes with Type I and II disease. In a sib with Type I disease, a single G to A transition at nucleotide 973 was detected and resulted in a glutamic acid to lysine substitution in residue 325 of the  $\alpha$ -GalNAc polypeptide. In the type II homozygote, a single C to T transition was identified at nucleotide 985 which resulted in an arginine to tryptophan substitution at position 329. These base substitutions were confirmed by dot-blot analyses of PCR-amplified genomic DNA from members of each family using allele-specific oligonucleotides. Furthermore, transient expression of  $\alpha$ -GalNAc constructs containing these nucleotide transitions resulted in the expression of an immunoreactive polypeptide which had no detectable  $\alpha$ -GalNAc activity. The isolation of an  $\alpha$ -GalNAc cDNA also facilitated the isolation and characterization of the  $\alpha$ -GalNAc structural gene. The  $\alpha$ -GalNAc gene was ~14 kb and contained nine exons. All exon/intron junctions conformed to the GT/AG rule. Analysis of 1432 bp of 5' flanking sequence revealed two Sp1 and one GC-box promoter elements in this lysosomal housekeeping gene. Six *Alu*-repetitive elements were identified and all were in the reverse orientation. Comparison of the  $\alpha$ -GalNAc gene with the  $\alpha$ -Gal A gene revealed homologous exonic placement of  $\alpha$ -GalNAc introns 2 through 7 with  $\alpha$ -Gal A introns 1 through 6. Predicted amino acid homology among  $\alpha$ -GalNAc exons 2-7 with  $\alpha$ -Gal A exons 1-6 ranged from 46.2% to 62.7% with a few short gaps. In contrast, there was little similarity between  $\alpha$ -Gal A exon 7 and  $\alpha$ -GalNAc exons 8 and 9 which had only 15.8% homology with numerous gaps. The high exonic homologies, in addition to homologous intron placement, suggest that these genes are evolutionarily related and arose through duplication and divergence from a common ancestral gene.

## ACKNOWLEDGEMENTS

I wish to express my sincerest thanks and gratitude to Dr. Robert J. Desnick for providing the environment in which these studies were conducted. His dedication to science and academic excellence and rare insight have served as a model and motivating force during these years. It is with his faith and perseverance that this thesis was realized.

A significant portion of these studies, as well as identification of a new lysosomal disorder, is due to Dr. Detlev Schindler. His commitment to clinical care, diagnosis and scientific pursuit has contributed significantly to the area of the inborn errors of metabolism. His influence, in addition to that of Drs. David F. Bishop and Gregory A. Grabowski, have enriched my educational experience.

I am also thankful for the excellent assistance of Thomas Fitzmaurice, Raman Reddy, Yiannis Ioannou and Richard Gotlib. Their expertise has leavened the gravity on many occasions.

Lastly, I am indebted to my parents, Pin and Chen Hsiu Wang, whose unfaltering love and encouragement have supported me through not only these studies, but my life.

**TABLE OF CONTENTS**

	<i>page</i>
<b>Abstract</b>	iv
<b>List of Tables</b>	viii
<b>List of Figures</b>	ix
<b>List of Abbreviations</b>	x
<b>Background</b>	1
<b>References</b>	14
<b>Objectives and Significance</b>	19
<b>Chapter One</b>	21
<b>References</b>	48
<b>Chapter Two</b>	51
<b>References</b>	76
<b>Chapter Three</b>	79
<b>References</b>	97
<b>Concluding Remarks</b>	99

**LIST OF TABLES**

	<i>page</i>
<b>Chapter Two</b>	
<b>Table 1</b>	<b>71</b>
<b>Chapter Three</b>	
<b>Table 1</b>	<b>89</b>
<b>Table 2</b>	<b>95</b>

**LIST OF FIGURES**

	<i>page</i>
<b>Chapter One</b>	
<b>Figure 1</b>	33
<b>Figure 2</b>	36
<b>Figure 3</b>	39
<b>Figure 4</b>	40
<b>Figure 5</b>	44
<b>Chapter Two</b>	
<b>Figure 1</b>	61
<b>Figure 2</b>	63
<b>Figure 3</b>	67
<b>Figure 4</b>	69
<b>Figure 5</b>	72
<b>Chapter Three</b>	
<b>Figure 1</b>	87
<b>Figure 2</b>	90
<b>Figure 3</b>	92

**LIST OF ABBREVIATIONS**

$\alpha$ -Gal A	$\alpha$ -galactosidase A
$\alpha$ -Gal B	$\alpha$ -galactosidase B
$\alpha$ -GalNAc	$\alpha$ - <i>N</i> -acetylgalactosaminidase
A	adenosine
bp	base pair(s)
C	cytosine
cDNA	complementary DNA
DNA	deoxyribonucleic acid
G	guanosine
kCal	kilocalorie
kDa	kilodalton
mRNA	messenger RNA
4-MU- $\alpha$ -Gal	4-methylumbelliferyl- $\alpha$ -D-galactopyranoside
4-MU- or pNP- $\alpha$ -GalNAc	4-methylumbelliferyl- or p-nitrophenyl- $\alpha$ - <i>N</i> -acetylgalactosaminide
NaDodSO <sub>4</sub> /PAGE	sodium dodecyl sulfate polyacrylamide gel electrophoresis
nt	nucleotide(s)
PCR	polymerase chain reaction
RNA	ribonucleic acid
T	thymidine

## BACKGROUND

### A. Biochemical Characterization of $\alpha$ -N-Acetylgalactosaminidase

It is now recognized that human lysosomal  $\alpha$ -galactosidase B ( $\alpha$ -Gal B) is an  $\alpha$ -N-acetylgalactosaminidase ( $\alpha$ -D-N-acetylgalactosamine-N-acetylgalactosaminohydrolase,  $\alpha$ -GalNAc; E.C 3.2.1.49). The confusion regarding the natural substrate specificity of  $\alpha$ -Gal B resulted from the use of artificial galactosides as substrates which detected two  $\alpha$ -galactosidase isozymes in human tissues (1, 2). The major isozyme, designated  $\alpha$ -Gal A, represented about 80-90% of total  $\alpha$ -galactosidase activity and was heat labile, while the other isozyme, designated  $\alpha$ -Gal B, represented the remaining activity and maintained 80% of enzymatic activity at 50°C for at least 4 hours (3-9). Myoinositol specifically inhibited  $\alpha$ -Gal A activity (5, 7) whereas the B isozyme was competitively inhibited by  $\alpha$ -N-acetylgalactosamine (3, 4, 6). The two "isozymes" were separable by starch and cellulose acetate gel electrophoresis (7, 10), isoelectric focusing (9, 11, 12), and ion exchange chromatography (5, 6, 9, 13, 14). These studies revealed the presence of multiple forms of  $\alpha$ -Gal A with pI values ranging from 4.3 to 5.1 while  $\alpha$ -Gal B had a single form with a pI of about 4.5. It was also clearly demonstrated that  $\alpha$ -Gal A activity was deficient in Fabry disease (1, 7, 8, 10-12, 15) and that the residual activity detected in classical hemizygotes with Fabry disease was due to  $\alpha$ -Gal B. Neuraminidase treatment of  $\alpha$ -Gal A and B converted the multiple  $\alpha$ -Gal A forms to a single activity band at pI 4.6, while the migration of  $\alpha$ -Gal B was unaffected (2, 16, 17). In the early 1970s, long before it was recognized that  $\alpha$ -Gal B was an  $\alpha$ -N-acetylgalactosaminidase, these findings were erroneously interpreted;  $\alpha$ -Gal A was thought to be a sialated form of  $\alpha$ -Gal B.

However, the two isozymes were immunologically distinguishable (5, 7, 15): anti- $\alpha$ -Gal A did not crossreact with  $\alpha$ -Gal B, and anti- $\alpha$ -Gal B did not crossreact with  $\alpha$ -Gal A. The residual activity in Fabry disease could be immunoprecipitated with anti- $\alpha$ -Gal B, but not with

anti- $\alpha$ -Gal A. These findings suggested that the two isozymes were indeed distinct proteins. Subsequently, these genes were mapped to separate chromosomes;  $\alpha$ -Gal B to 22q13→ter (18, 19) and  $\alpha$ -Gal A to Xq21.33→22 (20), confirming their independent nature.

In 1977, Dean and Sweeley in the United States, and Schram and colleagues in the Netherlands independently demonstrated that  $\alpha$ -Gal B functioned as an  $\alpha$ -*N*-acetylgalactosaminidase which hydrolyzed artificial and natural substrates with terminal  $\alpha$ -*N*-acetylgalactosaminyl moieties and was competitively inhibited by  $\alpha$ -*N*-acetylgalactosamine (3, 4, 6, 21, 22). While studies suggested that  $\alpha$ -Gal B was able to cleave terminal  $\alpha$ -galactosyl residues from glycolipid substrates (albeit at a very slow rate) and water soluble oligosaccharides,  $\alpha$ -Gal B was 2-5 fold more active as an  $\alpha$ -*N*-acetylgalactosaminidase on glycolipid and water soluble oligosaccharide substrates with terminal  $\alpha$ -*N*-acetylgalactosaminyl moieties (14). Furthermore,  $\alpha$ -Gal B had a higher affinity for *o*- or *p*-nitrophenyl- $\alpha$ -*N*-acetylgalactosaminide ( $K_m \sim 1$ -2 mM) than for *o*- or *p*-nitrophenyl- $\alpha$ -*D*-galactoside ( $K_m \sim 7$ -20 mM) (3, 4, 14).  $\alpha$ -Gal B also hydrolyzed the terminal  $\alpha$ -*N*-acetylgalactosaminyl residues from Forssman hapten, a canine intestinal neutral glycosphingolipid, and blood group A trisaccharide (6). Thus, it was concluded that  $\alpha$ -Gal B was actually an  $\alpha$ -*N*-acetylgalactosaminidase which did not cleave the glycolipid substrates accumulated in Fabry disease.

Human  $\alpha$ -GalNAc has been purified from a variety of sources, including placenta and liver (3-6, 9, 14). The enzyme is a homodimeric glycoprotein with a subunit molecular weight of about 48 kDa. Support for the homodimeric structure is based on biosynthetic studies of  $\alpha$ -GalNAc in human fibroblasts which indicated that the enzyme was synthesized as a 65 kDa precursor which was processed to a 48 kDa mature lysosomal form (23). Additionally, analysis of radiolabelled oligosaccharides from immunoprecipitated enzyme revealed a high-mannose type oligosaccharide structure ( $\text{Man}_{8,9}\text{GlcNAc}$ ) on both precursor and processed enzyme forms; only

the precursor had phosphorylated mannose residues. The absence of complex-type oligosaccharide moieties containing terminal sialic residues was supported by neuraminidase studies which did not alter the pI of the enzyme, and by qualitative oligosaccharide analysis of the mature, processed enzyme (24, 25). The range of substrates hydrolyzed by  $\alpha$ -GalNAc have not yet been identified, but any compound containing an  $\alpha$ -*N*-acetylgalactosaminyl residue is a potential substrate. These glycoconjugates include glycosphingolipids, mucins, *O*- and *N*-linked glycoproteins, and the mucopolysaccharide, keratin sulfate.

#### B. Lysosomal Biology

$\alpha$ -Gal B is representative of a class of enzymes present in lysosomes, virtually all of which are acid hydrolases (26). Genetic deficiencies of specific lysosomal hydrolases result in the different lysosomal storage diseases. To date, over 30 lysosomal storage diseases in humans have been identified and have generated wide interest in the cell biology of lysosomes and the biochemistry of lysosomal enzymes. Intense studies have led to the delineation of the unique properties of lysosomal enzymes, especially related to their biosynthesis, processing, and trafficking to the lysosome. Briefly, lysosomal enzymes are typically synthesized as preproenzymes which are converted to proenzymes in the endoplasmic reticulum and cis-Golgi, and then processed into mature forms within the lysosome (27). These extensive co- and post-translational modifications include glycosylation and sequential carbohydrate modification and phosphorylation. Transport to the lysosome is mediated by the mannose-6-phosphate residues on most of the enzymes. These residues are added in the Golgi and are recognized by mannose-6-phosphate receptors in membranes which are responsible for transporting the enzymes to the endosomes and ultimately to the lysosomes.

The study of lysosomal enzymes and lysosomal storage diseases has proceeded to the molecular level. In general, each lysosomal enzyme constitutes approximately 0.01% of total cellular protein (28). Despite this low abundance, several cDNA sequences have been cloned for lysosomal enzymes, including  $\alpha$ -fucosidase (29),  $\alpha$ -Gal A (30, 31), and  $\beta$ -hexosaminidase  $\alpha$ - (32-34) and  $\beta$ -chains (35). Several strategies have been employed in cloning these sequences, including screening cDNA libraries with radiolabelled probes which may be oligonucleotides constructed from peptide sequence (34), cDNA sequences from other species (36, 37), or an exonic fragment of a human gene from the same gene family (38). Screening of cDNA libraries enhanced for message by hybrid selection and polysome immunoabsorption (32) and of expression libraries constructed in  $\lambda$ gt11 with monospecific antisera (29, 30, 39-41) also have been successful. More recently, the polymerase chain reaction (PCR) and mixed oligonucleotide primed amplification of cDNA (MOPAC) has permitted the cloning of partial cDNA sequences using amino acid microsequence data. These partial cDNAs may then be used as specific probes for more efficient library screening (42). Analysis of these cDNA sequences has confirmed at the molecular level sequences corresponding to pre- and prepro-peptide sequences which were presumed cleaved during the transport of lysosomal enzymes. These pre- and prepro-peptide sequences were hydrophobic as expected (43).

To date, only five genomic lysosomal enzyme sequences have been reported:  $\alpha$ -Gal A (44), the  $\alpha$ - (45) and  $\beta$ -chains (46) of  $\beta$ -hexosaminidase,  $\beta$ -glucosidase (47) and acid phosphatase (48).  $\alpha$ -Gal A is 12 kb with seven exons and has been sequenced in its entirety (49). The first exon contains the 5' untranslated region, the signal peptide, and the first 33 amino acids of the mature enzyme. The gene is unusual in that it lacks a 3' untranslated region. The  $\beta$ -hexosaminidase  $\alpha$ - and  $\beta$ -chains are structurally similar, and extensive homology between the two genes suggests evolution of the two sequences from a common ancestral gene. The  $\alpha$ -

chain gene spans about 40 kb with 14 exons, and the  $\beta$ -chain gene spans about 35 kb and also has 14 exons. Twelve of the 13 introns in the  $\beta$ -chain were shown to interrupt the gene at homologous positions in the  $\alpha$ -chain.  $\beta$ -glucosidase is 7.6 kb with 11 exons. The gene and its tightly linked 5.8 kb pseudogene which has 96% nucleotide homology have been sequenced. Acid phosphatase spans nine kb with 11 exons. Like the  $\alpha$ -Gal A gene, exon 1 encodes the signal peptide and the first eight residues of the mature protein.

As a result of cloning cDNA and genomic sequences for many of the proteins associated with genetic diseases, it has been possible to determine the precise molecular pathology underlying these disorders and study their effects. Among the lysosomal disorders, a number of mutations in the  $\alpha$ -subunit of  $\beta$ -hexosaminidase A have been described in Tay-Sachs disease ( $\beta$ -hexosaminidase A deficiency). The major  $\alpha$ -subunit mutation in the Ashkenazi Jewish population is an insertion of four bases into exon 11 (50) which causes a frameshift and subsequent termination nine nucleotides downstream. A second mutation occurring within this population is a G to C transversion at the 5' border of intron 12 (51-53) which alters the consensus GT dinucleotide necessary for intron splicing and results in no detectable mRNA. A mutation responsible for the adult-form or late-onset  $G_{M2}$  gangliosidosis has been identified as a G to A transition in exon 7 which causes the substitution of serine for glycine at residue 269 (54, 55). This substitution results in poor association of the  $\alpha$ - and  $\beta$ -subunits in the active tetrameric configuration (56). Among the French-Canadian population, the prevalent mutation has been identified as a 7.6 kb deletion of exon 1 and flanking sequences of the  $\alpha$ -subunit gene which may have occurred by homologous recombination of misaligned *Alu* sequences (57, 58). Other mutations which cause  $\beta$ -hexosaminidase A deficiency have also been identified and include base substitutions and deletions which cause amino acid substitutions, frameshifts and premature termination (59-62). Analyses of these mutations have revealed the abnormal

processing and trafficking of the defective  $\alpha$ -subunits as well as altered kinetics and substrate specificities. As illustrated by the example of Tay-Sachs disease, the molecular lesions resulting in this disorder, as well as in other lysosomal disorders (e.g., Fabry disease, Gaucher disease), are heterogeneous in nature and include partial gene duplications (63).

A variety of methods have been used to delineate the specific nature of these mutations, including the construction and screening of cDNA and genomic libraries from the DNA of affected individuals. S1 nuclease protection studies (64, 65) and RNase A cleavage analysis (66) have been successfully used to detect single base substitutions as well as small insertions and deletions. Most recently, primer directed amplification of specific DNA sequences by the PCR (67) has allowed direct, rapid identification and sequencing of mutations. The effects of these mutations have been studied using a variety of techniques including transient expression, pulse-chase labelling and immunoprecipitation (61, 68).

### C. The Evolution of Human Gene Families

Comparative studies at the nucleotide and deduced amino acid levels have provided insight into the mechanisms of molecular evolution. Gene duplication from a single ancestral gene is the primary mechanism in the evolution of genes with new functions (69). This initial event may occur by several mechanisms: 1) nonhomologous chromosomal breakage and reunion, 2) homologous, but unequal recombination between repeated elements such as *Alu I* or *Kpn I* sequences, 3) RNA-mediated duplications, and 4) gene amplification. After this initial event, one copy of the gene may retain its original function while a duplicated copy may diverge by the accumulation of mutations which alter its function or regulation.

The most extensive work in the structure and molecular evolution of genes has involved the preproinsulin and the  $\beta$ -globin gene families. Perler and colleagues (70) established a

method for calculating evolutionary divergence time based on the nucleotide replacement substitution rate. Mutations have been estimated to occur in DNA sequences at a minimum rate of  $7 \times 10^{-9}$  per nucleotide site per year. When exonic regions of two homologous sequences are compared, two types of substitutions are evident: replacement substitutions which result in amino acid replacement, and silent substitutions which do not alter the amino acid. Nucleotide differences among homologous proteins have been used to estimate the replacement site divergence of the corresponding genes by determining the minimum number of base changes necessary to generate the observed amino acid replacements. Corrections have been devised for multiple base change events with individual codons. Within a protein family, replacement changes appear to accumulate linearly with divergence time; this is supported by the fossil record. In general, silent and intronic changes do not provide good evolutionary information as they accumulate more rapidly than replacement substitutions.

The Perler method for measuring evolutionary time was based on sequence analysis of the rat, human, and chicken preproinsulin genes. Most vertebrate species have one insulin gene, but rat, mice and three species of fish have two nonallelic insulin genes. The rat preproinsulin gene II contains a 499 bp intron in the C peptide that is not present in gene I. Thus, one gene either gained or lost the intron during duplication. To determine which gene was ancestral, the preproinsulin gene from chicken was cloned. The chicken gene contained a 3.5 kb intron in the C peptide, and therefore, the rat gene II was the ancestral gene.

Using the Perler method, the structure and evolution of the human  $\beta$ -globin gene family was delineated by Efstradiatis *et al.* (71). The complete sequences of embryonic  $\epsilon$ - (72), fetal  $\delta$ - and  $\gamma$ - (73), and adult  $\alpha$ - (74) and  $\beta$ - (75) globins have been determined. From these sequences, pairwise calculations of replacement and silent substitutions according to the Perler method were performed. This analysis allowed the construction of an evolutionary tree for the

human  $\beta$ -like globin genes based entirely on nucleotide sequence comparisons. The tree agreed with data available from the fossil record.

Nucleotide sequence data may also provide structural information regarding the evolutionary conservation of homologous proteins. Among the lysosomal enzymes, extensive homology has been found between the  $\beta$ -hexosaminidase  $\alpha$ - and  $\beta$ -chain cDNA and chromosomal genes (46). The  $\alpha$ - and  $\beta$ -chains are the polypeptide subunits of the two major isozymes of  $\beta$ -hexosaminidase; hexosaminidase A is composed of two  $\alpha$ - and two  $\beta$ -subunits, and hexosaminidase B, contains four  $\beta$ -subunits. A third minor isozyme, hexosaminidase S, is composed of  $\alpha$ -subunits. Comparison of the cDNA sequences for the  $\alpha$ - and  $\beta$ -chains revealed 55% nucleotide and 57% amino acid similarity, suggesting divergence of the two genes from a common ancestral gene (45). The finding of extensive similarity was not unexpected, as the subunits share several common features. First, homodimers of each chain have the ability to hydrolyze the same artificial substrates *in vitro*, and therefore must have similar substrate binding and active sites. Secondly, several independently prepared antisera against hexosaminidase S have shown crossreactivity with hexosaminidase B, suggesting that the two subunits share related epitopes. Characterization of the  $\alpha$ - and  $\beta$ - chain genomic sequences (45, 46) has shown that the structural organization of the genes were highly conserved. Comparison of the two genes revealed that 12 of the 13 introns interrupted the coding regions at homologous positions, clearly indicating that they were duplicated from a single ancestral gene and that the single difference in intron/exon structure was probably due to subsequent divergence.

Study of the genomic structure of genes in relation to molecular evolution has also provided support for hypotheses concerning the conservation of active site regions and for exons encoding functional domains. X-ray crystallography studies of the same protein purified from different species have shown that secondary structure is highly conserved (76, 77). In a study

of the phosphoglycerate kinase gene from six species (76), the amino acids forming the substrate binding cleft were very highly conserved. The homology between yeast and mammalian sequences for the binding cleft was greater than 94%, whereas the overall homology was less than 65%. Trypanosome phosphoglycerate kinase had only 44% overall homology with the mammalian enzyme, but showed 74% homology at the binding cleft. Analysis of the sequences of genes with functionally conserved regions has shown that these regions undergo replacement substitutions at lower rates than for nonconserved regions (77). The lysosomal hydrolase,  $\alpha$ -glucosidase, has also been shown to have significant homology with human isomaltase and the rabbit intestinal sucrase-isomaltase complex (78). Alignment of these sequences revealed 26% homology. Ten of 13 amino acids were identical in the region surrounding the known active site of sucrase and isomaltase, suggesting that the active site of  $\alpha$ -glucosidase may be in this region. From these sequences, it was proposed that the isomaltase-sucrase single gene precursor arose by duplication of an ancestral isomaltase gene. Since  $\alpha$ -glucosidase has the same degree of homology with isomaltase and sucrase,  $\alpha$ -glucosidase must have duplicated from the same ancestral gene. A more recently emerging gene family with potential active site conservation is the sulfatases. Protein homologies have been identified among the human arylsulfatases A and B, glucosamine-6-sulfatase and steroid sulfatase (42, 79). While these proteins share an overall amino acid homology of ~20%, the highest degree of homology among the coding regions of these genes was concentrated at the *N*-terminal third of the four mature polypeptides; studies with arylsulfatase A indicate that these regions may be involved in the enzyme's active site (79).

Conservation of active sites or functional protein domains by genomic structural organization has also been proposed (80). Structural characterization of the phosphoglycerol kinase gene (81) has provided evidence supporting the concept of exon shuffling. Comparative

analysis of human phosphoglycerol kinase, maize alcohol dehydrogenase and chicken glyceraldehyde-3-phosphate dehydrogenase revealed that the sequence for the nucleotide binding domain showed homologous intron placement in each gene. This domain is specified by five exons; three exons coded for the first mononucleotide domain and two for the second domain. The homology of the exon organization in structurally similar regions suggests that the nucleotide binding domain evolved by gene duplication and was subsequently dispersed to different proteins through a process of intron-mediated recombination. Thus, this gene family supports the concept of exon shuffling in the conservation of functional domains.

#### D. Schindler Disease

The first patients with  $\alpha$ -GalNAc deficiency were reported in the last year (82-85). Although these cases are due to markedly deficient  $\alpha$ -GalNAc activity, the clinical manifestations in the unrelated patients differed remarkably. Type I disease is characterized by a clinical course identical to that of an infantile form of neuroaxonal dystrophy, whereas the patient with Type II disease has no neurological involvement, but does have angiokeratoma corporis diffusum.

The first patients with Type I Schindler disease were two German brothers, the consanguineous offspring of fourth cousins. Development of the brothers was normal until about 15 months of age, after which each experienced rapid regression. Clinical onset was signaled by grand mal seizures in the younger sib at 9 months, and by clumsiness and falling episodes at 12 months in the older sib. From 15 months of age, both developed strabismus, nystagmus, optic atrophy, muscular hypotonia and frequent myoclonic seizures. They also exhibited profound psychomotor retardation. By three to four years of age, they were immobile, had decorticate postures, cortical blindness, and had little, if any, contact with the environment.

The deficient activity of  $\alpha$ -GalNAc was demonstrated using the artificial substrate p-nitrophenyl- $\alpha$ -*N*-acetylgalactosaminide (82, 83). This result was later confirmed using a newly synthesized, highly sensitive and specific fluorogenic substrate, 4-methylumbelliferyl- $\alpha$ -*N*-acetylgalactosaminide (4-MU- $\alpha$ -GalNAc) (84). This substrate allowed reliable detection of the enzyme defect in plasma, lymphoblasts, and fibroblasts from the affected brothers as well as the intermediate activity levels in their heterozygous parents, consistent with an autosomal recessive mode of inheritance.

An abnormal urinary oligosaccharide profile was also observed in the affected brothers (82-84). The excreted products were later shown to be *O*-linked glycopeptides (87) with the structures: NeuNAc $\alpha$ 2 $\rightarrow$ 3Gal $\beta$ 1 $\rightarrow$ 3(NeuNAc $\alpha$ 2 $\rightarrow$ 6)GalNAc $\alpha$ 1 $\rightarrow$ O-serine and -threonine, NeuNAc $\alpha$ 2 $\rightarrow$ 3Gal $\beta$ 1 $\rightarrow$ 4GlcNAc $\beta$ 1 $\rightarrow$ 6(NeuNAc $\alpha$ 2 $\rightarrow$ 3Gal $\beta$ 1 $\rightarrow$ 3)GalNAc $\alpha$ 1 $\rightarrow$ O-serine and -threonine, and GalNAc $\alpha$ 1 $\rightarrow$ O-serine and -threonine. Only the older affected sib, who was blood group type A positive, excreted the blood group A trisaccharide: GalNAc $\alpha$ 1 $\rightarrow$ 3(Fuc $\alpha$ 1 $\rightarrow$ 2)Gal.

Further clinical evaluation of the two brothers revealed several interesting findings. Ultrastructural evaluation of biopsied rectal tissue revealed abnormal tubulovesicular material free in the cytoplasm of a few preterminal and terminal axons in the myenteric plexus. Subsequent examination of the frontal lobe cortical biopsy showed the characteristic neuropathology of neuroaxonal dystrophy (88), which included the characteristic "spheroids" in terminal axons that contained tubulovesicular and other membranous arrays. These structures were essentially identical to those observed in the inherited neuroaxonal dystrophies, including Seitelberger disease and Hallervorden-Spatz disease. The morphologic hallmark in these disorders is generalized axonal dystrophy with peculiar swellings or "spheroids" in terminal autonomic axons. Each is inherited as an autosomal recessive trait (89).

Since the clinical and neuropathologic findings clearly classified this disease as a neuroaxonal dystrophy, efforts were directed to determine the  $\alpha$ -GalNAc activity in unrelated patients with Seitelberger disease (84). In eight biopsy- or autopsy-proven cases, the  $\alpha$ -GalNAc activities in plasma, cultured cells, or cortical tissue were within the normal range. Thus,  $\alpha$ -GalNAc was not the enzymatic defect in these unrelated patients with infantile neuroaxonal dystrophy. Moreover, it could be concluded that the biochemical defects underlying the neuroaxonal dystrophy phenotype were heterogeneous.

In the second reported case of Schindler disease, and the first with the Type II phenotype, (86) the propositus was a 46-year old Japanese woman, a daughter of a first cousin marriage (85). Development of this patient was normal with the notable exception of the eruption of disseminated angiokeratoma corporis diffusum during the third decade of life. The development of angiokeratoma is characteristic of several lysosomal disorders including Fabry disease ( $\alpha$ -Gal A deficiency), fucosidosis ( $\alpha$ -fucosidosis deficiency) and galactosialidosis ( $\beta$ -galactosidase and sialidase deficiency). However, the patient had normal levels of activity for these and a number of other lysosomal hydrolases. An abnormal urinary oligosaccharide profile was observed in the affected individual (85, 90) and the products were purified and determined to be essentially identical to the ones identified in Type I disease. This finding suggested the diagnosis of  $\alpha$ -GalNAc deficiency (86). The patient had less than 1% normal activity, while her two children had intermediate activity levels, consistent with the autosomal recessive mode of inheritance. Electron microscopy of a skin biopsy from this individual revealed vacuolation of the cytoplasm in endothelial cells. Similar vacuoles were also observed in sweat gland cells, dermal fibrocytes, vascular pericytes, smooth muscle cells, fat cells, Schwann cells and neuron axons. Such vacuoles were essentially absent in the brothers affected with Type I disease.

Thus far, two distinct subtypes of  $\alpha$ -GalNAc deficiency have been described. These subtypes are phenotypically distinct, but result from the same deficient enzymatic activity.

## REFERENCES

1. Beutler, E. and Kuhl, W. (1971) *J Lab Clin Med* 78:987.
2. Kint, J. A. (1971) *Arch Int Physiol Biochem* 78:633-634.
3. Dean, K. J., Sung, S. and Sweeley, C. C. (1977) *Biochem Biophys Res Com* 77:1411-1417.
4. Schram, A. W., Hamers, M. N. and Tager, J. M. (1977) *Biochim Biophys Acta* 482:138-144.
5. Beutler, E. and Kuhl, W. (1972) *J Biol Chem* 247: 7195-7200.
6. Callahan, J. W., Lassilla, E. L., DenTandt, W. and Philippart, M. (1973) *Biochem Med* 7:424-431.
7. Rietra, P. J. G. M., VanDenBergh, F. A. J. T. M. and Tager, J. M. (1975) *Clin Chim Acta* 62:401-413.
8. Johnson, D. L. and Desnick, R. J. (1978) *Biochim Biophys Acta* 538:195-204.
9. Kusiak, J. W., Quirk J. M. and Brady, R. O. (1978) *J Biol Chem* 253:184-190.
10. Beutler, E. and Kuhl, W. (1972) *Amer J Hum Genet* 24:237-249.
11. Salvayre, R., Maret, A., Negre, A. and Douste-Blazy, L. (1979) *Eur J Biochem* 100:377-383.
12. Salvayre, R., Negre, A., Maret, A., Lenoir, G. and Douste-Blazy, L. (1981) *Biochim Biophys Acta* 659:445-456.
13. Dean, K. J. and Sweeley, C. C. (1979) *J Biol Chem* 254:9994-10000.
14. Dean K. J. and Sweeley, C. C. (1979) *J Biol Chem* 254:10001-10005.
15. Schram, A. W., Hamers, M. N., Brouwer-Kelder, B., Donker-Koopman, W. E. and Tager, J. M. (1977) *Biochim Biophys Acta* 482:125-37.
16. Ho, M. W., Beutler, E., Tennant, L. and O'Brien, J. S. (1972) *Am J Hum Genet* 24:256-266.
17. Mapes, C. A. and Sweeley, C. C. (1973) *Arch Biochim Biophys* 158:297-304.
18. deGroot, P. G., Westerveld, A., Meera-Khan, P. and Tager, J. M. (1978) *Hum Genet* 44:305-312.

19. Mysczkiewicz, B. A., Wilkinson, F. E., Kung, H. J. and Sweeley, C. C. (1984) *Fed Proc* 43:1529.
20. Unpublished results.
21. Dean, K. J., Sung, S. and Sweeley, C. C. (1978) *Adv Exp Med Biol* 101:515-523.
22. Schram, A. W., Hamers, M. N. and Tager, J. M. (1978) *Adv Exp Med Biol* 101:525-529.
23. LeDonne, N. C., Fairley, J. L. and Sweeley, C. C. (1983) *Arch Biochim Biophys* 224:186-195.
24. Dean, K. J. and Sweeley, C. C. (1979) *J Biol Chem* 254:10006-10010.
25. Unpublished results.
26. deDuve, C. (1983) *Eur J Biochem* 137:391-397.
27. VonFigura, K. and Hasilik, A. (1986) *Ann Rev Biochem* 55:167-193.
28. Myerowitz, R., and Neufeld, E. F. (1981) *J Biol Chem* 256:3044-3048.
29. Fukushima, H., DeWet, J. R. and O'Brien, J. S. (1985) *Proc Natl Acad Sci USA* 82:7289-7293.
30. Calhoun, D. H., Bishop, D. F., Bernstein, H. S., Quinn, M., Hantzopoulos, P. and Desnick, R. J. (1986) *Proc Natl Acad Sci USA* 82:7364-7368.
31. Bishop, D. F., Calhoun, D. H., Bernstein, H. S., Hantzopoulos, P., Quinn, M. and Desnick, R. J. (1986) *Proc Natl Acad Sci USA* 83:4859-4863.
32. Myerowitz, R. and Proia, R. (1984) *Proc Natl Acad Sci USA* 81:5394-5398.
33. Myerowitz, R., Peikarz, R., Neufeld, E. F., Shows, T. B. and Suzuki, K. (1985) *Proc Natl Acad Sci USA* 82:7830-7834.
34. Korneluk, R. G., Mahuran, D. J., Neote, K., Klavins, M. H., O'Dowd, B. F., Tropak, M., Willard, H. F., Anderson, M., Lowden, J. A. and Gravel, R. A. (1986) *J Biol Chem* 261:8407-8413.
35. O'Dowd, B. F., Quan, F., Willard, H. F., Lanhonwah, A., Korneluk, R. G., Lowden, J. A., Gravel, R. A. and Mahuran, D. J. (1985) *Proc Natl Acad Sci USA* 82:1184-1188.
36. Guise, K. S., Korneluk, R. G., Waye, J., Lanhonwah, A. M., Quan, F., Palmer, R., Ganschow, R. E., Sly, W. S. and Gravel, R. A. (1985) *Gene* 34:105-110.

37. Chan, S. J., Segundo, B. S., McCormick, M. B. and Steiner, D. F. (1986) *Proc Natl Acad Sci USA* 83:7721-7725.
38. Faust, P. L., Kornfeld, S. and Chirgwin, J. M. (1985) *Proc Natl Acad Sci USA* 82:4910-4914.
39. Sorge, J., West, C., Westwood, B. and Beutler, E. (1985) *Proc Natl Acad Sci USA* 82:7289-7293.
40. Tsuji, S., Choudary, P. V., Martin, B. M., Winfield, S., Barranger, J. A. and Ginns, E. I. (1986) *J Biol Chem* 261:50-53.
41. Martiniuk, F., Mehler, M., Pellicer, A., Tzall, S., LaBadie, G., Hobart, C., Ellenbogen, A. and Hirschhorn, R. (1986) *Proc Natl Acad Sci USA* 83:9641-9644.
42. Schuchman, E. H., Jackson, C. E. and Desnick, R. J. (1990) *Genomics* 6:149-158.
43. vonHeijne, G. (1986) *Nucleic Acids Res* 14:4683-4690.
44. Bishop, D. F., Kornreich, R. and Desnick, R. J. (1988) *Proc Natl Acad Sci USA* 85:3903-3907.
45. Proia, R. L. and Soravia, E. (1987) *J Biol Chem* 262:5677-5681.
46. Proia, R. L. (1988) *Proc Natl Acad Sci USA* 85:1883-1887.
47. Horowitz, M., Wilders, S., Horowitz, A., Reiner, O., Gelbart, T. and Beutler, E. (1989) *Genomics* 4:87-96.
48. Geier, C., von Figura, D. and Pohlmann, R. (1989) *Eur J Biochem* 183:611-616.
49. Kornreich, R., Desnick, R. J. and Bishop, D. F. (1989) *Nuc Acids Res* 14:3301-3302.
50. Myerowitz, R. and Costigan, F. C. (1988) *J Biol Chem* 263:18587-18589.
51. Myerowitz, R. (1988) *Proc Natl Acad Sci USA* 85:3955-3949.
52. Arpaia, E., Dumbrille-Ross, A., Maler, T., Neote, K., Tropak, M., Troxel, C., Stirling, J. L., Pitts, J. S., Bapat, B., Lamhonwah, A. M., Mahuran, D. J., Schuster, S. M., Clarker, J. T. R., Lowden, J. A. and Gravel, R. A. (1988) *Nature* 333:85-86.
53. Ohno, K. and Suzuki, K. (1988) *Biochem Biophys Res Comm* 153:463-469.
54. Navon, R. and Proira, R. L. (1989) *Science* 243:1471-1474.
55. Paw, B. H., Kaback, M. M. and Neufeld, E. F. (1989) *Proc Natl Acad Sci USA* 86:2413-2417.

56. D-Azzo, A., Proira, R. L., Kolodny, E. H., Kaback, M. M. and Neufeld, E. F. (1984) *J Biol Chem* 259:11070-11074.
57. Myerowitz, R. and Hogikyan, N. D. (1985) *Science* 232:1646-1648.
58. Myerowitz, R. and Hogikyan, N. D. (1987) *J Biol Chem* 262:15396-15899.
59. Proira, R. L. and Neufeld, E. F. (1982) *Proc Natl Acad Sci USA* 79:6360-6364.
60. Nakano, T., Muscillo, M., Ohno, K., Hoffman, A. J. and Suzuki, K. (1988) *J Neurochem* 51:984-987.
61. Lau, M. M. and Neufeld, E. F. (1988) *J Cell Biol* 107:342A.
62. Ohno, K. and Suzuki, K. (1988) *J Neurochem* 50:316-318.
63. Kornreich, R., Bishop, D. F. and Desnick, R. J. (in press) *J Biol Chem*.
64. Prokop, D. J. (1984) *Am J Hum Genet* 36:499-505.
65. Oshima, A., Kyle, J. W., Miller, R. D., Hoffman, J. W., Powell, P. P., Grubb, J. H., Sly, W. S. Tropak, M., Guise, K. S. and Gravel, R. A. (1987) 84:685-589.
66. Chamberlain, M. and Ryan, T. (1982) In *The Enzymes*, Vol XV. Boyer, P. Editor. Academic Press, NY.
67. Saiki, R. K., Gelfand, D. H., Stoffel, S., Scharf, S. J., Higuchi, R., Horn, G. T., Mullis, K. B. and Erlich, H. A. (1987) *Science* 239:487-491.
68. Brown, C. A., Neote, K., Leung, A., Gravel, R. A. and Mahuran, D. J. (1989) *J Biol Chem* 264:21705-21710.
69. Maeda, N. and Smithies, O. (1986) *Ann Rev Genet* 20:81-108.
70. Perler, F., Efstradiatis, A., Lomeico, P., Gilbert, W., Kolodner, R. and Dodgson, J. (1980) *Cell* 20:555-566.
71. Efstradiatis, A., Posakony, J. W., Maniatis, T., Lawn, R. M., O'Connell, C., Spritz, R. A., DeRiel, J. K., Forget, B. G., Weissman, S. M., Slightom, J. L., Blechel, A. E., Smithies, O., Baralle, F. E., Shoulders, C. C. and Proudfoot, N. J. (1980) *Cell* 21:653-558.
72. Baralle, F. E., Shoulders, C. C. and Proudfoot, N. J. (1980) *Cell* 21:621-26.
73. Slightom, J. L., Blechel, A. E. and Smithies, O. (1980) *Cell* 21:627-38.
74. Spritz, R. A., DeRiel, J. K., Forget, B. G. and Weissman, S. (1980) *Cell* 21:639-46.

75. Lawn, R. M., Efstradiatis, A., O'Connell, C. and Maniatis, T. (1980) *Cell* 21:647-651.
76. Mori, N., Singer-Sam, J. and Riggs, A. D. (1986) *FEBS Letters* 204:313-317.
77. McDonald, J. D., Lin, F. and Goldwasser, E. (1986) *Mol Cell Biol* 6:842-848.
78. Hoefsloot, L. H., Hoogeveen-Westerveld, M., Kroos, M. A., Vanbeeumen, J., Reuser, A. J. J. and Oostra, B. A. (1988) *EMBO Journal* 7:1697-1704.
79. Peters, C., Schmidt, B., Rommerskrich, W., Rupp, K., Zuhlsdorf, M., Vingron, M., Meyer, H. E., Pohlmann, R. and von Figura, K. (1990) *J Biol Chem* 265:3374-3381.
80. Gilbert, W. (1978) *Nature*, 271:501.
81. Michelson, A. M., Blake, C. C., Evans, S. T. and Orkin, S. H. (1985) *Proc Natl Acad Sci USA* 82:6965-6969.
82. van Diggelen, O. P., Schindler, D., Kleijer, W. J., Huijmans, J. M. G., Galjaard, H., Linden, H. U., Peter-Katalinic, J., Egge, H., Dabrowski, U., Cantz, M. (1987) *Lancet* ii:804.
83. van Diggelen, O. P., Schindler, D., Willemson, R., Boer, M., Kleijer, W. J., Huijman, J. G. M., Blom, W. and Galjaard, H. (1988) *J Inher Met Dis* 11:349-357.
84. Schindler, D., Bishop, D. F., Wolfe, D. E., Wang, A. M., Egge, H., Lemieux, R. U. and Desnick, R. J. (1989) *N Engl J Med* 320:1735-1740.
85. Kanzaki, T., Yokota, M., Mizuno, N., Matsumoto, Y., and Hirabayashi, Y. (1989) *Lancet* April 22, 1989:875-877.
86. Wang, A. M., Kanzaki, T., Schindler, D. and Desnick, R. J. (1989) *Am J Hum Genet* 43:4 A228.
87. Linden, H. U., Klein, R. A., Egge, H., Peter-Katalinic, J., Dabrowski, J. and Schindler, D. (1989) *Biol Chem Hoppe-Seyer* 370:661-672.
88. Wolfe, D. E., Schindler, D., Desnick, R. J. and Perl, D. (1989) *J Neuropathol Exp Neurol* 48:349.
89. Seitelberger, F. (1986) In: *Handbook of Clinical Neurology* Vol. 49. Elsevier, Amsterdam. p. 391-415.
90. Hirabayashi, Y., Matsumoto, Y., Matsumoto, M., Toida, T., Iida, N., Matsubara, T., Kanzaki, T., Yokota, M. and Ishizuka, I. (1990) *J Bio Chem* 265:1693-1701.

## OBJECTIVES AND SIGNIFICANCE

The overall objectives of this research are to 1) determine the molecular genetics of  $\alpha$ -*N*-acetylgalactosaminidase ( $\alpha$ -GalNAc), and 2) elucidate the nature of the molecular lesions in Schindler disease ( $\alpha$ -GalNAc deficiency). To achieve these objectives, each of the following were successfully accomplished: 1) Purification of homogeneous human  $\alpha$ -GalNAc enzyme for amino acid microsequencing of *N*-terminal and tryptic peptides, and for quantitative tryptic and chymotryptic mapping; 2) Selection and synthesis of mixed oligonucleotide probes corresponding to amino acid sequences with minimal codon redundancy; 3) Isolation of a full-length cDNA encoding human  $\alpha$ -GalNAc and characterization of the sequence, including 5' and 3' untranslated regions, pre- and prepro-peptide sequences, determination of the transcriptional initiation site(s), and analysis of the  $\alpha$ -GalNAc transcripts; 4) Transient expression of the full-length  $\alpha$ -GalNAc cDNA to demonstrate its functional integrity; 5) Isolation of genomic sequences for  $\alpha$ -GalNAc and characterization of the genomic structure including intron/exon organization, 5' regulatory elements, and 3' flanking region. Both the cDNA and genomic sequences of  $\alpha$ -GalNAc were compared to  $\alpha$ -Gal A for identification of evolutionary relationships; and 6) Characterization of the molecular defects in Schindler disease Types I and II.

These studies relate directly to our understanding of lysosomal enzymes and provide insight into their function and possible evolution. The isolation and characterization of the cDNA and genomic sequences encoding  $\alpha$ -GalNAc and their comparison with those of  $\alpha$ -Gal A should provide insight into the potential structural and functional domains of these enzymes as well as into their evolutionary relationships. These findings along with those of other lysosomal enzymes known to share homology (e.g.,  $\beta$ -hexosaminidase  $\alpha$ - and  $\beta$ -chains, and

arylsulfatases A and B) may lead to the identification of potential "lysosomal domains" which may be conserved among lysosomal enzymes.

In addition, these studies are related to the most recently identified lysosomal disorder, Schindler disease. The Type I disorder is a new form of infantile neuroaxonal dystrophy and is the first neuroaxonal dystrophy in which the biochemical basis is known. The Type II disorder appears to be a new cause of angiokeratoma corporis diffusum. The isolation and characterization of the  $\alpha$ -GalNAc cDNA and genomic sequences has allowed the determination of the molecular lesions underlying these two subtypes in the first individuals identified with the Schindler disease subtypes.

These studies have provided the molecular basis for the further study of human  $\alpha$ -GalNAc and lysosomal enzymes in general. The characterization of the  $\alpha$ -GalNAc cDNA and genomic sequences, as well as the identification of the specific base substitutions resulting in two distinct clinical phenotypes due to  $\alpha$ -GalNAc deficiency, have provided the framework toward further delineating the function and specificity of  $\alpha$ -GalNAc in man and the pathology of its deficiency.

**CHAPTER ONE:**

**HUMAN  $\alpha$ -*N*-ACETYLGALACTOSAMINIDASE: MOLECULAR CLONING,  
NUCLEOTIDE SEQUENCE, AND EXPRESSION OF A FULL-LENGTH cDNA**

**Homology with Human  $\alpha$ -Galactosidase A Suggests Evolution**

**From a Common Ancestral Gene**

## SUMMARY

Human  $\alpha$ -*N*-acetylgalactosaminidase ( $\alpha$ -GalNAc; E.C. 3.2.1.49), the lysosomal glycohydrolase that cleaves  $\alpha$ -*N*-acetylgalactosaminyl moieties from glycoconjugates, is encoded by a gene localized to chromosome 22q13→qter. The deficient activity of  $\alpha$ -GalNAc is the enzymatic defect in Schindler disease, an inherited neuroaxonal dystrophy. To isolate a full-length cDNA, the enzyme from human lung was purified to homogeneity, 129 non-overlapping amino acids were determined by microsequencing the *N*-terminus and seven tryptic peptides, and four synthetic oligonucleotide mixtures were used to screen a human fibroblast cDNA library. A full-length cDNA, pAGB-3, isolated from a placental  $\lambda$ gt11 cDNA library, had a 2158 bp insert with an open reading frame which predicted an amino acid sequence that was colinear with all 129 microsequenced residues of the purified enzyme. The pAGB-3 insert had a 344 bp 5' untranslated region, a 1236 bp open reading frame encoding 411 amino acids, a 514 bp 3' untranslated region, and a 64 bp poly(A) tract. A signal peptide sequence of 17 amino acids as well as six *N*-glycosylation sites were predicted. The pAGB-3 cDNA was subcloned into the p91023(B) mammalian expression vector and human  $\alpha$ -GalNAc activity was transiently expressed in COS-1 cells, demonstrating the functional integrity of the full-length cDNA. Northern hybridization analysis of mRNA revealed two transcripts of about 3.6 and 2.2 kb, and primer extension studies indicated a cap site at nt -347. A second full-length cDNA, pAGB-35, was isolated from a human retinal library and corresponded to the 3.6 kb transcript. The pAGB-35 cDNA insert was 3597 bp and had 468 bp and 1893 bp 5' and 3' untranslated sequences, respectively. Isolation of a genomic clone, gAGB-1, and sequencing the 2048 bp region including pAGB-3 revealed an 1754 bp intron between codons 319 and 320, which also was the site of a 70 bp insertion and a 45 bp deletion in other cDNA clones. Notably, the  $\alpha$ -GalNAc cDNA had remarkable amino acid homology with the human  $\alpha$ -galactosidase A ( $\alpha$ -Gal A)

cDNA suggesting the evolutionary relatedness of these genes. The  $\alpha$ -GalNAc cDNA had 46.9 to 64.7% amino acid identity in sequences (codons 1-319) corresponding to  $\alpha$ -Gal A exons 1 through 6, while the comparable exon 7 sequence (pAGB-3 codons 320-411) had only 15.8% homology with numerous gaps. These findings implicate the genomic region at and surrounding codon 319 as a potential site for the abnormal processing of  $\alpha$ -GalNAc transcripts as well as for a recombinational event in the evolution and divergence of  $\alpha$ -Gal A and  $\alpha$ -GalNAc. The availability of the full-length cDNA for human  $\alpha$ -GalNAc will permit studies of the genomic organization and evolution of this lysosomal gene, as well as the characterization of the molecular lesions causing Schindler disease.

## INTRODUCTION

In the early 1970's, several investigators demonstrated the existence of two  $\alpha$ -galactosidase isozymes, designated A and B, which hydrolyzed the  $\alpha$ -galactosidic linkages in 4-MU- and/or p-NP- $\alpha$ -D-galactopyranosides (1-7). In tissues, about 80-90% of total  $\alpha$ -galactosidase ( $\alpha$ -Gal) activity was due to a thermolabile, myoinositol-inhibitable  $\alpha$ -Gal A isozyme, while a relatively thermostable isozyme,  $\alpha$ -Gal B, accounted for the remainder. The two "isozymes" were separable by electrophoresis, isoelectric focusing, and ion exchange chromatography. After neuraminidase treatment, the electrophoretic migrations and pI values of  $\alpha$ -Gal A and B were very similar (1), initially suggesting that the two enzymes were the differentially glycosylated products of the same gene. The finding that the purified glycoprotein enzymes had similar physical properties including subunit molecular weight (~46 kDa), homodimeric structures, and amino acid compositions also indicated their structural relatedness (8-14). However, the subsequent demonstration that polyclonal antibodies against  $\alpha$ -Gal A or B did not cross-react with the other enzyme (8, 11), that only  $\alpha$ -Gal A activity was deficient in hemizygotes with Fabry disease (1-8) and that the genes for  $\alpha$ -Gal A and B mapped to different chromosomes (7, 15), clearly demonstrated that these enzymes were genetically distinct. Thus, it was not surprising when  $\alpha$ -Gal B was shown in 1977 to be an  $\alpha$ -N-acetylgalactosaminidase ( $\alpha$ -GalNAc), a homodimeric glycoprotein which hydrolyzed artificial and natural substrates with terminal  $\alpha$ -N-acetylgalactosaminyl moieties (10, 11, 14) including various O- and N-linked glycopeptides and glycoproteins, glycosphingolipids, and the proteoglycan, cartilage keratin sulfate II (7).

Purified  $\alpha$ -GalNAc has reported native and subunit molecular weights of 90 to 117 kDa and 46 to 48 kDa, respectively (8-14). Kinetic studies demonstrated that the enzyme was inhibited by  $\alpha$ -N-acetylgalactosamine ( $K_i \sim 2.1$  mM) and hydrolyzed synthetic substrates with

either terminal  $\alpha$ -*N*-acetylgalactosaminide ( $K_m \sim 1$ -2 mM) or  $\alpha$ -*D*-galactoside moieties ( $K_m \sim 7$ -10 mM) (8-14). Biosynthetic studies performed with cultured fibroblasts indicated that the human enzyme was synthesized as a 65 kDa glycosylated precursor which was processed to a mature 48 kDa lysosomal form; both the precursor and mature forms had high-mannose type oligosaccharide chains, but only the precursor's mannose residues were phosphorylated (16).

The deficient activity of  $\alpha$ -GalNAc was demonstrated in two brothers with Schindler disease (17, 18), a newly recognized form of infantile neuroaxonal dystrophy (18). The affected brothers excreted increased amounts of *O*-linked glycopeptides and oligosaccharides containing  $\alpha$ -*N*-acetylgalactosaminyl moieties which were detectable in urinary screening profiles (17, 19). Biochemical and immunologic studies revealed that neither  $\alpha$ -GalNAc activity or enzyme protein was present in fibroblast lysates from the affected sibs (18). Thus, efforts were undertaken to isolate and express a full-length  $\alpha$ -GalNAc cDNA in order to determine the nature of the molecular lesions in patients with Schindler disease, and to characterize the genomic organization and expression of the human gene encoding this lysosomal hydrolase.

While expression studies of a hybrid  $\alpha$ -GalNAc sequence were in progress, Tsuji *et al.* reported the isolation of a human  $\alpha$ -GalNAc cDNA (20). Unlike the full-length pAGB-3  $\alpha$ -GalNAc cDNA sequence reported here, their clone, pcD-HS1204, contained a 70 bp insertion after pAGB-3 nt 957 which altered the reading frame for pAGB-3 residues 330 to 411 and resulted in a truncated polypeptide of only 358 residues. Although their predicted amino acid sequence did not include our tryptic peptide containing residues 335 to 344, the 70 bp insertion may have resulted from alternative splicing. Thus, efforts were directed to determine if such an alternatively spliced transcript occurs in man. In this communication, we report the isolation, nucleotide sequence and transient expression of a full-length cDNA encoding  $\alpha$ -GalNAc. Genomic sequencing did not reveal the presence of the putative 70 bp insertion, thereby

affirming that the expressible pAGB-3 transcript is authentic. In addition, remarkable homology between the predicted  $\alpha$ -GalNAc and  $\alpha$ -Gal A amino acid sequences was identified, suggesting the evolutionary relatedness of the autosomal and X-linked genes encoding these lysosomal hydrolases.

## EXPERIMENTAL PROCEDURES

*Affinity Purification, Microsequencing and Antibody Production.* Human lung  $\alpha$ -GalNAc was purified to homogeneity, polyclonal rabbit anti-human  $\alpha$ -GalNAc antibodies were produced and purified, and cell supernatants were immunoblotted as described previously (18, 21, 22). For the isolation of tryptic peptides, the concentrated post-hydroxylapatite fraction was subjected to preparative 10% NaDodSO<sub>4</sub>/PAGE and the 48 and 117 kDa  $\alpha$ -GalNAc species were electroeluted separately (23), digested with tosylphenylalaninechloromethylketone-treated trypsin, and the resulting tryptic peptides from each species were separated by C8 reversed-phase HPLC (24). The *N*-terminal amino acid sequences of both the 48 and 117 kDa species and the sequences of selected tryptic peptides from the 48 kDa species were determined by automated gas-phase microsequencing and HPLC identification of phenylthiohydantoin amino acid derivatives (25).

*Construction of Synthetic Oligonucleotide Probes.* Mixed and unique oligonucleotides were synthesized on an Applied Biosystems Model 380B oligonucleotide synthesizer. Four oligonucleotide mixtures were constructed to regions of minimal codon redundancy for the following *N*-terminal and tryptic peptide sequences: *N*-terminus [5'-CA(AG)ACNCCNCCNATGGG-3']; peptide T-106A [AA(TC)AT(TCA)GA(TC)GA(TC)-TG(TC)TGGAT(TCA)GGNGG-3']; peptide T-72 [5'-ACNTT(TC)GCNGA(AG)TGGAA-3']; and peptide T-133, [5'-TGGCCNGCNTA(TC)GA(AG)GG-3']. Oligonucleotide probes for library screening were 5' end-labelled with [ $\gamma$ -<sup>32</sup>P]ATP using T4 polynucleotide kinase (26). Unique sequence oligonucleotides (17-mers) were synthesized and used as primers in sequencing reactions. To determine the cap site, two unique, overlapping 30-mers were synthesized for primer extension, 5'-TCGGGACTCCCAGCACTGCAGAGGGTGTGA-3' and 5'-CTG CAGAGGGTGTGAGGTCTGACATCCAGG-3'. To detect alternatively spliced transcripts,

PCR sense and antisense primers for the exonic region flanking the putative 70 bp insertion had the sequences 5'-AGTCGAATTCTGATGTCCACAGACCTGCGT-3', and 5'-AGTCGTCGAGCATATCGGTCCTGCAGCTGA-3', respectively. The four PCR primer sequences for the construction of the  $\alpha$ -Gal A and  $\alpha$ -GalNAc hybrid cDNA were  $\alpha$ -Gal A sense, 5'-TGGGGAGTAGATCTGCTAAAA-3';  $\alpha$ -Gal A antisense, 5'-GATGAGAGATTTTTT-CCTGTCTAAGCTGGTACCC-3';  $\alpha$ -GalNAc sense, 5'-TACCAGCTTAGACAGGAAAAA-TCTCTCATCGAA-3'; and  $\alpha$ -GalNAc antisense, 5'-AAGAGGTCAGATCTCTCTACT-3'.

*Isolation and Characterization of cDNA and Genomic Clones.* The pcD human fibroblast cDNA library, kindly provided by Dr. Hiroto Okayama (NIH), was screened with the radiolabelled 26-mer oligonucleotide mixture corresponding to tryptic peptide T-106A by colony hybridization (24). Plasmid DNA isolation and Southern hybridization analyses of positive clones were performed as previously described (26). A  $\lambda$ gt10 human retinal library, kindly provided by Dr. Jeremy Nathans (Johns Hopkins University) was screened with the radiolabelled pAGB-3 cDNA insert using the conditions previously described. For isolation of a full-length cDNA, a 0.9 kb *Bam*HI fragment corresponding to the 5' portion of the pAGB-1 insert was then isolated, nick-translated, and used to screen recombinants from a  $\lambda$ gt11 human placental library (Clontech Laboratories, Palo Alto, CA) by plaque hybridization (26). To isolate genomic clones containing the entire  $\alpha$ -GalNAc sequence,  $1 \times 10^6$  recombinants from a human genomic cosmid library were screened with the radiolabelled pAGB-3 cDNA insert using the conditions described above for cDNA library screening. The genomic library was prepared from size-selected human lymphoblast DNA and kindly was provided by Dr. Henrik Vissing (Mount Sinai School of Medicine).

*DNA Sequencing and Computer-Assisted Analyses.* The *Bam*HI inserts from pAGB-1, a *Eco*RI-*Bam*HI restriction fragment of pAGB-3 and the *Eco*RI insert from pAGB-35 were

subcloned into M13 mp18 and mp19. All DNA sequencing reactions were carried out by primer extension using either M13 universal primers or  $\alpha$ -GalNAc-specific synthetic oligonucleotide primers by the dideoxy method in both orientations (27). Searches for nucleotide and amino acid sequence similarity were carried out with the University of Wisconsin Genetics Computer Group Software (28). Computer-assisted RNA folding was performed with the PCFOLD program (29).

*Transient Expression Assays.* The human  $\alpha$ -GalNAc full-length pAGB-3 cDNA insert was subcloned into the p91023(B) eukaryotic expression vector (30), kindly provided by Dr. R. J. Kaufmann (Genetics Institute, Boston, MA). Plasmid DNA from the construct (designated p91-AGB-3) was purified and COS-1 monkey kidney cells were transfected with 10  $\mu$ g of the p91-AGB-3 plasmid DNA by calcium-phosphate precipitation (31). Cells were harvested at 24 hr intervals after transfection and assayed for  $\alpha$ -GalNAc activity as previously described (18). One unit (U) of enzymatic activity is equal to that amount of enzyme required to hydrolyze 1 nmol of 4-MU- $\alpha$ -GalNAc per hour. Protein concentrations were determined by the fluorescamine method (21).

*Northern Hybridization and Cap Site Analyses.* Total RNA was isolated from human lymphoblasts, fibroblasts, and placentae and northern hybridization was performed using the nick-translated pAGB-3 insert as probe (26). Alternatively, the pAGB-3 insert was subcloned into pGEM-4Z (Promega, Madison, WI) and radiolabelled  $\alpha$ -GalNAc riboprobe, rAGB-3, was generated using the Promega riboprobe system and used for northern hybridization. For identification of the  $\alpha$ -GalNAc cap-site, two unique, overlapping 30-mer oligonucleotide primers were synthesized corresponding to regions 60 and 75 bp from the 5' end of the pAGB-3 cDNA and end-labelled (26). Each primer (100 ng) was used to extend 10  $\mu$ g of total placental RNA with the BRL cDNA Synthesis Kit (BRL, Gaithersburg, MD). First-strand synthesis was

terminated by phenol extraction and ethanol precipitation. The pellet was washed three times with 70% ethanol, resuspended in 6  $\mu$ l of H<sub>2</sub>O, and then mixed with 6  $\mu$ l of loading dye (0.3% xylene cyanol, 0.3% bromophenol blue, 0.37% EDTA, pH 7.0). The RNA/DNA heteroduplexes were denatured at 65 °C for 3 min and an aliquot was electrophoresed on a standard 8 M urea, 8% polyacrylamide sequencing gel.

*Construction of p91 $\alpha$ -GalA6/ $\alpha$ -GalNAc7.* A plasmid containing  $\alpha$ -Gal A exons 1 through 6 from pcDAG-126 (32) was fused to the 3' region of pAGB-3  $\alpha$ -GalNAc insert which corresponded in position to  $\alpha$ -Gal A exon 7. The hybrid cDNA, designated p $\alpha$ -GalA6/ $\alpha$ -GalNAc7 was constructed with the sense and antisense primers indicated above using a PCR-based method (33) and sequenced. The p $\alpha$ -GalA6/ $\alpha$ -GalNAc7 insert was subcloned into the expression vector, p91023(B), and the construct was transiently expressed in COS-1 cells as described above. The  $\alpha$ -Gal A and  $\alpha$ -GalNAc enzymatic activities and enzyme proteins were detected with 4-MU substrates and by immunoblotting with the respective polyclonal antibodies as described above.

*Primer Extension and PCR Amplification of cDNA and Genomic Sequences.* For PCR amplification of the putative alternatively spliced region, the 30-mer sense and antisense primers (described above) were used to amplify the 1) reverse-transcribed mRNA from various human sources, 2) cDNA inserts from clones pAGB-4 to 34, and 3) the gAGB-1 genomic sequence. DNAs from pAGB-4 to 34 cDNA clones and the gAGB-1 genomic clone were isolated as described (24, 26). cDNA was synthesized from 10  $\mu$ g of lymphoblast, fibroblast, and placental total RNA or 2.5  $\mu$ g of brain poly (A)<sup>+</sup> mRNA (Clontech, Palo Alto, CA) using the BRL cDNA Synthesis Kit. Bacteriophage DNA (~0.1  $\mu$ g) and reverse-transcribed mRNA (~0.1  $\mu$ g) or genomic cosmid DNA (~ 1  $\mu$ g) was PCR-amplified using 20  $\mu$ M of each primer and the GeneAmp DNA Amplification Reagent Kit (Perkin Elmer Cetus, Norwalk, CT). Each PCR

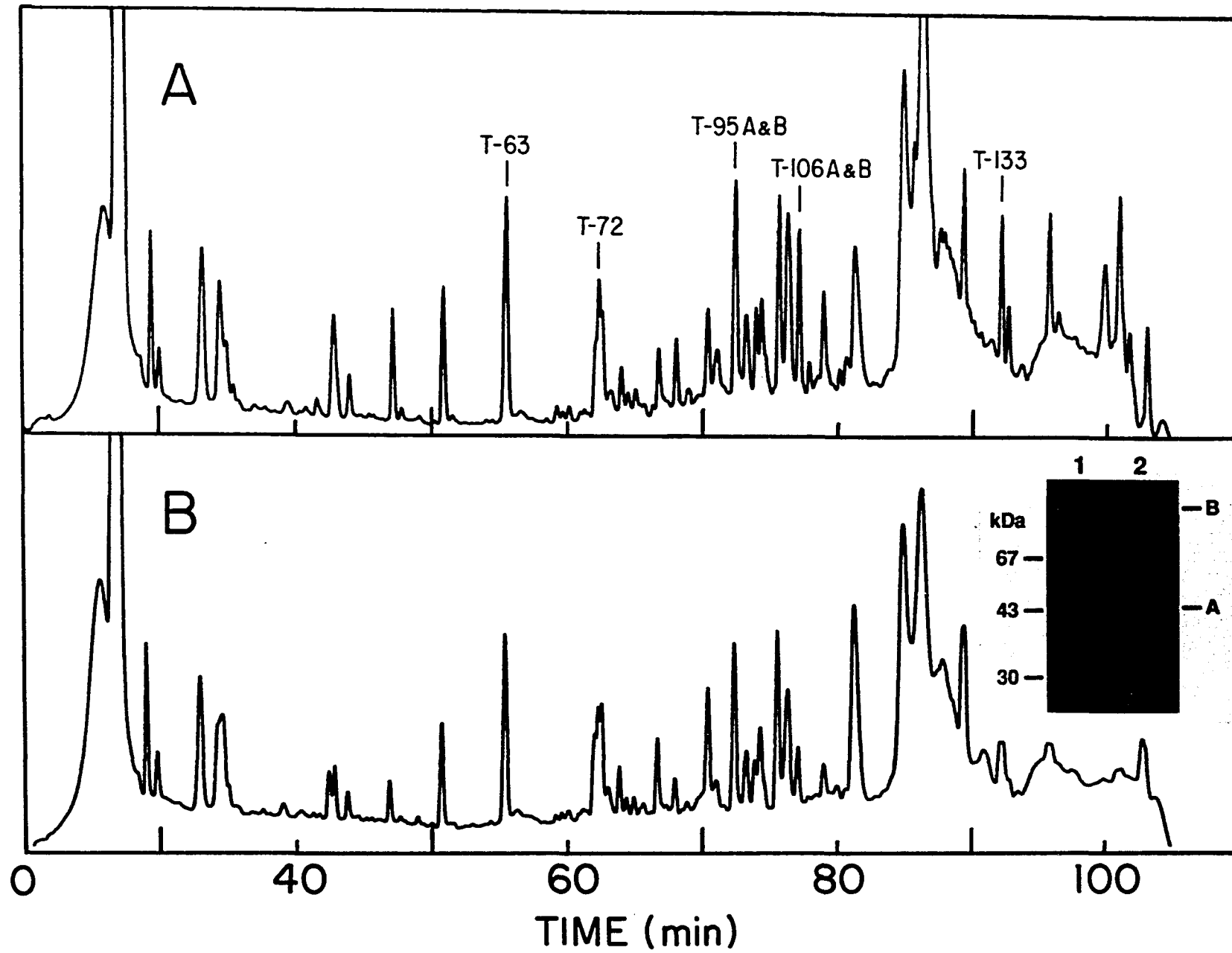
cycle consisted of 1 min denaturation at 94 °C; 2 min annealing at 37 °C; and a 7 min extension at 60 °C. The PCR products were phenol extracted, ethanol precipitated and resuspended in 20 µl of H<sub>2</sub>O. An aliquot (2 µl) of each PCR reaction was analyzed by electrophoresis on agarose gels using *Hind*III digested lambda and *Hae*III digested φX174 DNAs as size standards. For identification of potential stops during reverse transcription of the region surrounding the pcD-HS1204 insertion, a unique 32-mer, 5'-AGTAGTAAGCTTTCATATATCACAGACCCGGT-3', was used to extend 10 µg of total placental RNA or 1 µg of rAGB-3 generated *in vitro* by the Promega riboprobe system as described above.

## RESULTS AND DISCUSSION

*Purification and Characterization of Human  $\alpha$ -GalNAc.* Human  $\alpha$ -GalNAc was purified to homogeneity (specific activity =  $\sim 370,000$  U/mg protein) as assessed by the presence of only the 48 and 117 kDa species on NaDodSO<sub>4</sub>/PAGE (Fig. 1, inset). The 117 kDa was not reduced by boiling or by dialysis against 8 M urea in the presence of  $\beta$ -mercaptoethanol. The 27 microsequenced *N*-terminal residues of the electroeluted 117 kDa species were identical to those of the 48 kDa species. Further evidence that the 117 kDa species was a homodimer of the 48 kDa glycoprotein subunit was the finding that the tryptic digests (and chymotryptic digests; not shown) of both species had essentially identical HPLC profiles (Fig. 1). Microsequencing of the *N*-terminus and seven tryptic peptides from the 48 kDa species identified a total of 129 non-overlapping  $\alpha$ -GalNAc residues. For library screenings, synthetic oligonucleotide mixtures (17- to 26-mers) were constructed to contain all possible codons for selected amino acid sequences from the *N*-terminus and three internal tryptic peptides (Figs. 1 and 2).

*Isolation, Characterization and Expression of a Full-Length cDNA.* Screening of  $2 \times 10^6$  recombinants from the pcD human fibroblast cDNA library with a 26-mer oligonucleotide mixture of 576 species corresponding to internal peptide T-106A detected two putative positive clones. pAGB-1, which hybridized with all four oligonucleotide mixtures, had a 1.8 kb insert with an open reading frame of 1242 bp, a 514 bp 3' untranslated region, and a poly(A) tract, but no apparent 5' untranslated sequence. Authenticity was established by colinearity of the pAGB-1 insert's predicted amino acid sequence with 129 microsequenced residues of the purified protein. In order to isolate a full-length cDNA, the 0.9 kb 5' *Bam*HI fragment from the pAGB-1 insert was radiolabelled and used to screen a human placental cDNA library. Of 32 putative positive clones (pAGB-3 to 34), pAGB-3 contained the longest insert and was sequenced in both

*Figure 1.* Reversed-phase HPLC separation of tryptic peptides from electroeluted 117 kDa (A) and 48 kDa (B) species of purified human  $\alpha$ -GalNAc. The indicated peptides were microsequenced. (Inset) NaDodSO<sub>4</sub>/PAGE of purified  $\alpha$ -GalNAc. See text for details.



orientations. As shown in Fig. 2, the 2158 bp pAGB-3 insert had a 344 bp 5' untranslated region, a 1236 bp open reading frame which encoded 411 amino acids, a 514 bp 3' untranslated region and a 64 bp poly(A) tract. An upstream, inframe ATG occurred at -192 nt, but there were inframe termination codons at -141, -135, and -120 nt, indicating that the -192 ATG was non-functional. A single consensus polyadenylation signal (AATAAA) and a consensus recognition sequence (CACTG) for the U4 small nuclear ribonucleoprotein (34) were located 16 and 65 bp from the poly(A) tract, respectively. In retrospect, the partial cDNA, pAGB-1 had the entire 1236 bp coding region as well as 6 bp of 5' untranslated sequence.

Analysis of the deduced amino acid sequence of pAGB-3 indicated a signal peptide sequence of 17 residues since Leu-18 was the *N*-terminal residue of the microsequenced mature enzyme. When the weight-matrix method of von Heijne (35) was used to predict the peptidase cleavage site, the preferred site, between Ala-13 and Gln-14, had a score of 4.34, whereas cleavage after Met-17 had a score of 2.38. The predicted molecular mass of the 394 residue mature, unglycosylated enzyme subunit ( $M_r=44,700$ ) was consistent with that (48 kDa) estimated by NaDodSO<sub>4</sub>/PAGE of the purified glycosylated enzyme. These findings suggest that the mature glycoprotein subunit had at least two *N*-linked oligosaccharide chains, although there were six putative *N*-glycosylation sites at Asn residues 124, 177, 201, 359, 385 and 391 (Fig. 2).

For transient expression, the pAGB-3 full-length cDNA insert was subcloned into the eukaryotic expression vector p91023(B) and the construct, p91-AGB-3, was transfected into COS-1 monkey kidney cells. Compared to the endogenous mean  $\alpha$ -GalNAc activity in mock transfected COS-1 cells (35 U/mg; range: 23-50 U/mg; n=6), the transfected cells had a mean activity of 600 U/mg (range: 104-2,400 U/mg; n=6) 72 hr after transfection, or about 17 times the endogenous activity. The expressed human enzyme protein also was detected by

*Figure 2.* Nucleotide and predicted amino acid sequences of the pAGB-3 and pAGB-35 cDNA inserts containing the complete coding region for human  $\alpha$ -N-GalNAc. The A of the initiation ATG is nt +1 and the N-terminal Met of the signal peptide is amino acid 1. Bold underlines indicate colinear amino acid sequence obtained by microsequencing the N-terminal (N-Ter) and tryptic peptides (T) of the purified enzyme. CHO indicates potential sites of N-glycosylation. Overlines indicate the polyadenylation signals (AATAAA) and the pentanucleotide sequence (CACTG) recognized by the U4 small nuclear ribonucleoprotein.

-469 CCCTTCGGTACCATAGTCCGGTGGGAGGGGGCTCTCGGTTTGGGGACCCAGGGGGGAACCCGGACCGGGCTGGAAGTCCGGAGCCGCCGCCAGCCCCCCCCCTCCGCT -358  
-357 CTTTCTTGGTGACCTTAAGCCAGTGGCTGCCCTTTTCTGAGCCCGGGCGGGCCGAAGGGCCCGTAGGCCCTCGGGACTCCAGCACTGCAGAGGGTGTGAGGTCTGACATCCAAGAC -239  
-238 ACGTGTGTTTGGTATTTCGAAGGAAGAATCAAGCTCCGGGAAGTATGGCTGGGGATGGGGGGGCAACTTGGGACCGAGTGTACGATCCAGCCCTAAGGTTGAGGGGGGGCGAGCT -120  
-119 AGCCAGCCAGCCGTGACCCAGTGCCTTTTCAGAGCTTTCTTAGCTTCCAGACCCCAACACATACAGCTGATACCCGACAGCCAGATCTGGTCAGGTCTCCGAAGCTGAGTCCAGAGCC - 1

1 ATG CTG CTG AAG ACA GTG CTC TTG CTG GGA CAT GTG GCC CAG GTG CTG ATG CTG GAC AAT GGG CTC CTG CAG ACA CCA CCC ATG GGC TGG 90  
1 Met Leu Leu Lys Thr Val Leu Leu Leu Gly His Val Ala Gln Val Leu Met Leu Asp Asn Gly Leu Leu Gln Thr Pro Pro Met Gly Trp 30  
M-Ter

91 CTG GCC TGG GAA CCC TTC CGC TCC AAC ATT AAC TGT GAT GAG GAC CCA AAG AAC TGC ATA AGT GAA CAG CTC TTC ATG GAG ATG GCT GAC 180  
31 Leu Ala Trp Glu Arg Phe Arg Cys Asn Ile Asn Cys Asp Glu Asp Pro Lys Asn Cys Ile Ser Glu Gln Leu Phe Met Glu Met Ala Asp 60  
T-106B

181 CGG ATG GCA CAG GAT GGA TGG CGG GAC ATG GCC TAC ACA TAC CTA AAC ATT GAT GAC TGC TGG ATC GGC GGT CGC GAT GCC AGT GGC CGC 270  
61 Arg Met Ala Gln Asp Gly Trp Arg Asp Met Gly Tyr Thr Tyr Leu Asn Ile Asp Asp Cys Trp Ile Gly Gly Arg Asp Ala Ser Gly Arg 90  
T-106A

271 CTG ATG CCA GAT CCC AAG CGC TTC CCT CAT GGC ATT CCT TTC CTG GGT GAC TAC GTT CAC TCC CTG GGC CTG AAG TTG GGT ATC TAC GGC 360  
91 Leu Met Pro Asp Pro Lys Arg Phe Pro His Gly Ile Pro Phe Leu Ala Asp Tyr Val His Ser Leu Gly Leu Lys Leu Gly Ile Tyr Ala 120

361 GAC ATG GGC AAC TTC ACC TGC ATG GGT TAC CCA GGC ACC ACA CTG GAC AAG GTG GTC CAG GAT GCT CAG ACC TTC GCC GAG TGG AAG GTA 450  
121 Asp Met Gly Asn Phe Thr Cys Met Gly Tyr Pro Gly Thr Thr Leu Asp Lys Val Val Gln Asp Ala Gln Thr Phe Ala Gly Trp Lys Val 150  
CBO

451 GAC ATG CTC AAG CTG GAT GGC TCC TTC TCC ACC CCC GAG GAG CGG CGC GAG GGC TAC CCC AAG ATG GCT GCT GCC CTG AAT GCC ACA GGC 540  
151 Asp Met Leu Lys Leu Asp Gly Cys Phe Ser Thr Pro Glu Glu Arg Ala Cln Gly Tyr Pro Lys Met Ala Ala Ala Leu Asn Ala Thr Gly 180  
CBO

541 GGC CCC ATC GCC TTC TCC TGC AGC TGG CCA GCC TAT GAA GGC GGC CTC CCC CCA AGG GTG AAC TAC AGT CTG CTG GCG GAC ATC TCC AAC 630  
181 Arg Pro Ile Ala Phe Ser Cys Ser Trp Pro Ala Tyr Glu Gly Gly Leu Pro Pro Arg Val Asn Tyr Ser Leu Leu Ala Asp Ile Cys Asn 210  
T-95B

631 CTC TGG CGT AAC TAT GAT GAC ATC CAG GAC TCC TGG TGG AGC GTG CTC TCC ATC CTG AAT TGG TTC GTG GAG CAC CAG GAC ATA CTG CAG 720  
211 Leu Trp Arg Asn Tyr Asp Asp Ile Gln Asp Ser Trp Trp Ser Val Leu Ser Ile Leu Asn Trp Phe Val Glu His Gln Asp Ile Leu Gln 240

721 CCA GTG GCC GGC CCT GGG CAC TGG AAT GAC CCT GAC ATG CTG CTC ATT GGG AAC TTT GGT CTC AGC TTA GAG CAA TCC CGG GCC CAG ATG 810  
241 Pro Val Ala Gly Pro Gly His Trp Asn Asp Pro Asp Met Leu Leu Ile Gly Asn Phe Gly Leu Ser Leu Glu Gln Ser Arg Ala Gln Met 270  
T-133

811 GCC CTG TGG ACG GTG CTG GCA GCC CCC CTC TTG ATG TCC ACA GAC CTG CGT ACC ATC TCC GCC CAG AAC ATG GAC ATT CTG CAG AAT CCA 900  
271 Ala Leu Trp Thr Val Leu Ala Ala Pro Leu Leu Met Ser Thr Asp Leu Arg Thr Ile Ser Ala Gln Asn Met Asp Ile Leu Gln Asn Pro 300

901 CTC ATG ATC AAA ATC AAC CAG GAT CCC TTA GCC ATC CAG GGA CCC AGG ATT CAC AAG GAA AAA TCT CTC ATC GAA GTG TAC ATG CGG CCT 990  
301 Leu Met Ile Lys Ile Asn Gln Asp Pro Leu Gly Ile Gln Gly Arg Arg Ile His Lys Glu Lys Ser Leu Ile Glu Val Tyr Met Arg Pro 330  
T-63

991 CTG TCC AAC AAG GCT AGC GCC TTA GTC TTC TTC AGC TGC AGG ACC GAT ATG CCT TAT CCC TAC CAC TCC TCC CTT GGC CAG CTG AAC TTC 1080  
331 Leu Ser Asn Lys Ala Ser Ala Leu Val Phe Phe Ser Cys Arg Thr Asp Met Pro Tyr Arg Tyr His Ser Ser Leu Gly Gln Leu Asn Phe 360  
T-95A

1081 ACC GGG TCT GTG ATA TAT GAG GCC CAG GAC GTC TAC TCA GGT GAC ATC ATC AGT GGC CTC CGA GAT GAA ACC AAC TTC ACA GTG ATC ATC 1170  
361 Thr Gly Ser Val Ile Tyr Glu Ala Gln Asp Val Tyr Ser Gly Leu Arg Asp Glu Thr Asn Phe Thr Val Ile Ile 390  
CBO

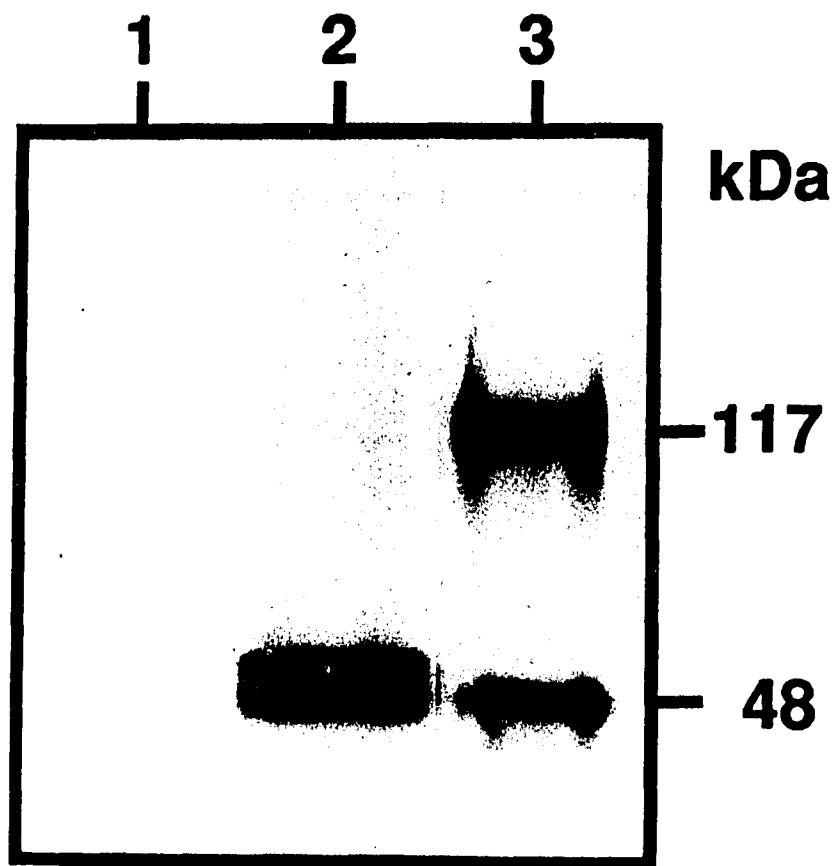
1171 AAC CCT TCA GGG GTA GTG ATG TGG TAC CTG TAT CCC ATC AAG AAC CTG GAG ATG TCC CAG CAG TGA GGAGCTGGGACATGTGACAGGCTGTGGTGGC 1267  
391 Asn Pro Ser Gly Val Val Met Trp Tyr Pro Ile Lys Asn Leu Asp Met Ser Gln Gln End 411  
CBO

1268 ACCACTGAGCTAGACCATGGAGCCCTGGCCATGCCAGGGCAAGTGGGGAGGTTCTCTGCTCCCCAGCCCTCTCGGTGACTGACCCCATCATACCCAAAGTGCAATCTCACGCCAGG 1386  
1387 TTCATAGCCCTGTCCAAAGGTTAAAGCCCTCTTGGAAACTCTTTTGGGGCAATTTTCCTGTGGCTTCTCGGCCCTTACTTCCATGTGGCCAGCCCCACAGAGCTTGGCTGAGCAACTGCC 1505  
1506 CAGCCCTCTGAGCTCCATGGCCATCAGGACTTAGCCCTTGCCTGTGACCTTGCCTTTGACCTGTGAAATCAGGATTTGGAAGTTTTOGAATTAGGAGTAGAGAGATCTGACCTCTTGGCCAGGAA 1624  
1625 GCGCATGCATGATGGCTTTTCTAACCATAGAGGGCCCTTGCAGCCCTGATACCCACTGGGAGTGGAGCTCACAAAGGAGACCTTGGCTCCCTCAGGTCACCAATAAACCTGTCTT 1743  
1744 TAATCAACTAGTGTAAATAGGACCTGGGGTACCTTTTAACTTAGGGACTCTGTGGGACTACAGGGGGTCTGGATTCAAAATCTCACCTCTATGGATGATGTCTGGCATGCTCTGGCCCC 1862  
1863 TCCCTTTTCCAAAGTACCAGAGATGTACTTCTAGAGATGCCAGACTTGCCTCACACCAGGAGCTGCCCTTGCCTGTGGCTGGTATGCTGGCAGTGGTATGGTGACAACACAATGA 1981  
1982 GCCCAGTAGCAATGCTCCAGGCTCCAGCTGTGCTCTGGCCAGCAGTCCCTGGCCAGCAGCTACGCTAGCCCTGCTTCTCTGCTGGTGGCTCCACCCCAACCCCTACTCTTTTCCTT 2100  
2101 TATTCTCTCCCAACTGGCTGCAAAAGGAGTTCATAATCCAGATGTGTCCACAACTCTCCAGGAGTGTGCTTAAAATGATTTCTGCTCCCTCTTCTCAAGGACAGGGTAAAC 2219  
2220 AGTCAAAGTGAAGCAATGATGGGTTCTGCCCTGTCTGCTTTGGTCCGATTTTCCCAAGCCCTCCCTCTTTTAAAGCCACTACAGGCGACTGTCTTCCCTGGTGGCTCTCCCTTTGG 2338  
2339 AGAGACTAGTGGCAGAGAACTAGCTAGGCTTCCCTGGCCCTTGAATGCCAGGACCCAGCTGCCCAAGCCCTTGCAGTGTCTGCCCAAGGAGTATTTCTGGTGGATCTCCCA 2457  
2458 CTCTCAGGGTCTGAGGTCTTGTATTTAGGGTATGGGTAAAAACTTACATATGACACTTTTGTACATCATTAAAGATGATGAAGCTTTACTGGTTAAGTACAGGGTCCACCAG 2576  
2577 AAAGTTCCTTTCCCCAGTGTATGTCTTGGGCTGTGAGAAACTTAGGACAGAGCCCTACTTATCCAGCCCTGGCCCTCCGTGCAATGCTCTGCCCTTTGGGCTTGTTCATAGTCT 2695  
2696 GTTCAGGCAAAAATGGAAATAGTCTCAGGAGGGAACCTTCTCTTCTAGCTCCAGGAGGCTTTGTTGGGAAATGAGTGTAGACCAATGGACTTGGGAGTGGGGGCCAGAGTGGAGCTG 2814  
2815 GACTTGACCAATTAATTCCTTACCTTTCAGCCCAAGATAGCTACATGCTTACTTCTGTATAAGTTCCTTTTGTCTGGGGGTGGCTGATGATCTATCTCTTCCCTTTCAT 2933  
2934 TCTTTCAACAAACATGCCATCTATCTCCAGGACATTTATTCCTTCTATTCAAGGCCAGTTAAGAAATCAATTTATTTTTTTCAGTTAAATACAGACTTGTTTGAGCCCAAGGTAC 2952  
3053 CCTTTCCTTTTTAAAACCTTTTATCATGAAAACCTTCAAAATCAAAAATAAAATGGTAGAGAGTGAACCTCCATG

immunoblot analysis using rabbit-anti-human  $\alpha$ -GalNAc antibodies, whereas the endogenous monkey enzyme was variably visible as a faint band at ~40 kDa or was not detectable (Fig. 3). The expressed human enzyme subunit had a molecular weight of ~48 kDa, indicating that it was glycosylated.

*Northern Hybridization and Cap-Site Analyses.* Northern hybridization analyses revealed two transcripts in total, cytoplasmic, or poly(A)<sup>+</sup> RNA of about 2.2 and 3.6 kb, which were present in similar amounts (not shown). The cap-site was determined to be at -347, or 3 nt beyond the 5' end of the pAGB-3 cDNA insert by primer extension of total placental RNA using two overlapping oligonucleotide probes. The 3.6 kb transcript was determined to be due to a second polyadenylation signal of AATAAA 1366 nt downstream of the first hexamer identified in the pAGB-3 cDNA. This 3.6 kb cDNA encoding  $\alpha$ -GalNAc was isolated from a human retinal library and was identical to the pAGB-3 cDNA sequence with the exception of an additional 125 bp and 1379 bp of 5' and 3' untranslated sequences respectively (Fig. 2). The pAGB-35 cDNA confirmed the alternative use of polyadenylation hexamers in producing two mRNA transcripts, and also suggests that the  $\alpha$ -GalNAc gene has multiple transcriptional start sites, and that the -347 cap-site is only one of at least two such sites.

*Sequence Homology Between  $\alpha$ -GalNAc with  $\alpha$ -Gal A.* Computer-assisted searches of nucleic acid and protein data bases revealed no significant amino acid sequence similarities between  $\alpha$ -GalNAc and that of any other DNA or protein sequence except for human  $\alpha$ -Gal A (32). Comparison of the nucleic acid and deduced amino acid sequences of the full-length  $\alpha$ -GalNAc and  $\alpha$ -Gal A cDNAs revealed 55.8% and 46.9% overall homology, respectively. Since the intron/exon junctions and the entire genomic sequence encoding human  $\alpha$ -Gal A have been determined (32, 36), it was possible to compare the  $\alpha$ -GalNAc amino acid sequence with those deduced from each of the seven  $\alpha$ -Gal A exons (Fig. 4). Notably, there was remarkable identity



*Figure 3.* Immunoblot of human  $\alpha$ -GalNAc expressed in COS-1 cells. Lanes: 1, mock-transfection; 2, p91-AGB-3 transfection; 3, purified human lung  $\alpha$ -GalNAc.

*Figure 4.* Alignment of amino acid sequences deduced from the full-length cDNAs encoding human  $\alpha$ -GalNAc ( $\alpha$ -Gal B),  $\alpha$ -Gal A, yeast *Mel 1*, and *E. coli Mel A*. Colons, identical residues; single dots, isofunctional amino acids; and boxes, identical or isofunctional residues in  $\alpha$ -GalNAc,  $\alpha$ -Gal A, *Mel 1* and/or *Mel A*. Gaps were introduced for optimal alignment. Numbered vertical lines indicate exon boundaries for  $\alpha$ -Gal A (32).



(56.4%) between the  $\alpha$ -GalNAc sequences corresponding to those of  $\alpha$ -Gal A exons 1 through 6. For example, all eight cysteine residues in  $\alpha$ -GalNAc were present in the identical positions in  $\alpha$ -Gal A. Of the 14 proline and 23 glycine residues in  $\alpha$ -Gal A, 10 and 20 were conserved in identical positions in  $\alpha$ -GalNAc, respectively. Furthermore, all four of the  $\alpha$ -Gal A *N*-glycosylation sites were conserved in  $\alpha$ -Gal B. Putative functional domains were suggested by shorter stretches of amino acid homology shared by  $\alpha$ -GalNAc,  $\alpha$ -Gal A, yeast  $\alpha$ -galactosidase (*Mel 1*) (37) and/or *E. coli*  $\alpha$ -galactosidase (*Mel A*) (38) in  $\alpha$ -Gal A exons 1 through 6. In contrast, there was little, if any, similarity in the predicted  $\alpha$ -GalNAc carboxy-terminal amino acid sequence after residue 319 which corresponded to  $\alpha$ -Gal A exon 7 (15.8% homology with numerous gaps). In addition, there were no significant similarities for the cDNAs encoding other human lysosomal polypeptides, with the exception of a short  $\alpha$ -GalNAc sequence (residues 365 to 371) in which six out of seven amino acids were identical to residues 194 to 200 in the  $\beta$ -hexosaminidase  $\alpha$ -chain, a lysosomal polypeptide with *N*-acetylgalactosaminidase specificity (39). These findings suggested that a cDNA construct containing  $\alpha$ -Gal A exons 1-6 joined to  $\alpha$ -GalNAc exon 7 might express a hybrid protein with both  $\alpha$ -Gal A and B activities. Therefore, a hybrid cDNA containing  $\alpha$ -Gal A exons 1 through 6 (nt -60-1029) and  $\alpha$ -GalNAc exon 7 (nt 958-1258) was constructed and expressed in COS-1 cells; however, neither immunoreactive protein nor enzymatic activity for either  $\alpha$ -Gal A or  $\alpha$ -GalNAc were detected.

The finding of extensive homology between  $\alpha$ -GalNAc and  $\alpha$ -Gal A suggested that they evolved by duplication and divergence of an ancestral sequence for  $\alpha$ -Gal A exons 1 through 6. Although there is little, if any, homology among the other lysosomal amino acid sequences (i.e., no "lysosomal domains"), there are notable examples of lysosomal enzyme subunits, pseudogenes or gene families which presumably evolved by duplication and divergence (e.g., 39-41). Future comparison of the  $\alpha$ -GalNAc and  $\alpha$ -Gal A intron/exon boundaries should

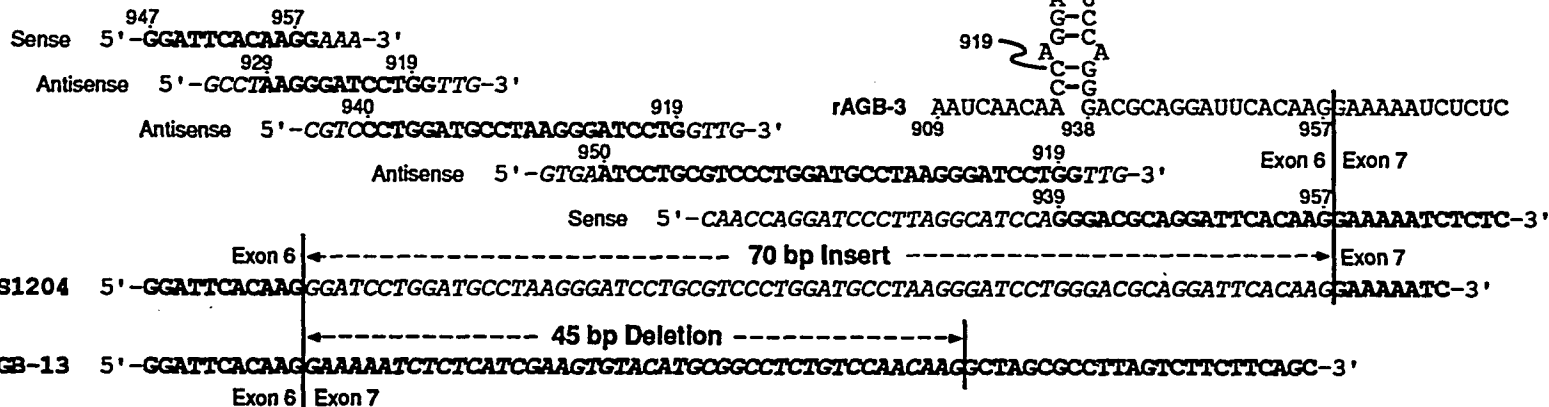
provide further information on the evolution of these lysosomal genes which encode structurally related, but functionally specific glycohydrolases.

*Primer Extension and PCR and Sequence Analyses of cDNA and Genomic Sequences.*

During the course of these studies, Tsuji *et al.* reported a similar human  $\alpha$ -GalNAc cDNA sequence (20) which differed from pAGB-3 by a 70 bp insertion after nt 957 (Fig. 5A) and by several substitutions (nt 493, 494, 524, 614, and 667). The 70 bp insertion consisted of three inverted repeats (nt 919-926, 919-936, and 919-944) and a direct repeat (nt 940-957) from the pAGB-3 coding sequence nt 919 to 957. Analysis of the pAGB-3 cDNA sequence from nt 760-1053 using an RNA folding program (29) predicted a stem and loop structure from nt 918 to 937 (Fig. 5A) which could stall or stop reverse transcription of the  $\alpha$ -GalNAc mRNA during cDNA synthesis. To determine if this secondary structure could cause cDNA synthesis errors in library construction, a 32-mer oligonucleotide primer was used to extend total placental RNA and  $\alpha$ -GalNAc transcripts generated *in vitro* with the riboprobe construct, rAGB-3. Stops of varying intensity were observed from nt 903 to 1009, including two weak stops at the 3' base (nt 940) and 5' end (nt 921) of the stem and loop structure (Fig. 5A). However, there were no strong stops in this region. Although the actual mechanism is unknown, these findings were consistent with the 70 bp insertion resulting from a complex abnormality involving an RNA-DNA duplex in cDNA library construction (42). Another possibility would be an insertion due to a complex strand-switching event involving DNA polymerase I (43). Alternatively, this 70 bp insertion may have resulted from alternative splicing, although the insertion predicts a truncated  $\alpha$ -GalNAc polypeptide of 358 residues. To investigate the possible occurrence of  $\alpha$ -GalNAc transcripts with a 70 bp insertion after pAGB-3 nt 957, PCR was used to amplify this region in 1) reverse-transcribed mRNA from various sources, 2) the cDNA inserts from clones pAGB-4 to 34, and 3) the gAGB-1 genomic clone. If the cDNA inserts or reverse-transcribed RNAs

*Figure 5.* Partial genomic sequence of human  $\alpha$ -GalNAc including an intron between coding nt 957 and 958. (A) rAGB-3: partial  $\alpha$ -GalNAc RNA sequence (nt 909 to 969) corresponding to the 3' end of  $\alpha$ -Gal A exon 6 (32). The indicated stem and loop structure between nt 918 and 937 had a  $\Delta G$  of -11.6 (29). The overlapped antisense and sense sequences shown in bold are inverted and direct repeats derived from nt 919 to 957 of pAGB-3 that are in the 70 bp insertion in pcD-HS1204 (20). The 45 bp deletion in clone pAGB-13 is indicated in italics. (B) The genomic  $\alpha$ -GalNAc sequence from coding nt 760 to 1053 (upper case) includes a 1754 nt intron between nt 957 and 958 which corresponds in position to the  $\alpha$ -Gal A exon 6 and 7 boundaries. Dashed line, the 5' splice donor sequence; solid underlines, putative branch point sequences; dotted underlines, putative polypyrimidine tracts at the 3' acceptor sites for the normal gene and mutant pAGB-13; and asterisks, differences from the consensus sequence (44).

**A**



**B**

760 CTG CTC ATT GGG AAC TTT GGT CTC AGC TTA GAG CAA TCC CGG GCC CAG ATG GCC CTG TGG ACG GTG CTG GCA GCC CCC CTC TTG ATG TCC 848  
 Leu Leu Ile Gly Asn Phe Gly Leu Ser Leu Glu Gln Ser Arg Ala Gln Met Ala Leu Trp Thr Val Leu Ala Ala Pro Leu Leu Met Ser

849 ACA GAC CTG CGT ACC ATC TCC GCC CAG AAC ATG GAC ATT CTG CAG AAT CCA CTC ATG ATC AAA ATC AAC CAG GAT CCC TTA GGC ATC CAG 938  
 Thr Asp Leu Arg Thr Ile Ser Ala Gln Asn Met Asp Ile Leu Gln Asn Pro Leu Met Ile Lys Ile Asn Gln Asp Pro Leu Gly Ile Gln

939 GGA CGC AGG ATT CAC AAG gtactagggtgtggaggggaaggaaggggagggctgaggaactgggttctcctgagagaaaggctgccagctccctggggcaacacctggcgagg 957  
 Gly Arg Arg Ile His Lys  
 Exon 6

tacaggagtgcgccagtcaccaaccagggtacccttctggttcttataggttgaggactctgatgggagctgctccaactgtcctcctcttctgctgggtgagagcagggtgagcagg  
 acagctcaagggagtcgggatgagaggtgtcagccacataagtgacacatagcaagggtaggacagagcttctatacaccctgatggcctgcagagagcttggacttccctccaga  
 gcaggagagctgggttgggttttggagacaggtctcactctgtcaccaggtggagtgagtgagcacaatttcgactcactgcaatctctacctgccaggttcaagcaattctc  
 gtgcctcagcctcctgagtagctggcactacaggcgcctgcaccacaccagctaatTTTTgtatTTTTtagtagagacaccatgttggccaggttctctcgaactcctggcctcaggt  
 gatccaccctgatcagcctcccaagtgctgggattacagggcatgagcaccgcaactcggccaggagaagctggtatagccaaggaatactacgactactggtggctgctatTTATtgag  
 tacctaccatgtctgggagtttagataatTTTTctcagcaaggtagtatatctgcccatttacaatgagaaaaatgaaacttcgagagctgagtaactttatcccaaggctacac  
 agttggtaaaaacaagactggacttcaggtgtcactcaaaagcctTTTTTTTTTTTTTTTTTTTTTgagatggagctcagctgtagccaggtggagtgagcaccatctcagct  
 cactgcaacctctggctcccaggttcaagcgatttctcctcagcctcccaggtagctgggatacaggtgtgcccaccaccccggctaattTTTTTTTgtatTTTTTcagtaga  
 gacagggTTTTcaccatgttggccaggtactctcaaaactcctgacgtcagctgatccactgcctcggcctcacaagtaaatgggattacagcatgagccactgtgcctgtctgcctt  
 gctcttaccaaatcctggattctggtaaaaaagaacactacagaactatggaaggcactatagaactggtgatgccagaggaagtaacaattccctgccagaggggctgatggtgga  
 gctggcctggaaaaccttctggaggtggagttcacatccagctccactcaccctcctggaacagagttcaactgttcccactggacagcaccctccaggccagcactggcagct  
 gtttggggccagcactcatacgtgtactgttgttgcgcttccctgttctgcgttatccctcccggtgtcctatgagcttctggggcagggtcatgcagcactgtctcagtggtc  
 tagcatagggccgggctcagagtaggtgttgatgagatctcgtgagtcaggaaggtggcagatagggtagataagctggggtgctggaggccctgctcctccctaaacctg  
 958 tgtgacatggagctgtgaactgggggaccagaactcagggagggccagggagggcaatggtaggtcctgtctgagcaagggacccagccagtagccacctctgtgccacGAA AAA 963  
 Glu Lys  
 Exon 7

964 TCT CTC ATC GAA GTG TAC ATG CGG CCT CTG TCC AAC AAG GCT AGC GCC TTA GTC TTC TTC AGC TGC AGG ACC GAT ATG CCT TAT CGC TAC 1053  
 Ser Leu Ile Glu Val Tyr Met Arg Pro Leu Ser Asn Lys Ala Ser Ala Leu Val Phe Phe Ser Cys Arg Thr Asp Met Pro Tyr Arg Tyr  
 Exon 7 - pAGB-13

contained the 70 bp insert, a 290 bp PCR product would be observed, whereas the absence of the insert would result in a 220 bp PCR product. Only the 220 bp product was observed in PCR-amplified reverse-transcribed total RNA from human lymphoblasts, fibroblasts, and placenta, or in Poly(A)<sup>+</sup> mRNA from brain (not shown). Thus, these analyses did not detect longer or shorter transcripts. All of the pAGB-4 through 34 cDNA inserts had only the 220 bp PCR product with the exception of pAGB-13, which had an inframe 45 bp deletion after pAGB-3 nt 957 (*i.e.* deleted nt 958 to 993). A short direct repeat (ACAAG) was present at both breakpoint junctions. Notably, the deletion occurred at the identical 5' site of the 70 bp insertion in pcD-HS1204 (14) (Fig. 5A).

Subsequent sequencing of the region including pAGB-3 codons 254 to 351 in the genomic clone, gAGB-1, revealed a 2048 bp sequence containing a 1754 bp intron between pAGB-3 nt 957 and 958. The intronic sequence had no homology with  $\alpha$ -Gal A intron 6, contained two *Alu*-repetitive sequences in reverse orientation and did not have the 70 bp insertion in either orientation (Fig. 5B). It was remarkable that both the pAGB-13 deletion and the pcD-HS1204 insertion occurred at the 5' donor splice site, nt 957, of this intron. Perhaps the location of the consensus lariat branch point sequences in the intron far upstream (94 and 199 bp) from the 3' splice site may impair splicing (44). This concept is supported by the pAGB-13 deletion in which the more closely positioned cryptic lariat branch point and 3' splice site were used. Thus, this intron or surrounding region may have a unique sequence and/or secondary structure that impairs the fidelity of hnRNA processing. Since the intron/exon junction after coding nt 957 also is the site of divergence between the  $\alpha$ -Gal A and B sequences, this region also may be mechanistically important in the evolution of human  $\alpha$ -GalNAc.

In conclusion, the availability of an authentic full-length cDNA encoding human  $\alpha$ -GalNAc should permit the characterization of the structure/function and evolutionary

relationships of  $\alpha$ -GalNAc and  $\alpha$ -Gal A as well as the identification of the molecular lesions that cause Schindler disease.

## REFERENCES

1. Kint, J. A. (1971) *Arch. Int. Physiol. Biochem.* **79**, 633-644
2. Beutler, E. and Kuhl, W. (1972) *Amer. J. Hum. Genet.* **24**, 237-249
3. Romeo, G., Childs, B. and Migeon B. R. (1972) *FEBS Lett.* **27**, 161-166
4. Wood, S. and Nadler, H. L. (1972) *Am. J. Hum. Genet.* **24**, 250-255
5. Ho, M. W., Beutler, S., Tennant, L. and O'Brien, J. (1972) *Am. J. Hum. Genet.* **24**, 256-266
6. Desnick, R. J., Allen, K. Y., Desnick, S. J., Raman, M. K., Bernhlohr, R. W. and Krivit, W. (1973) *J. Lab. Clin. Med.* **81**, 157-171
7. Desnick, R. J. and Bishop, D. F. (1989) in *The Metabolic Basis of Inherited Disease*. (Scriver, C. R., Beaudet, A. L. Sly, W. S. and Valle, D., eds), pp. 1751-1796, McGraw Hill, New York.
8. Beutler, E. and Kuhl, W. (1972) *J. Biol. Chem.* **247**, 7195-7200
9. Callahan, J. W., Lasilla, E. L., Den Tandt, W. and Philippart, M. (1973) *Biochem. Med.* **7**, 424-431
10. Dean, K. J., Sung, S. and Sweeley, C. C. (1977) *Biochem. Biophys. Res. Comm.* **77**, 1411-1417
11. Schram, A. W., Hamers, M. N. and Tager, J. M. (1977) *Biochim. Biophys. Acta.* **482**, 138-144
12. Kusiak, J. W., Quirk, J. M. and Brady, R. O. (1978) *J. Biol. Chem.* **253**, 184-190
13. Dean K. J., and Sweeley, C. C. (1979) *J. Biol. Chem.* **254**, 10001-10005
14. Bishop, D. F., Dean, K. J., Sweeley, C. C. and Desnick, R. J. (1980) in *Enzyme Therapy in Genetic Disease:2*, (Desnick, R.J., ed.), pp. 17-32, Alan R. Liss, Inc., New York
15. deGroot, P. G., Westerveld, A., Meera-Khan, P. and Tager, J. M. (1978) *Hum. Genet.* **44**, 305-312
16. Sweeley, C. C., LeDonne, N. C. and Robbins, P. W. (1983) *Arch. Biochim. Biophys.* **223**, 158-165
17. van Diggelen, O. P., Schindler, D., Willemsen, R., Boer, M., Kleijer, W. J., Huijmans, J. G. M., Blom, W. and Galjaard, H. (1988) *J. Inher. Met. Dis.* **11**, 349-357

18. Schindler, D., Bishop, D. F., Wolfe, D. E., Wang, A. M., Egge, H., Lemieux, R. U. and Desnick, R. J. (1989) *N. Engl. J. Med.* **320**, 1735-1740
19. Schindler, D. Kanzaki, T. and Desnick, R. J. (in press) *Clin. Chim. Acta.*
20. Tsuji, S., Yamanchi, T., Hiraiwa, M., Isobe, T., Okuyama, T., Sakimura, K., Takahashi, Y., Mishizawa, M., Uda, Y. and Miyatake, T. (1989) *Biochem. Biophys. Res. Comm.* **163**, 1498-1504
21. Bishop, D. F. and Desnick, R. J. (1981) *J. Biol. Chem.* **256**, 1307-1316
22. Calhoun, D., Bishop, D. F., Berstein, H. S., Quinn, M., Hantzopoulos, P. and Desnick, R. J. (1985) *Proc. Natl. Acad. Sci. USA* **82**, 7364-7368
23. Hunkapillar, M. W., Lujan, E., Ostrander, F. and Hood, L. E. (1973) *Methods Enzymol.* **91**, 227-236
24. Tsai, S. F., Bishop, D. F. and Desnick, R. J. (1988) *Proc. Natl. Acad. Sci. USA* **55**, 7049-7053
25. Hunkapillar, M. W. and Hood, L. E. (1983) *Science* **219**, 650-669
26. Maniatis, T., Fritsch, E. F. and Sambrook, J. (eds.) (1982) in *Molecular Cloning: A Laboratory Manual*, pp. 309-328, Cold Spring Harbor Laboratories, NY, NY
27. Sanger, F., Coulson, A. R., Barrell, B. G., Smith, A. J. H. and Roe, B. A. (1980) *J. Mol. Biol.* **143**, 161-178
28. Wolf, H., Modrow, S., Motz, M., Jameson, B. A., Hermann, G., and Fortsch, B. (1988) *CABIOS*. **4**, 187-191
29. Zuker, M. (1989) *Methods Enzymol.* **180**, 262-88
30. Wong, G. G., Witek, J. S., Temple, P. A., Wilkens, K. M., Leary, A. C., Luxenberg, D. P., Jones, S. S., Brown, E. L., Kay, R. M., Orr, E. C., Shoemaker, C. S., Golde, D. W., Kaufman, R. J., Hewick, R. M., Wang, E. A. and Clark, S. C. (1985) *Science* **228**, 810-813
31. Chen, C. and Okayama, H. (1987) *Mol. Cell. Biol.* **7**, 2745-2752
32. Bishop, D. F., Kornreich, R. and Desnick, R. J. (1988) *Proc. Natl. Acad. Sci. USA* **85**, 3903-3907
33. Ho, S. N., Hunt, H. D., Horton, R. M., Pullen, J. K. and Pease, L. R. (1989) *Gene* **77**, 51-59
34. Berget, S. M. (1984) *Nature (London)* **309**, 179-182

35. von Heijne, G. (1986) *Nucleic Acids Res.* **14**, 4683-4960
36. Kornreich, R., Desnick, R. J. and Bishop, D. F. (1989) *Nucleic Acids Res.* **17**, 3301-3302
37. Liljeström, P. L. (1985) *Nucleic Acids Res.* **13**, 7257-7268
38. Liljeström, P. L. and Lijeström, P. *Nucleic Acids Res.* **15**, 2213-2220
39. Proira, R. L. (1988) *Proc. Natl. Acad. Sci. USA* **85**, 1883-1887
40. Horowitz, M., Wilder, S., Horowitz, A., Reiner, O., Gelbart, T. and Beutler, E. (1989) *Genomics* **4**, 87-96
41. Schuchman, E. H., Jackson, C. E. and Desnick, R. J. (1990) *Genomics* **6**, 149-158
42. Roberts, J. D., Preston, B. D., Johnston, L. A., Soni, A., Loeb, L. A. and Kunkel, T. A. (1989) *Mol. Cell. Biol.* **9**, 468-476
43. Papanicolaou, C. and Ripley, L. S. (1989) *J. Mol. Biol.* **207**, 335-353
44. Reed, R. and Maniatis, T. (1988) *Genes Dev.* **2**, 1268-1276

**CHAPTER TWO:**

**SCHINDLER DISEASE: THE MOLECULAR LESIONS IN THE  
 $\alpha$ -N-ACETYLGALACTOSAMINIDASE GENE  
THAT CAUSE TYPE I AND II DISEASE**

**ABSTRACT**

Schindler disease is a recently recognized disorder resulting from the deficient activity of the lysosomal hydrolase,  $\alpha$ -N-acetylgalactosaminidase ( $\alpha$ -GalNAc). Two clinical phenotypes of  $\alpha$ -GalNAc deficiency have been described. The Type I disease is similar to the infantile form of neuroaxonal dystrophy, whereas the Type II disorder is characterized by disseminated angiokeratoma corporis diffusum with no neurological involvement. The recent isolation and expression of the full-length cDNA encoding  $\alpha$ -GalNAc has facilitated the identification of the molecular lesions in the first affected homozygotes with each of the remarkably distinct forms of this autosomal recessive disease. Southern and northern hybridization analyses of DNA and total RNA from the Type I and II affected homozygotes revealed normal  $\alpha$ -GalNAc gene structures and normal transcript sizes and amounts. Therefore, the reverse-transcribed  $\alpha$ -GalNAc transcripts from affected homozygotes with Type I and Type II disease were amplified individually by the polymerase chain reaction (PCR) for DNA sequence analysis. In a sib from Family D with Type I disease, a single G to A transition at nucleotide 973 was detected in multiple subclones of the PCR-amplified  $\alpha$ -GalNAc transcript. This point mutation results in a glutamic acid to lysine substitution in residue 325 of the  $\alpha$ -GalNAc polypeptide. In the Type II homozygote from Family K, a single C to T transition was identified at nucleotide 985 which results in an arginine to tryptophan substitution at position 329. These base substitutions were confirmed by dot-blot analyses of PCR-amplified genomic DNA from members of each family using allele-specific oligonucleotides. Furthermore, transient expression of  $\alpha$ -GalNAc constructs containing the 973 G to A transition and the 985 C to T transition each resulted in the expression of an immunoreactive polypeptide which had no detectable  $\alpha$ -GalNAc activity.

## INTRODUCTION

Schindler disease is a recently recognized disorder resulting from the deficient activity of the lysosomal glycohydrolase,  $\alpha$ -*N*-acetylgalactosaminidase (E.C. 3.2.1.29;  $\alpha$ -GalNAc) (1-5). The enzymatic defect, inherited as an autosomal recessive trait, leads to the tissue accumulation and increased urinary excretion of glycopeptides and oligosaccharides containing  $\alpha$ -*N*-acetylgalactosaminyl moieties (6-8). Two subtypes, designated Type I and II, with variant phenotypes have been described (2, 5). The Type I disorder is an infantile form of neuroaxonal dystrophy first described in two brothers, the offspring of a consanguineous couple of German descent (Family D). The affected infants had normal births and early developments for the first 9 to 12 months of age. Then, they experienced a period of developmental delay which was followed by rapid regression in the second year of life. By three to four years of age, the progressive neurologic impairment led to cortical blindness, myoclonic seizures, spasticity, decorticate posturing and profound psychomotor retardation. A cortical biopsy revealed the presence of "spheroids" in terminal axons, the unique neuropathology which classified this disease as a neuroaxonal dystrophy (3, 9). The clinical and morphologic manifestations of the two affected brothers with this disease are essentially identical to those in patients with Seitelberger disease (9); however, the finding of normal  $\alpha$ -GalNAc activity in eight unrelated patients with Seitelberger disease demonstrated that Schindler disease is biochemically distinct and that the infantile neuroaxonal dystrophies are genetically heterogeneous (3).

The Type I affected homozygotes had  $\alpha$ -GalNAc activity levels that were less than 1% of normal in various sources as assayed with the newly synthesized substrate, 4-methylumbelliferyl- $\alpha$ -*N*-acetylgalactosaminide (4MU- $\alpha$ -GalNAc) (3, 10). Their parents had intermediate levels of activity consistent with being obligate heterozygotes for this autosomal recessive disorder. Using monospecific rabbit anti-human  $\alpha$ -GalNAc antibodies, immunoblots

of fibroblast extracts from the affected brothers revealed no detectable immunoreactive enzyme protein, while in extracts from their parents and a brother, there were two immunoreactive peptides of 48 and 117 kDa, the monomeric and dimeric forms of  $\alpha$ -GalNAc observed in cultured fibroblasts from normal individuals and the purified human enzyme (3). These findings suggested that the  $\alpha$ -GalNAc mutation in this family markedly altered the enzyme's activity and stability.

Type II Schindler disease was described in a 46-year old Japanese woman, the daughter of a first-cousin marriage (Family K). This affected homozygote had no neurological involvement, but had disseminated angiokeratoma similar to that observed in males with Fabry disease (4). Urinary oligosaccharide analysis revealed that this individual excreted the same glycopeptides as those characterized in the urine from the two brothers with Type I Schindler disease (6, 7). These findings suggested that this patient may have  $\alpha$ -GalNAc deficiency. In fact, the affected homozygote had less than 1% of normal  $\alpha$ -GalNAc activity in various sources as assayed with 4-MU- $\alpha$ -GalNAc. Her children had intermediate levels of activity consistent with being obligate heterozygotes for this autosomal recessive disorder. Immunoblotting of fibroblast extracts from this individual also revealed no immunoreactive enzyme protein, while the two peptide species were observed in extracts from her son. These findings suggested that the  $\alpha$ -GalNAc mutation in this family also altered the enzymes's activity and stability.

The recent isolation, sequencing and expression of the full-length cDNA encoding human  $\alpha$ -GalNAc (11) has facilitated the investigation of the molecular lesions in the  $\alpha$ -GalNAc gene that cause Schindler disease in these families. The  $\alpha$ -GalNAc cDNA encodes a polypeptide of 411 amino acids including a signal peptide of 17 residues. The precursor polypeptide is co-translationally glycosylated and, following carbohydrate modifications in the Golgi, the mature polypeptide of 394 amino acids dimerizes to form the active enzyme of ~117

kDa. Studies of the biosynthesis of human  $\alpha$ -GalNAc in fibroblasts indicated that both the 65 kDa precursor and 48 kDa mature lysosomal forms had only high mannose-type oligosaccharide structures (12). In this communication, the specific base substitutions in the  $\alpha$ -GalNAc coding region of each family are identified by sequencing the PCR-amplified products of the reverse-transcribed mutant  $\alpha$ -GalNAc mRNAs. Confirmation of this mutation was demonstrated by dot-blot analyses of PCR-amplified genomic DNA from family members, and by transient expression studies of the mutant cDNAs. Of interest was the fact that both of the expressed mutant enzymes had no activity, but were immunologically detectable, whereas no immunoreactive enzyme protein was present in cultured cells from the patients. These findings suggest that the enzyme polypeptide is synthesized *in vivo*, even though sensitive immunologic analyses may reveal the absence of crossreactive immunological material in cultured cells or tissue extracts.

## METHODS

*Cell lines.* Primary cultures of fibroblasts and lymphoblasts were established from skin biopsies and peripheral blood samples obtained from the Family D and K members with Schindler disease and from normal controls with informed consent. The COS-1 cells lines were obtained from the American Type Tissue Collection (Rockville, MD). The fibroblasts, lymphoblasts and COS-1 cell lines were grown in RPMI with 10% fetal bovine serum and 1% penicillin-streptomycin and 2 mM glutamine (Gibco, Grand Island, NY) by standard procedures (13).

*Assays of  $\alpha$ -GalNAc activity and protein.* The  $\alpha$ -GalNAc activity in cultured cells and COS-1 cells was determined with the synthetic fluorogenic substrate 4-methylumbelliferyl- $\alpha$ -N-acetylgalactosaminide (10) as previously described (3). One unit (U) of enzymatic activity is equal to that amount of enzyme required to hydrolyze 1 nmol of 4MU- $\alpha$ -GalNAc per hour. Protein concentration was determined by the fluorescamine assay (14). Immunoblot analyses of  $\alpha$ -GalNAc in transfected COS-1 cells was performed using monospecific rabbit anti-human antibodies as described (3, 12).

*Southern and northern hybridization analyses.* For analysis of the  $\alpha$ -GalNAc gene integrity, genomic DNA was isolated from at least  $10^6$  cultured fibroblasts or lymphoblasts (15). The DNA was digested with restriction endonucleases (e.g., *BamH* I, *Pst* I, *EcoR* I, *Taq* I) (New England Biolabs, Beverly, MA), electrophoresed in 1% agarose, transferred to BioTrace RP charge-modified nylon 66 binding matrix (Gelman Sciences, Inc., Ann Arbor, MI) (16) and analyzed with nick-translated  $\alpha$ -GalNAc cDNA (17). For analysis of the relative sizes and amount of the  $\alpha$ -GalNAc transcripts, total RNA was isolated from at least  $10^8$  cultured fibroblasts or lymphoblasts (18). RNA samples were electrophoresed in 1% agarose/formaldehyde denaturing gels (19), transferred to BioTrace (Gelman Sciences), and

analyzed with radiolabelled riboprobe synthesized with SP6 polymerase from the  $\alpha$ -GalNAc cDNA cloned into the pGEM4Z vector (Promega, Madison, WI).

*DNA amplification and sequencing of the mutant allele.* Sense and anti-sense oligonucleotide primers designed to amplify the entire coding region of the  $\alpha$ -GalNAc gene in two overlapping fragments were synthesized on a model 380B DNA synthesizer (Applied Biosystems, Foster City, CA). As shown in Fig. 1, the 5' region of the  $\alpha$ -GalNAc transcript [cDNA nt -54 to 689 (11)] was amplified using the 32-mer sense primer, P1 (5'-AGTAGTGAATTCCTGATACACGCAGACCAGAT-3') corresponding to  $\alpha$ -GalNAc nt -34- -53 with an additional 12 nt which included an *EcoR* I site and the 32-mer antisense primer, P2 (5'-AGTAGTAAGCTTTTCAGGATGGAGAGCTCGCT-3') corresponding to  $\alpha$ -GalNAc cDNA nt 670-689 with a *Hind* III site for forced cloning. The 3' sequence of the coding (cDNA nt 595 to 1292) was amplified using the 32-mer sense primer, P3 (5'-AGTAGTGAATTCAGGGTGAAGTACTACAGTCTGCT-3') corresponding to nt 595-614 with an *EcoR* I site and a 32-mer antisense primer, P4 (5'-AGTAGTAAGCTTGCTCCATGGTCTAGGCTCAG-3') corresponding to nt 1273-1292 and contained a *Hind* III site. Total RNA (10  $\mu$ g) was reverse-transcribed to cDNA using the BRL cDNA Synthesis Kit (Bethesda Research Laboratories, Gaithersburg, MD). One-fourth of the cDNA product was PCR-amplified (20) using the GeneAmp DNA Amplification Reagent Kit (Perkin Elmer Cetus, Norwalk, CT) and 1  $\mu$ M of each primer. Each of the 30 PCR cycles consisted of denaturation at 94 °C for 1 min; annealing at 37 °C for 2 min; and extension at 60 °C for 7 min. The PCR products were phenol extracted, ethanol precipitated, resuspended in 20  $\mu$ l H<sub>2</sub>O and 2  $\mu$ l were analyzed on a 2% agarose gel. The remaining 18  $\mu$ l was digested with *EcoR* I and *Hind* III and subcloned into M13mp18 and 19 vectors (21). Clones containing the PCR products were identified by plaque hybridization using the nick-translated  $\alpha$ -GalNAc cDNA (17). Single-

stranded template was isolated (21) from six separate M13 clones, and each was sequenced in both orientations by the dideoxy chain termination method (22) using universal and sequence-specific primers.

*Oligonucleotide hybridization of genomic DNA.* To confirm the mutation at the genomic level, 1  $\mu$ g of genomic DNA from the Family D and K affected homozygotes, obligate heterozygotes, and normal individuals was PCR-amplified by the conditions described above using primers that flanked the region of the mutations. The 5' sense primer was a 20-mer (5'-TGATGTCCACAGACCTGCGT-3') corresponding to nt 842-861 of the full-length cDNA (11) and the 3' antisense primer was a 20-mer (TCAGCTGCAGGACCGATATG) corresponding to nt 1005-1024. Normal and mutation specific oligonucleotide probes (21-mers) were synthesized with the normal sequence 5'-TCTCTCATCGAAGTGTACATG-3' and the mutation-specific sequence 5'-TCTCTCATCAAAGTGTACATG-3' for detection of the Type I base substitution, and 5'-GTGTACATGCGGCCTCTGTCC-3' and 5'-GTGTACATGTGGCCTCTGTCC-3', for identification of the Type II lesion. The calculate  $T_m$  values for the Type I normal and mutation specific oligonucleotides were 58 °C and 55 °C, respectively, while the corresponding values for the Type II allele-specific oligonucleotides were 65 °C, and 63 °C, respectively. Dot blots were denatured, hybridized, and washed as described (23).

*Construction of the mutant cDNAs and transient expression of the Type I and II mutations in COS-1 cells.* A *Bam*H I-*Kpn* I fragment containing the Type I G to A or Type II C to T transitions from the respective M13 clones were used as cassettes to replace the corresponding sequence in the pAGB-3 cDNA. The Type I and II constructs were individually subcloned into the eukaryotic expression vector p91023(B) (27) and the sequence of each construct was confirmed by dideoxy sequencing. COS-1 cells were transfected (28) with the

Type I or II full-length construct, harvested at 72 h post-transfection, assayed and immunoblotted as previously described (2).

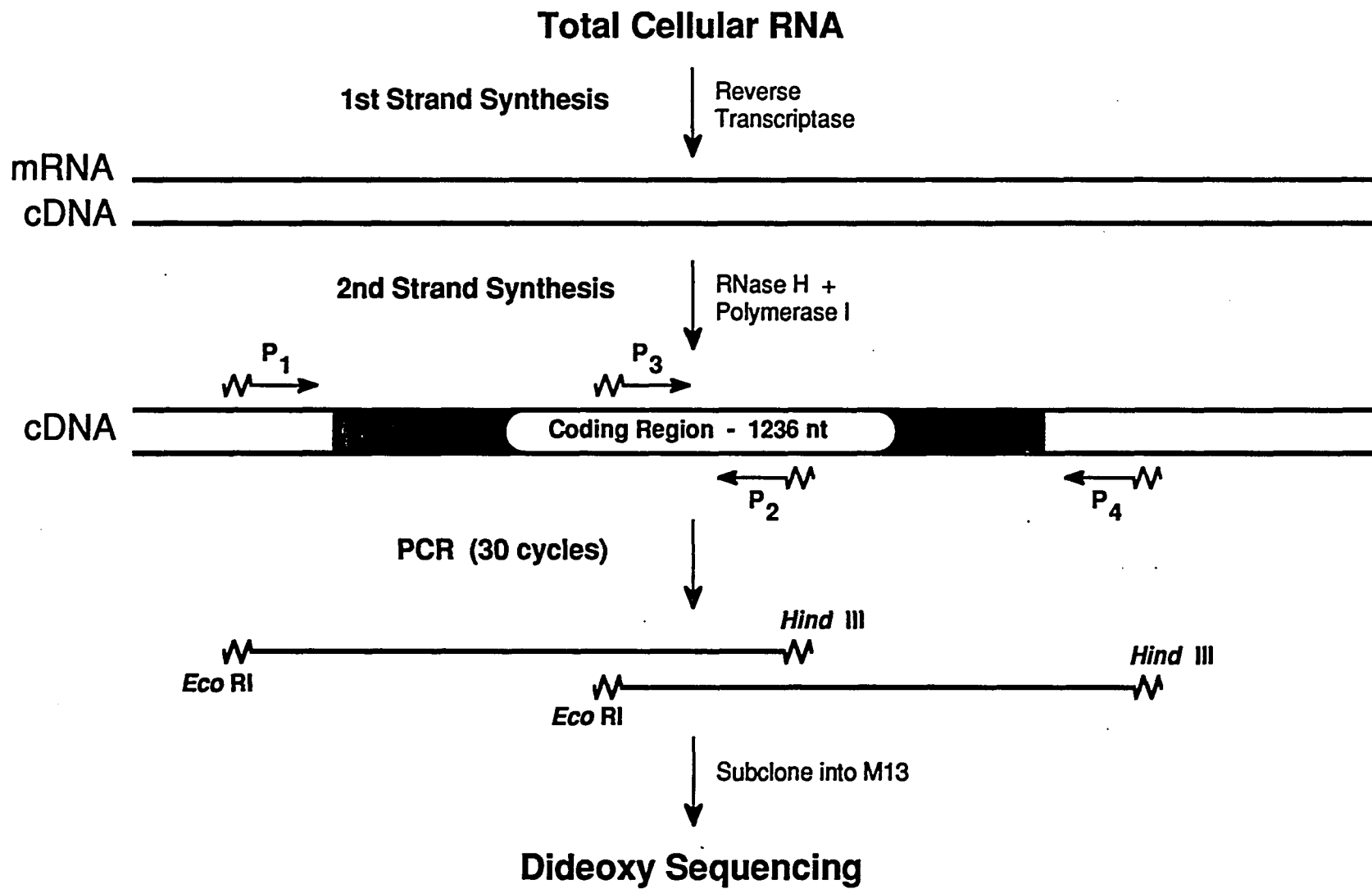
*Computer-assisted analyses.* Local secondary structure for both the normal and mutant enzymes were predicted by the algorithms of Chou and Fasman (24) and Garnier *et al.* (25) using the University of Wisconsin Computer Group software (26). A region of 20 amino acids surrounding the site of the Type I and II amino acid substitutions were used for the prediction of the local secondary structure.

## RESULTS

*Molecular characterization of Families D and K.* In order to characterize the nature of the mutations in Families D and K, Southern analysis was performed with several restriction endonucleases of genomic DNA from the affected homozygotes and obligate heterozygotes. No gross gene rearrangements were observed (data not shown). Interestingly, Family D was polymorphic for two restriction endonucleases, *BamH* I and *Taq* I. In a normal Caucasian population of 38 and 58 alleles, the polymorphisms were found at a frequency of about 44% and 2.6%, respectively. Northern analysis of total RNA from the Family D and K homozygotes and obligate heterozygotes revealed the presence of the two  $\alpha$ -GalNAc transcripts of ~2.2 kb and ~3.0 kb. Each of the transcripts were present in normal amounts in all individuals studied (data not shown), thereby excluding the possibility of splice, promoter, or mRNA stability mutations.

These findings, and the fact that no  $\alpha$ -GalNAc immunoreactive protein was detected in fibroblasts from Type I or II homozygotes (2, 5), suggested that the mutations causing Type I and II Schindler disease in these families were due to point mutations or small insertions or deletions within the coding region of the  $\alpha$ -GalNAc gene. To identify the mutations in the  $\alpha$ -GalNAc coding region of these Type I and II homozygotes, their  $\alpha$ -GalNAc transcripts were reverse-transcribed, PCR-amplified and sequenced in both orientations (Fig. 1). In the sequence analyses of the Type I and II  $\alpha$ -GalNAc cDNAs, a single unique substitution was found in all Type I or II subclones. The transcript from the Type I homozygote had a G to A transition at nucleotide 973 which resulted in a glutamic acid to lysine substitution at residue 325 (Fig. 2). In all the Type II cDNAs, a C to T transition occurred at nt 985 which would substitute a tryptophan for an arginine at residue 329 (Fig. 2). These substitutions were confirmed by dot-blot hybridization analyses with PCR-amplified genomic DNA from the respective Family D and K affected homozygotes, obligate heterozygotes with allele-specific oligonucleotides.

***Figure 1.*** Strategy for the identification of the specific molecular lesions in Schindler disease Types I and II. Total RNA was isolated from lymphoblasts of affected homozygotes, reverse-transcribed to cDNA, PCR-amplified in two overlapping reactions, cloned into M13 sequencing vectors, and sequenced in both orientations.



***Figure 2.*** The mutation in the Type I homozygote was identified by DNA sequencing to be a G to A transition at nucleotide 973 which caused a glutamic acid to lysine substitution at residue 325. The mutation in the Type II homozygote was identified by DNA sequencing to be a C to T transition at nucleotide 985 which caused an arginine to tryptophan substitution at residue 329.

<b>322</b>	<b>323</b>	<b>324</b>	<b>325</b>	<b>326</b>	<b>327</b>	<b>328</b>	<b>329</b>	<b>330</b>	<b>331</b>	<b>332</b>
Ser -	Leu -	Ile -	<b>Glu</b> -	Val -	Tyr -	Met -	<b>Arg</b> -	Pro -	Leu -	Ser
TCT	CTC	ATC	<b>GAA</b>	GTG	TAC	ATG	<b>CGG</b>	CCT	CTG	TCC
			↓				↓			
			<b>AAA</b>				<b>TGG</b>			
			Lys				Trp			
			Family D				Family K			

As shown in Figs 3A and 3B, the respective normal probes hybridized with to normal and heterozygote amplified DNA products, while the respective mutation-specific probes hybridized only with the amplified genomic sequences from the respective Type I or II homozygotes indicating that each homozygote was homoallelic for the respective Type I or II mutations. In support of these findings, PCR-amplified genomic DNA from both obligate heterozygous parents of the Type I homozygotes and both obligate heterozygous children of the Type II homozygote hybridized with the respective mutant probes.

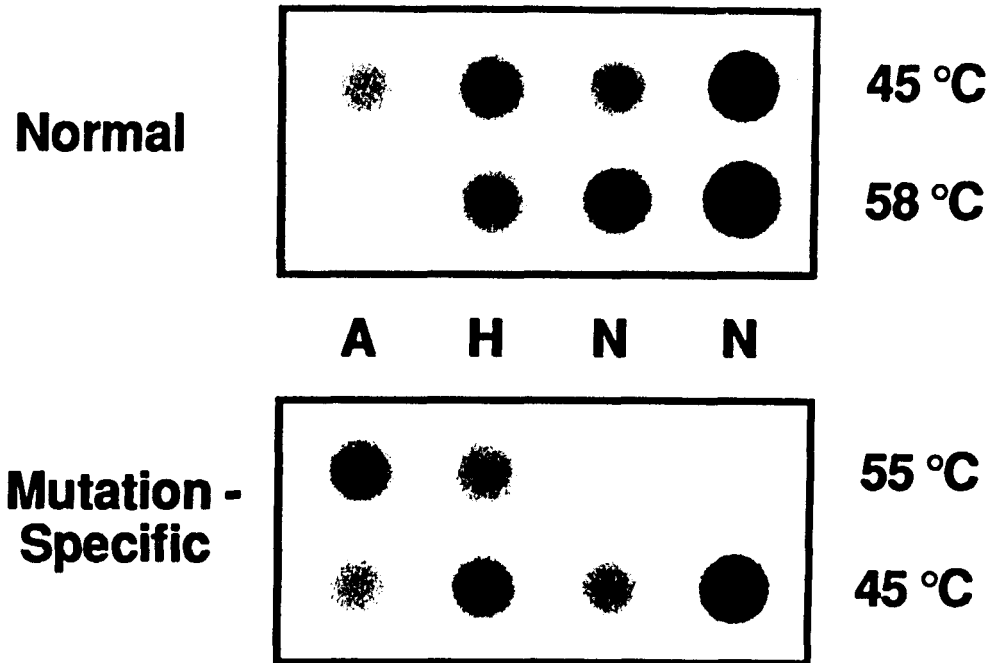
*Expression of the  $\alpha$ -GalNAc mutations in COS-1.* To further characterize the Type I and II mutations, each was introduced as a *Bam*H I-*Kpn* I cassette (nt 921-1193) into the full-length pAGB-3 cDNA. The mutated cDNAs, designated pAGB-973A and pAGB-985T, were then subcloned into the eukaryotic expression vector p91023(B) and sequenced. Each construct, p91-AGB-973A or p91-AGB-985T was transfected into COS-1 cells, harvested at 72 h post-transfection, and the extracts assayed for  $\alpha$ -GalNAc activity and immunoreactive protein. Immunoblot analysis revealed the presence of immunologically detectable proteins with subunit molecular weights of ~48 kDa expressed from the p91-AGB-973A and p91-AGB-985T constructs (Fig. 4A and B). However, neither mutant protein expressed 4MU- $\alpha$ -GalNAc activity above endogenous COS-1 levels, whereas the normal construct p91-AGB-3 expressed immunologically detectable protein which had  $\alpha$ -GalNAc activity that was 20-fold greater than in mock-transfected COS-1 cells (Table 1).

*Secondary structure analyses.* Computer-assisted regional secondary structural analysis of residues 310 to 330 was carried out using the algorithm of Garnier *et al.* (22). Compared to the predicted secondary structure of the normal sequence, the glutamic acid to lysine substitution in the Type I homozygote extended a  $\beta$ -pleated sheet region 5' to the mutation, a created a region of  $\alpha$ -helicity around and including the mutation (Fig. 5). The arginine to tryptophan

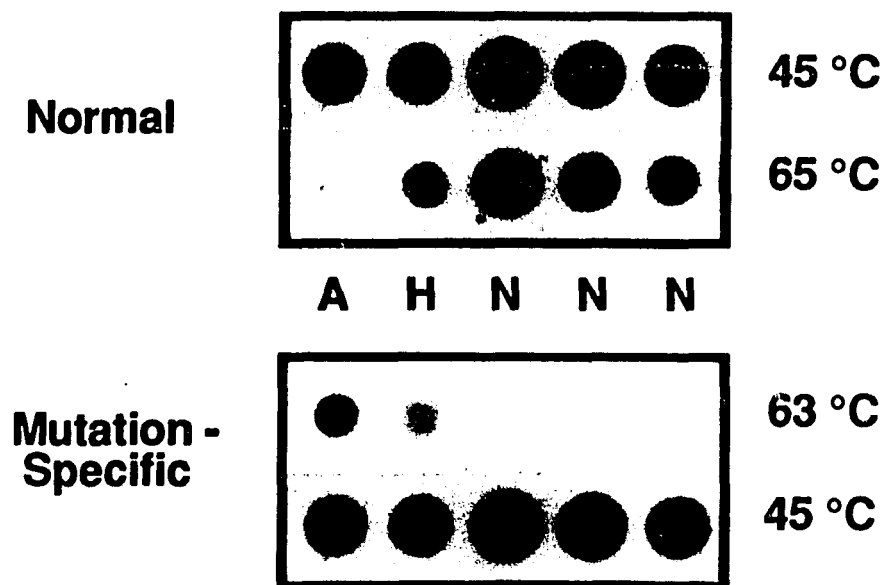
substitution in the Type II homozygote predicted a shortening of  $\beta$ -pleated sheet 3' to the mutation (Fig. 5).

*Figure 3.* Allele specific oligonucleotide probes were constructed and hybridized to PCR-amplified genomic DNA by dot blot analysis and washed at the differential temperatures as indicated. A, affected homozygote; H, obligate heterozygote; and N, normal control. (A) The dot blots for Family D with Type I Schindler disease. (B) The dot blots for Family K with Type II Schindler disease.

A

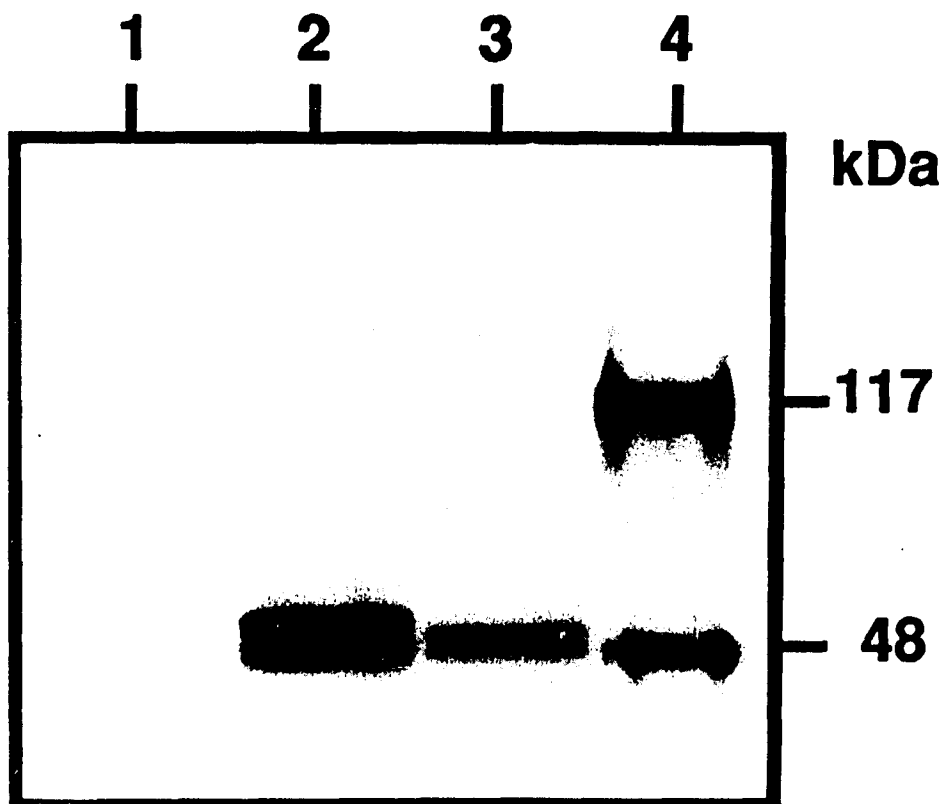


B

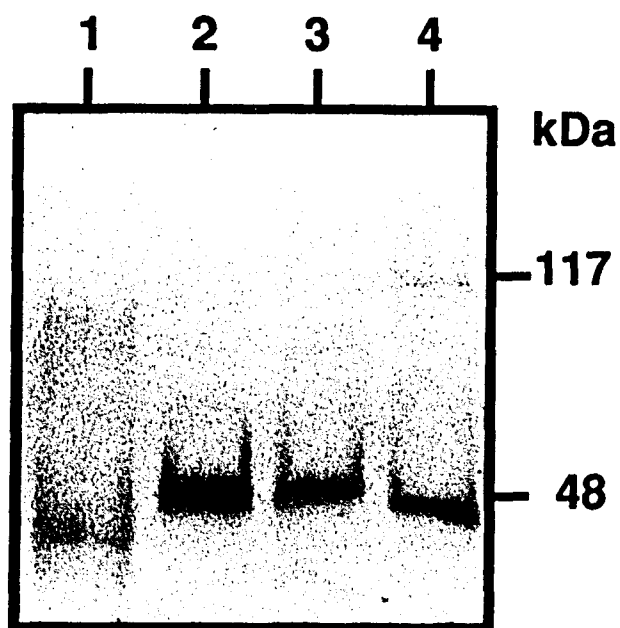


*Figure 4.* Immunoblot of transiently expressed  $\alpha$ -GalNAc in COS-1 cells. (A) Family D: lanes 1, mock transfection; 2, p91-AGB-3; 3, p91-AGB-973A; 4, human fibroblast control. (B) Family K: lanes 1, mock transfection; 2, p91-AGB-3; 3, p91-AGB-985T; 4, human fibroblast control.

A



B

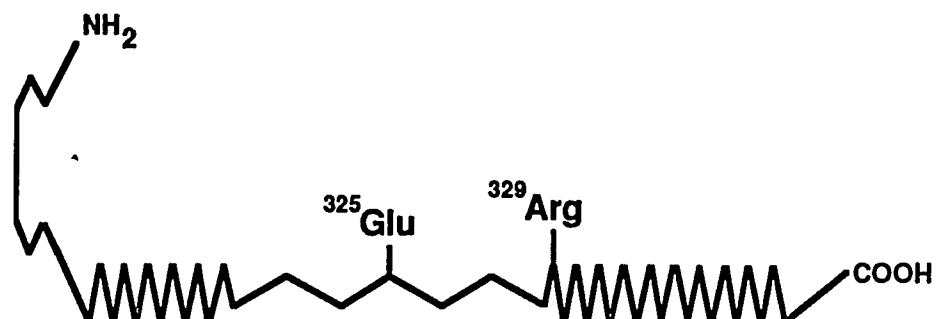
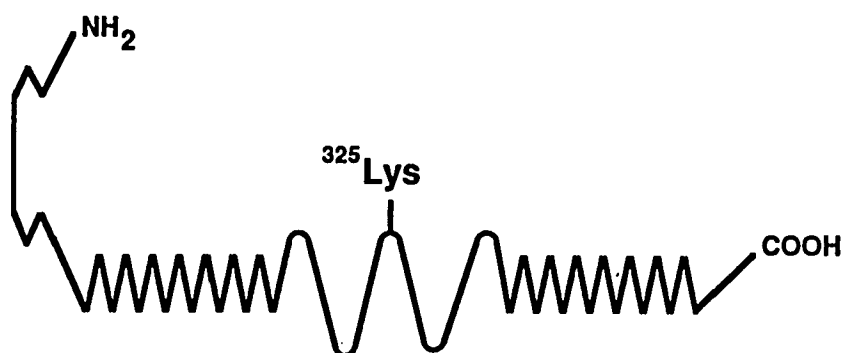
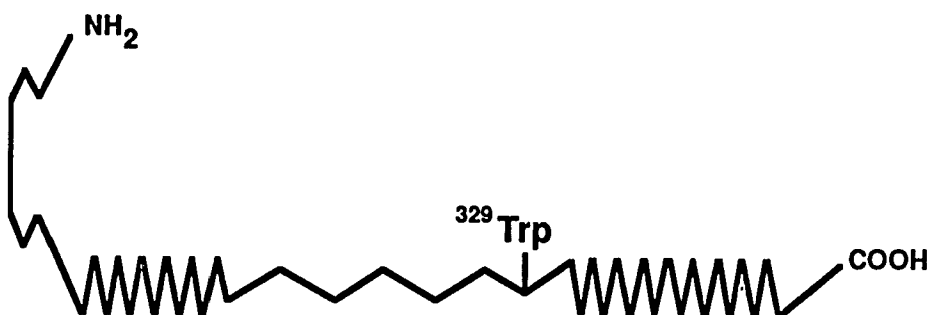


**Table 1. Transient expression of normal and mutant  $\alpha$ -GalNAc constructs in COS-1 cells.**

Construct	$\alpha$ -GalNAc Activity*
	U/mg
p91-AGB-3	
Mean	661
Range (n=6)	105-2415
p91-AGB-973A	
Mean	44
Range (n=6)	27-65
p91-AGB-985T	
Mean	36
Range (n=5)	31-64
Mock Transfection (No DNA)	
Mean	35
Range (n=6)	23-50

\*As defined using 0.8 mM 4MU- $\alpha$ -GalNAc as substrate.

**Figure 5.** Computer assisted secondary structure was determined for the residues flanking and including the mutations for the normal and mutant sequences. Top: normal sequence; middle: Type I substitution of glutamic acid to lysine at residue 325; bottom: Type II substitution of arginine to tryptophan at residue 329.

**NORMAL****FAMILY D****FAMILY K****Amino Acids 310-330**

## DISCUSSION

The recent identification of two remarkably distinct phenotypes with severely deficient  $\alpha$ -GalNAc activity stimulated studies to investigate the molecular nature of these enzymatic defects. Using the full-length  $\alpha$ -GalNAc cDNA as a probe, Southern and northern analyses revealed no gross chromosomal rearrangements and transcripts of normal size and abundance in the DNA and RNA from the affected homozygotes with Type I and II Schindler disease. Subsequent reverse-transcription, amplification, and sequencing of the  $\alpha$ -GalNAc transcripts from these affected individuals revealed specific base substitutions. In the two brothers with Type I disease, a G to A transition at nt 973 resulted in a glutamic acid to lysine substitution at residue 325 of the  $\alpha$ -GalNAc polypeptide. In the affected woman with Type II disease, a C to T transition at nt 985 was identified which would substitute a tryptophan for an arginine at residue 329. These nucleotide changes were confirmed in genomic DNA from the consanguineous parents of the Type I homozygotes and the children of the Type II homozygote using allele-specific oligonucleotide probes. The mutations causing Type I and II Schindler disease were homoallelic, consistent with the fact that the Type I and II homozygotes were the products of consanguineous marriages. Computer-assisted secondary analysis of this region revealed that these substitutions changed the local secondary structure; the Type I mutation causing an extended  $\beta$ -pleated sheet region 5' to the substitution and a region of  $\alpha$ -helicity around and including the mutation, whereas the Type II mutation appeared to alter the secondary structure to a lesser degree by decreasing the  $\beta$ -pleated sheet regions 3' to the mutation.

The fact that fibroblast extracts from the affected homozygotes were CRIM-negative suggested that the two mutations caused the enzyme polypeptide to be unstable and rapidly degraded in cultured cells and, presumably, *in vivo*. Similar findings of CRIM-negative mutant proteins in fibroblasts which were immunologically detectable when transiently expressed in

COS-1 cells have been observed for mutations in the vitamin D receptor gene (29) and the cystathionine  $\beta$ -synthase gene (L. E. Rosenberg, personal communication). The detection of immunoreactive enzyme in the COS-1 cells most likely reflects the "pulsed" synthesis of the altered human polypeptide which then is degraded due to instability and/or the inability of the polypeptide subunits to associate into the more stable multimeric forms. This may be the case for the  $\alpha$ -GalNAc polypeptide since the homodimeric form is extremely stable (11). Thus, the classification of a mutation as CRIM-negative based on studies of the patient's cells or tissues does not imply that the enzyme is not synthesized. Subsequent stable expression of normal and mutant constructs may provide sufficient quantities of these enzyme proteins for comparison of their physical and kinetic properties. For example such studies may determine if the mutations alter subunit association, the interaction with different substrates, and/or other factors in the catalytic complex.

## REFERENCES

1. van Diggelen, O. P., D. Schindler, R. Willemson, M. Boer, W. J. Kleijer, J. G. M. Huijmans, W. Blom, and H. Galjaard. 1988.  $\alpha$ -N-acetylgalactosaminidase deficiency, a new lysosomal storage disorder. *J. Inher. Met. Dis.* 11:349-357.
2. Schindler, D., D. F. Bishop, D. E. Wolfe, A. M. Wang, H. Egge, R. U. Lemieux, and R. J. Desnick. 1989. Neuroaxonal dystrophy due to lysosomal  $\alpha$ -N-acetylgalactosaminidase deficiency. *N. Engl. J. Med.* 320:1735-1740.
3. Desnick, R. J., and D. F. Bishop. 1989. Schindler disease:  $\alpha$ -N-acetylgalactosaminidase deficiency. *In* The Metabolic Basis of Inherited Disease. C.R. Scriver, A.L. Beaudet, W.S. Sly and D. Valle, editors. McGraw-Hill, New York. 1751-1796.
4. Kanzaki, T., M. Yokota, M. Mizuno, Y. Matsumoto, and Y. Kirabayashi. 1989. Novel lysosomal lycoamino acid storage disease with angiokeratoma corporis diffusum. *Lancet* April 22, 1989:875-877.
5. Wang, A. M., T. Kanzaki, D. Schindler, and R. J. Desnick. 1989. Schindler disease: Molecular defects in the infantile- and recently identified adult-onset forms. *Am. J. Hum. Genet.* 43:4A228.
6. Linden, H. U., R. A. Klein, H. Egge, J. Peter-Katalinic, J. Dabrowski, and D. Schindler. 1989. Isolation and characterization of sialic acid-containing glycopeptides of the O-glycosidic type from the urine of two patients with hereditary deficiency in  $\alpha$ -GalNAc activity. *Biol Chem. Hoppe-Seyle*r. 370:661-672.
7. Hirabayashi, Y., Y. Matsumoto, M. Matsumoto, T. Toida, N. Iida, T. Masubara, T. Kanzaki, M. Yokota, and I. Shizuka. 1990. Isolation and characterization of major urinary amino acid O-glycosides and a dipeptide O-glycoside from a new lysosomal storage disorder (Kanzaki disease). *J. Biol. Chem.* 265:1693-1701.
8. Schindler, D., T. Kanzaki, and R. J. Desnick. In press. *Clin Chim Acta*.
9. Seitelberger, F. 1986. Neuroaxonal dystrophy: its relation to aging and neurological disease. In Handbook of Clinical Neurology, Vol 49. P. J. Vinken, G. W. Bruyn, and H. L. Klawans, editors. Elsevier, Amsterdam. 391-415.
10. Lemieux, R. U., and R. M. Ratcliffe. 1979. The azidonitration of tri-O-acetyl-D-galactal. *Can. J. Chem.* 57:1244-1251.
11. Wang, A. M., D. F. Bishop and R. J. Desnick. In press. Human  $\alpha$ -N-acetylgalactosaminidase: molecular cloning, nucleotide sequence and expression of a full-length cDNA. *J. Biol. Chem.*
12. Sweeley, C. C., N. C. LeDonne, and P. W. Robbins. 1983. Post-translational processing reactions involved in the biosynthesis of lysosomal  $\alpha$ -N-

- acetylgalactosaminidase in cultured human fibroblasts. *Arch. Biochem. Biophys.* 223:158-165.
13. Bernstein, H. S., D. F. Bishop, K. H. Astrin, R. Kornreich, C. M. Eng, H. Sakuraba, and R. J. Desnick. 1989. Fabry disease: Six gene rearrangements and an exonic point mutation in the  $\alpha$ -galactosidase gene. *J. Clin. Invest.* 83:1390-1399.
  14. Bishop, D. F., and R. J. Desnick. 1981. Affinity purification of  $\alpha$ -galactosidase A from human spleen, placental and plasma with elimination of pyrogen contamination. *J. Biol. Chem.* 256:1307-1316.
  15. Aldridge, J., L. Kunkel, G. Bruns, U. Tantravahi, M. Lalande, T. Brewster, E. Moreau, M. Wilson, W. Bromley, T. Roderick, and S. A. Latt. 1984. A strategy to reveal high frequency RFLPs along the human X chromosome. *Am. J. Hum. Genet.* 36:546-564.
  16. Southern, E. M. 1975. Detection of specific sequences among DNA fragments separated by gel electrophoresis. *J. Mol. Biol.* 98:503-517.
  17. Maniatis, T., E. F. Fritsch, and J. Sambrook. 1982. *Molecular Cloning: A Laboratory Manual.* Cold Spring Harbor Laboratories, New York.
  18. Chrigwin, J. M., A. E. Pryzbyla, R. J. MacDonald, and W. J. Rutter. 1979. Isolation of biologically active RNA from sources enriched in ribonuclease. *Biochemistry.* 18:5294-5298.
  19. Lehrach, H., D. Diamond, J. M. Wozney, and H. Doedtke. 1977. RNA molecular weight determinations by gel electrophoresis under denaturing conditions, a critical look. *Biochemistry.* 16:4743-4751.
  20. Saiki, R. K., S. Scharf, F. Faloona, K. B. Mullis, G. T. Horn, H. A. Erlich, and N. Arnheim. 1985. Enzymatic amplification of  $\beta$ -globin genomic sequences and restriction site analysis for diagnosis of sickle cell anemia. *Science (Wash. DC).* 230:1350-1359.
  21. Messing, J., and J. Vieira. 1982. A new pair of M13 vectors for selecting either DNA strand of double digest restriction fragments. *Gene (Amst.).* 19:269-276.
  22. Sanger, F., A. R. Coulson, B. G. Garrell, A. J. H. Smith and B. A. Roe. 1980. Cloning in single-stranded bacteriophage as an aid to rapid DNA sequencing. *J. Mol. Biol.* 143:161-178.
  23. Theophilus, B., T. Latham, G. A. Grabowski, and F. I. Smith. 1989. Gaucher disease: Molecular heterogeneity and phenotype-genotype correlations. *Am. J. Hum. Genet.* 45:212-225.
  24. Chou, P.Y., and G. D. Fasman. 1978. Prediction of the secondary structure of proteins from their amino acid sequence. *Adv. Enzymol.* 47:45-147.

25. Garnier, J., D. J. Osguthorpe, and B. Robson. 1978. Analysis of the accuracy and implications of simple methods for predicting the secondary structure of globular proteins. *J. Mol. Biol.* 120:97-120.
26. Wolf, H., S. Modrow, M. Motz, B. A. Jameson, G. Hermann, and B. Fortsch. 1988. An integrated family of amino-acid sequence-analysis programs. *CABIOS*. 4:187-191.
27. Wong, G. G., J. S. Witek, P. A. Temple, K. M. Wilkins, A. C. Leary, D. P. Luxenberg, S. S. Jones, E. L. Brown, R. M. Kay, E. C. Orr, C. S. Shoemaker, D. W. Golde, R. J. Kaufman, R. M. Hewick, E. A. Wang and S. C. Clarke. 1985. Human GM-CSF: Molecular cloning of the complementary DNA and purification of the natural and recombinant proteins. *Science (Wash, DC)* 228:810-813.
28. Chen, C., and H. Okayama. 1987. High efficiency transformation of mammalian cells by plasmid DNA. *Mol. Cell. Biol.* 7:2745-2752.
29. Ritchie, H. H., M. R. Hughes, E. T. Thompson, P. J. Malloy, Z. Hochberg, D. Feldman, J. W. Pike, and B. W. O'Malley. 1989. An ochre mutation in the vitamin D receptor gene causes hereditary 1,25-dihydroxyvitamin D<sub>3</sub>-resistant rickets in three families. *Proc. Natl. Acad. Sci. USA.* 86:9783-9787.

**CHAPTER THREE:****STRUCTURAL ORGANIZATION OF THE HUMAN  
 $\alpha$ -N-ACETYL GALACTOSAMINIDASE GENE: CONSERVATION OF  
GENOMIC STRUCTURE WITH HUMAN  $\alpha$ -GALACTOSIDASE A**

**ABSTRACT**

Human  $\alpha$ -N-acetylgalactosaminidase ( $\alpha$ -GalNAc; E.C. 3.2.1.49), the lysosomal glycohydrolase that cleaves  $\alpha$ -N-acetylgalactosaminyl moieties from glycoconjugates, is encoded by a gene localized to chromosome 22q13→ter. The deficient activity of this enzyme results in Schindler disease, an autosomal recessive disorder characterized by the increased urinary excretion of glycopeptides and oligosaccharides containing  $\alpha$ -N-acetylgalactosaminyl moieties. The cDNA sequence encoding  $\alpha$ -GalNAc has been isolated and its characterization revealed significant nucleotide and amino acid homology (55.8% and 46.9%, respectively) to  $\alpha$ -galactosidase A ( $\alpha$ -Gal A), a lysosomal hydrolase localized to Xq21.33→22. To determine the structural organization of the  $\alpha$ -GalNAc gene and to compare its structural features to that of  $\alpha$ -Gal A, the  $\alpha$ -GalNAc chromosomal gene was isolated from a genomic cosmid library using the radiolabelled full-length  $\alpha$ -GalNAc cDNA. One clone, gAGB-1, which contained the entire  $\alpha$ -GalNAc gene, was characterized by restriction mapping and sequencing. The  $\alpha$ -GalNAc gene was ~14 kb and contained nine exons. All exon/intron junctions conformed to the GT/AG rule. Analysis of 1432 bp of 5' flanking sequence revealed two Sp1 and three GC-box promoter elements in this lysosomal housekeeping gene. The region was GC-rich (56%), but no HTF-island (Hpa II tiny fragments) was identified. Six *Alu*-repetitive elements were identified and all were in the reverse orientation. Comparison of the  $\alpha$ -GalNAc gene with the  $\alpha$ -Gal A gene revealed homologous exonic placement of  $\alpha$ -GalNAc introns 2 through 7 with  $\alpha$ -Gal A introns 1 through 6. Two additional introns, 1 and 8, were identified in the  $\alpha$ -GalNAc gene. Predicted amino acid homology among  $\alpha$ -GalNAc exons 2-7 with  $\alpha$ -Gal A exons 1-6 ranged from 46.2% to 62.7% with a few short gaps. In contrast, there was little, if any, similarity between the remaining coding sequence ( $\alpha$ -Gal A exon 7 and  $\alpha$ -GalNAc exons 8 and 9) which had only 15.8% homology with numerous gaps. In addition, no similarities were found between  $\alpha$ -

GalNAc and  $\alpha$ -Gal A intronic or flanking sequences. The high exonic homologies, in addition to homologous intron placement, suggest that these genes are evolutionarily related and arose through duplication and divergence from a common ancestral gene.

## INTRODUCTION

Human  $\alpha$ -*N*-acetylgalactosaminidase (E.C. 3.2.1.29;  $\alpha$ -GalNAc) is a lysosomal glycohydrolase that cleaves glycoconjugates with  $\alpha$ -*N*-acetylgalactosaminyl moieties (1-3). The human enzyme, encoded by a gene mapped to chromosome 22q13→ter, has been purified to homogeneity from various sources, and its physical and kinetic properties have been determined (1-6). The subunits of this homodimeric enzyme are synthesized as 65 kDa precursor glycopeptides, and following modifications in the Golgi, the mature 48 kDa subunits dimerize to form the active enzyme of 117 kDa (7).

The deficient activity of  $\alpha$ -GalNAc results in Schindler disease, an autosomal recessive disorder characterized by the increased urinary excretion of glycopeptides and oligosaccharides containing  $\alpha$ -*N*-acetylgalactosaminyl moieties (8). Two subtypes of this enzymatic deficiency have been identified: Type I disease has been described in two children who have the phenotype of an infantile neuroaxonal dystrophy (9), whereas Type II disease occurred in an adult who is neurologically asymptomatic, but has angiokeratoma corporis diffusum (10, 11). The patients with Type I or II disease have less than 1% of normal enzymatic activity and no immunoreactive enzyme protein in fibroblast extracts. Subsequent delineation and identification of the molecular lesions in these two subtypes revealed that each was due to a point mutation that resulted in an amino acid substitution (11, 12). However, the physiologic basis for the remarkable clinical heterogeneity of the two distinct subtypes is unknown.

Recently, we reported the isolation and nucleotide sequence of a cDNA encoding human  $\alpha$ -GalNAc (13). The predicted amino acid sequence of the cDNA had significant homology with that of human  $\alpha$ -galactosidase A ( $\alpha$ -Gal A) (13), suggesting the evolutionary relatedness of these two genes. In this communication, we report the isolation, characterization and sequence of the human  $\alpha$ -GalNAc chromosomal gene. The structure, 5' regulatory elements,

**3' flanking sequence, repetitive elements, and structural similarities with the human  $\alpha$ -Gal A gene are described.**

## EXPERIMENTAL PROCEDURES

*Construction of Synthetic Oligonucleotide Probes.* Unique oligonucleotides for sequencing the human  $\alpha$ -GalNAc gene and for use as polymerase chain reaction (PCR) primers were synthesized on an Applied Biosystems Model 380B.

*Isolation of Genomic Clones.* To isolate genomic clones containing the entire  $\alpha$ -GalNAc sequence,  $1 \times 10^6$  recombinants from a human genomic cosmid library were screened by colony hybridization with the full-length cDNA, pAGB-3, as radiolabelled probe (13). The genomic library was prepared from size-selected human lymphoblast DNA and kindly provided by Dr. Henrik Vissing, Mount Sinai School of Medicine.

*Characterization of Genomic Clones.* Cosmid DNA was isolated from purified positive clones, digested with various restriction endonucleases, separated by agarose gel electrophoresis, transferred to nylon membranes and then probed with end-labelled oligonucleotides corresponding to the pAGB-3 cDNA sequence to identify exonic sequences (14). Fragments which contained the entire genomic sequence were subcloned into the pGEM4Z vector (Promega, Madison, WI). Approximate intronic sizes were determined from the PCR products generated by amplification (15) of selected regions of the cosmid DNA with the appropriate flanking exonic primers. The GeneAmp DNA Amplification Reagent Kit (Perkin Elmer Cetus, Norwalk, CT) was used to amplify 1  $\mu$ g of cosmid DNA using 1  $\mu$ M of each primer. Each of 30 PCR cycles consisted of 1 min denaturation at 94  $^{\circ}$ C, 2 min annealing at 37  $^{\circ}$ C; and 7 min extension at 60  $^{\circ}$ C. An aliquot (10-20  $\mu$ l) was analyzed by agarose gel electrophoresis using *Hind* III digested lambda and *Hae* III digested  $\phi$ X174 DNAs as size standards.

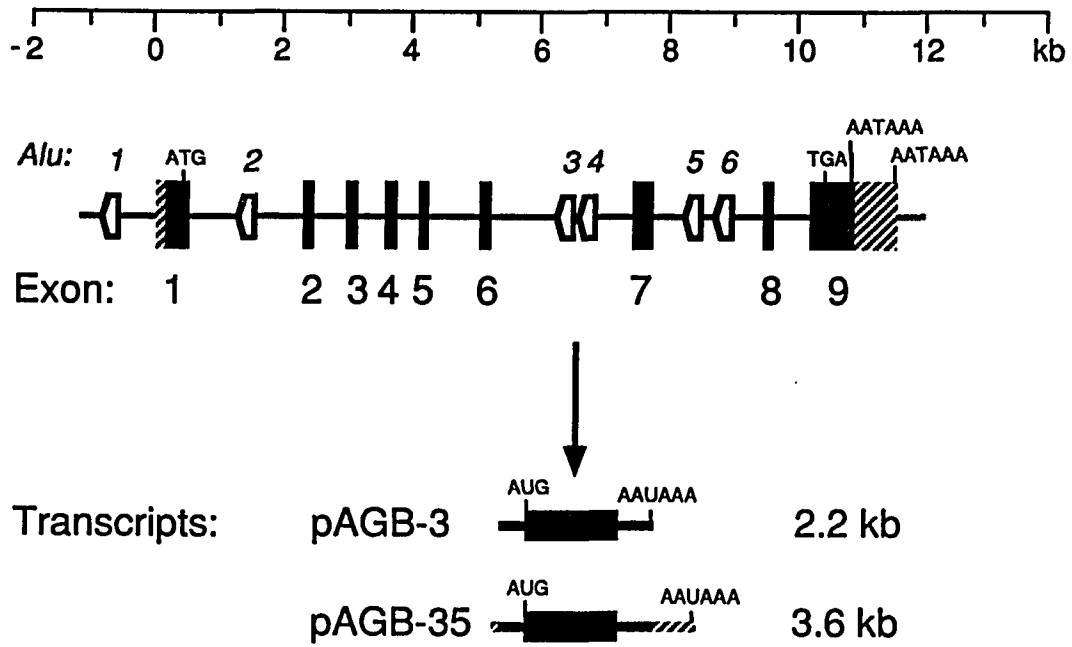
*DNA Sequencing and Computer-Assisted Analyses.* All DNA sequencing reactions were carried out by the dideoxy method (16) using universal or  $\alpha$ -GalNAc-specific synthetic oligonucleotide primers in both orientations using Sequenase (U.S. Biochemical Corp, Cleveland,

**OH). Searches for nucleotide and amino acid sequence similarities were carried out with the Microgenie DNA analysis program (Beckman, Fullerton, CA).**

## RESULTS AND DISCUSSION

*Isolation and Restriction Mapping of Genomic Clones.* A total human genomic DNA cosmid library was screened with the nick-translated pAGB-3 cDNA insert. Of  $1 \times 10^6$  independent recombinants screened, six putative positive clones were isolated, purified, and subjected to Southern hybridization with end-labelled oligonucleotides corresponding to  $\alpha$ -GalNAc cDNA sequences to identify exonic fragments. One clone, gAGB-1, which had an ~35 kb insert and appeared to contain the entire genomic sequence, was selected for further analysis. The entire  $\alpha$ -GalNAc genomic sequence occurred in a ~14 kb *EcoR* I fragment which was subcloned into the pGEM4Z vector and designated gAGB-1E.

*Intron-Exon Boundaries.* Since the  $\alpha$ -GalNAc and  $\alpha$ -Gal A cDNA sequences were highly homologous and suggested that the two genes may have evolved from a common ancestral sequence (13), it was expected that the genomic structure of the two genes would be conserved. Thus, initial efforts were directed to determine the approximate size and location of  $\alpha$ -GalNAc introns. Exonic oligonucleotides corresponding in position to the  $\alpha$ -Gal A intron/exon junctions (14, 17) were used to PCR-amplify the homologously located  $\alpha$ -GalNAc intron/exon regions. Sequencing of the PCR products revealed that all six corresponding  $\alpha$ -Gal A intron/exon junctions were conserved in the  $\alpha$ -GalNAc gene. This was analogous to the previously demonstrated conservation of the gene structures for the  $\beta$ -hexosaminidase  $\alpha$ - and  $\beta$ -chain genes, where 12 of 13 introns in the  $\beta$ -chain interrupted the gene at homologous positions in the  $\alpha$ -chain (18). Additional sequencing of exonic regions revealed the presence of two additional introns interrupting the  $\alpha$ -GalNAc gene, one 5' and one 3' to the six previously identified introns. To further characterize the  $\alpha$ -GalNAc gene, the entire 14 kb gAGB-1E fragment was sequenced in both orientations (Fig 1). The nine  $\alpha$ -GalNAc exons



*Figure 1.* Organization of the  $\alpha$ -GalNAc gene.

ranged in length from 95 to 649 bp, with exon 1 containing the 347 nt 5' untranslated region and the sequence encoding the first six amino acids of the 17-residue signal peptide. The eight introns ranged in length from 0.3 to 2.7 kb. All splice junctions followed the GT/AG rule (19) and were consistent with the 5' and 3' consensus sequences for splice junctions of RNA polymerase II (20). Putative lariat branch points were identified between -40 and -18 nt from the 3' splice junction for all eight introns by similarity to the less well conserved consensus sequence (C/T)N(C/T)T(A/G)A(A/C/T) (21). All three codon phases were observed at the exon junctions in the  $\alpha$ -GalNAc gene (Table 1). The  $\alpha$ -GalNAc genomic sequences also confirmed at the genomic level the presence of two polyadenylation hexamers of AATAAA whose alternative use results in the two  $\alpha$ -GalNAc transcripts of 2.2 and 3.6 kb (13) (Fig. 1).

*Regulatory Elements.* Computer-assisted analysis of the  $\alpha$ -GalNAc 1432 bp 5' flanking region identified two Sp1 binding sites at -26 and -258 in sense and antisense orientations, respectively, and three GC boxes (at -23, -68, and -257) from the transcription start site (22) (Fig. 2). No CAAT or TATA-like elements were found. The 5' flanking region was GC-rich (~56%) but there was not an HTF-island. The CpG/GpC ratio was 0.625, whereas an HTF-island has a ratio greater than 1.0 (23, 24). The absence of a TATA box and the high GC content are known features in the promotor region of other housekeeping genes (25). Imperfect inverted repeats were identified at nt -108 to -127 and at -27 to -17 which could pair with a predicted  $\Delta G$  of -34.8 kCal/mol (Fig. 2). The structure predicted by this pairing may be important to promotor function. Further evaluation of these putative promotor elements will require footprinting and functional assays with deletion/mutation mapping.

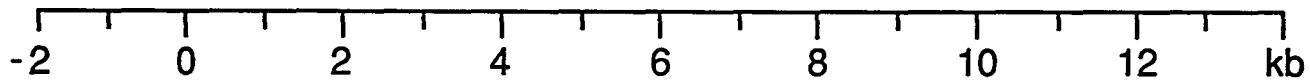
*Alu-Repetitive Elements.* Six *Alu*-repetitive elements (26) were identified in the 14 kb sequence: one each in the 5' flanking region and intron 1, and two in introns 6 and 7 (Fig. 3).



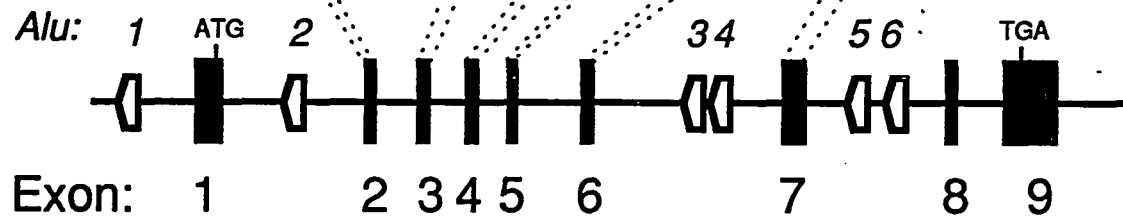
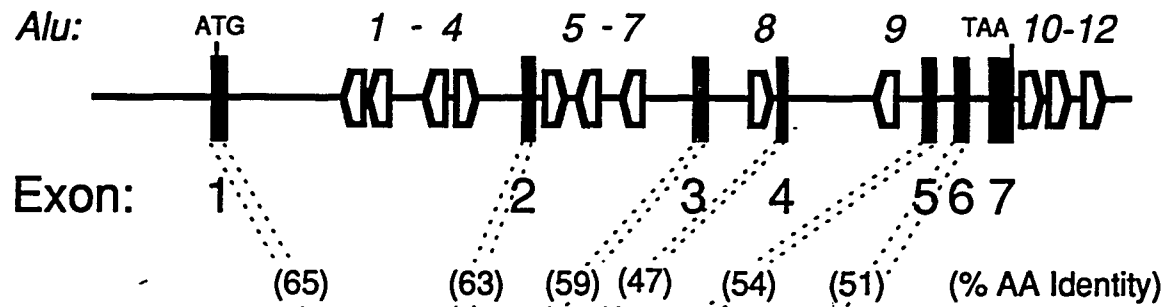
*Figure 2.* The nucleotide sequence of the  $\alpha$ -GalNAc gene 5' flanking region. The 1532 nt sequence was determined in both orientations from genomic clone, gAGB-1. The sequence includes 1432 nt of 5' flanking sequence and 100 nt of 5' untranslated sequence from exon 1. The G at nt +1 indicates an transcriptional initiation site as determined by primer extension (13). Sp1 binding sites (Sp1) are underlined with wavy arrows. A sequence corresponding to a GC box is boxed. An imperfect inverted repeat (IR) is underlined. An *Alu*-repetitive element from nt -958- -1234 is italicized.

-1432 GA GCTCGCCCGG GGATCTACCT TGCCTGATAC -1401  
 -1400 TAGGTCATAA GACTCTCATT CCAGAAGGGG TCCTGCCCTC CTCCCAGCAG GAGGGAATGC CACACAGAGA GAACAAGAAG AATCCGAATC GAACAGACTA -1301  
 -1300 TTGCTGGGTT TCCACACTCA TTAAATCATT CCCTTTTGTC TAATCGCGTT TGAGCGCAGT TGCCATTCC TTATTCCTTT TTTTTTTTTC CTATTGCCCG -1201  
 -1200 GGCTGGAGTA CAGTGACACA ATCACAGCTC ACTGCAACCT CGACCTCTGG ACTCAATCAA TCCTCCCACC TCAGCCTCCC GAGTAGCTGG GACCACAGGC -1101  
 -1100 ACCTGCCTGG CTTTTAATT TTTGTAGAG ACAGGTCTC ACTATGTTGC CCAGACTGTC GTTAAATGTC TGGCCTCAAG CGATCCCCCT GCCTCAGCCT -1001  
 -1000 TCCAAAGTTC TGGGATTACA GGCATGAGCC ACCAGACCCG GCCAGGGCTG TCCATTCTTC ATCAAACCGA GACATAAAAA TAAAGATTC CCTGAGTATT -901  
 -900 TGTGTCTTCA TTCCTGAAGG CTCTCAAGTC ACATAAACA TTGAGTAAAT CAATCTATTA TGCTTTTCTC TTGCTAACCT ATATTTTATT TGTAGGAGTG -801  
 -800 TGCGCTGTGA CCTTGGGAT GAGTAAGGGA AAATATCACA CCTTCTGCC CCTACAATGC CCAGAACAAA GCATTTGGCA CTTAGCAGTT TCTCAATAAA -701  
 -700 TACTTGATGA ATGAATGAAG AAGGAAGTAG CCCAGTCTGT GCGAGTAGTC CCTGAGCGTC AGGCGGGGAA GGAGAGCCGG CTGCCTTTC GGGAGCTGTG -601  
 -600 ACCGCCTGCT TTGGTGCCTC CCACGCACCC TCCTTGTAGA GACGATTAAT CCCAGACTTC CATTTTCTTA GCTAGCTCAC TAGGCGGGCC TTGCTTCAGA -501  
 -500 GCAGTCTCTC ATTCCTCTG CACAAAGCCC TCAACCGTCC TCTGCACGGC CGCTTACTGA CCCACAGAG GTGAAGACAG GCTCGCGGGG TCCGGGCCAC -401  
 -400 GAAGGACAGC CCCGAGTCTC GCGCCACCG TGAGCCTGGC ACGATCCTTG CCCTCTTGGT TTCTGCCAAT GCCAGCCTCT CTAGCACGCA GCGGGAGCC -301  
 -300 CAATGGGGAA GGGCTTAGGG AGGCTGGGAG GCCGGACCCA TATGGGCGGG GCACAGCGTT GCTAGGGGGG ATTCTGATAG GTTGC GCGAA GGGTGCCTT -201  
 -200 CGTGAGGCC CCCCTCACTG GAACGCTTCG GACTTATCAG GTTACCGGAT TCGAGTCAGA AGCGGGGCA GGTCTGAACG CTTCCGTAGC ATAAGTCCG -101  
 -100 TGGGAGGGG CTCTGCGGTT TTGGGACCC AGGGCGGAA CCCGACCCG GCTGGAAGTC CGGAGCCTGC CGCAGCCCCG CCCCTCCGCT CTTCTTGGT -1  
 IR GC Box Sp1 IR Sp1  
 1 GACCTTAAGC CAGTGGCTGC CTTTTCTGA GCCCGGGCGG GGCCGAAGGC GCCCGTAGGC CCTCGGGACT CCCAGCACTG CAGAGGGTGT GAGGTCTGAC 100  
 Exon 1

*Figure 3.* Comparison of the  $\alpha$ -Gal A and  $\alpha$ -GalNAc structural genes. Exons are denoted by solid rectangles. The positions and orientations of the *Alu* repeats are indicated. The homology between  $\alpha$ -Gal A exons 1-6 and  $\alpha$ -GalNAc exons 2-7 are denoted by dashed lines and the amino acid identity indicated in parentheses.



**$\alpha$ -GALACTOSIDASE A GENE**



**$\alpha$ -N-ACETYL GALACTOSAMINIDASE GENE**

All six were in the reverse orientation. The percent homology of these repetitive elements with the consensus 300 nt sequence ranged from 54.1 to 87.4%.

*Structural Homology of the  $\alpha$ -GalNAc and  $\alpha$ -Gal A Genes; Potential Evolutionary Relationships.* Sequencing of the  $\alpha$ -GalNAc structural gene revealed that six of the eight  $\alpha$ -GalNAc introns (2 through 7) interrupted the gene at positions that were identical to those interrupting the  $\alpha$ -Gal A coding sequence. Characterization of the exonic sequences encoding these two functionally related genes revealed 55.8% and 46.9% overall homology at the nucleotide and deduced amino acid levels. Analysis of the amino acid identity between homologous exons demonstrated 46.2% to 62.7% identity with only a few gaps between  $\alpha$ -GalNAc exons 2 through 7 and  $\alpha$ -Gal A exons 1 through 6 (Table 2). In contrast, the homology between  $\alpha$ -GalNAc exons 8 and 9 with  $\alpha$ -Gal A exon 7 dropped significantly to 15.8% with multiple gaps. The homologous intron placement and exonic amino acid identity suggested that the  $\alpha$ -GalNAc and  $\alpha$ -Gal A genes had a common origin, presumably they were duplicated from the same ancestral gene and underwent subsequent divergence to acquire their specific functions (18, 27).

A possible mechanism for the duplication of these genes is a recombinational event, perhaps involving *Alu*-repetitive elements. Although the  $\alpha$ -GalNAc and  $\alpha$ -Gal A genes have five and 12 *Alu*-repetitive elements, respectively, *Alu-Alu* mediated recombination probably was not the mechanism for duplication since the *Alu* sequences were not homologously placed in these genes, nor was there any homology in intron size or sequence or flanking sequences (Table 2) of  $\alpha$ -GalNAc and  $\alpha$ -Gal A. Comparison of the deduced amino acid sequences in  $\alpha$ -galactosidase cDNAs from human, mouse, yeast, and *E. coli*, and from the human  $\alpha$ -GalNAc cDNA (13, 28), suggested that  $\alpha$ -GalNAc exons 2 through 7 were conserved due to functional

**Table 2.** Comparison of the exonic and intronic homologies between  $\alpha$ -GalNAc and  $\alpha$ -Gal A.

$\alpha$ -GalNAc exon number and (size)	$\alpha$ -Gal A exon number and (size)	% Amino acid identity	$\alpha$ -GalNAc intron number and (size)	$\alpha$ -Gal A intron number and (size)
(nt)	(nt)		(kb)	(kb)
1 (363)	—	—	1 (1.7)	—
2 (136)	1 (254)	46.2	2 (0.5)	1 (3.8)
3 (172)	2 (175)	62.7	3 (0.5)	2 (2.0)
4 (178)	3 (178)	59.0	4 (0.3)	3 (1.5)
5 (95)	4 (92)	46.9	5 (0.8)	4 (1.2)
6 (162)	5 (162)	53.7	6 (2.7)	5 (0.2)
7 (198)	6 (198)	51.3	7 (1.7)	6 (0.4)
8 (144)	7 (291)	15.8	8 (0.8)	—
9 (649)				

domains common to these enzymes. In support of this concept are regions of up to 15 consecutive identical residues and stretches of 21 of amino acids with 20 identical residues in the human  $\alpha$ -GalNAc and  $\alpha$ -Gal A polypeptides. These domains may have been under strict functional constraints, while other domains ( $\alpha$ -GalNAc exons 8 and 9, and  $\alpha$ -Gal A exon 7) may have more freely diverged to generate unique functional specificities such as *N*-acetylgalactosamine recognition. Because the homology between  $\alpha$ -GalNAc and  $\alpha$ -Gal A falls sharply after  $\alpha$ -GalNAc exon 7, it is intriguing to speculate on the origin of  $\alpha$ -GalNAc exons 8 and 9. When the predicted amino acid sequence encoded by  $\alpha$ -GalNAc exons 8 and 9 were used to search protein sequences (13), the only homology identified was a region of seven  $\alpha$ -GalNAc residues (365-371) which was homologous to residues 194-200 in the  $\beta$ -hexosaminidase  $\alpha$ -chain, a subunit of the human lysosomal  $\beta$ -*N*-acetylgalactosaminidase (18). Thus, the functional specificity for a possible *N*-acetylgalactosaminyl recognition, binding, or catalytic site may have arisen by "exon shuffling" (2, 30). This concept would require that both  $\alpha$ -GalNAc exons 8 and 9 were shuffled together since  $\alpha$ -GalNAc intron 8 interrupted the gene between the nucleotides encoding residues 367 and 368.

Speculations regarding the divergence time and conserved structural domains between  $\alpha$ -GalNAc and  $\alpha$ -Gal A will require  $\alpha$ -Gal A and  $\alpha$ -GalNAc sequences from other species. The isolation and characterization of these genes will allow the calculation of their divergence times. X-ray crystallographic studies of purified recombinant  $\alpha$ -Gal A and  $\alpha$ -GalNAc should lend insight into the conserved structural domains between the two lysosomal glycohydrolases.

## REFERENCES

1. Dean, K. J., Sung, S. and Sweeley, C. C. (1977) *Biochem. Biophys. Res. Comm.* **77**, 1411-1417
2. Schram, A. W., Hamers, M. N. and Tager, J. M. (1977) *Biochim. Biophys. Acta.* **482**, 138-144
3. Bishop, D. F., Dean, K. J., Sweeley, C. C. and Desnick, R. J. (1980) in *Enzyme Therapy in Genetic Disease:2*, (Desnick, R. J., ed.), pp. 17-32. Alan R. Liss, Inc. New York
4. Beutler, E. and Kuhl, W. (1972) *J. Biol. Chem.* **247**, 7195-7200
5. Callahan, J. W., Lasilla, E. L., DenTandt, W. and Philippart, M. (1973) *Biochem. Med.* **7**, 424-431
6. Kusiak, J. W., Quirk, J. M. and Brady, R. O. (1978) *J. Biol. Chem.* **253**, 184-190
7. Sweeley, C. C., LeDonne, N. C. and Robbins, P. W. (1983) *Arch. Biochim. Biophys.* **223**, 158-165
8. van Diggelen, O. P., Schindler, D., Willemsen, R., Boer, M., Kleijer, W. J., Huijmans, J. G. H., Blom, W. and Galjaard, H. (1988) *J. Inher. Met. Dis.* **11**, 349-357
9. Schindler, D., Bishop, D. F., Wolfe, D. E., Wang, A. M., Egge, H., Lemieux, R. U. and Desnick, R. J. (1989) *N. Engl. J. Med.* **320**, 1735-1740
10. Kanzaki, T., Yokota, M., Mizuno, N., Matsumoto, Y. and Hirabayashi, Y. (1989) *Lancet* April 22, 1989, 875-877
11. Wang, A. M., Kanzaki, T., Schindler, D. and Desnick, R. J. (1989) *Am. J. Hum. Genet.* **43**, 4 A228
12. Wang, A. M. and Desnick, R. J. (in review) *J. Clin. Invest.*
13. Wang, A. M., Bishop, D. F. and Desnick, R. J. (in review) *J. Biol. Chem.*
14. Bishop, D. F., Kornreich, R. and Desnick, R. J. (1988) *Proc. Natl. Acad. Sci. USA* **85**, 3903-3907
15. Saiki, R. K., Gelfand, D. H., Stoffel, S., Scharf, S. J., Higuchi, R., Horn, G. T., Mullis, K. B. and Erlich, H. A. (1987) *Science* **239**, 487-491
16. Sanger, F., Nicklen, S. and Coulson, A. R. (1977) *Proc. Natl. Acad. Sci. USA* **74**, 5463-5467
17. Kornreich, R., Desnick, R. J. and Bishop, D. F. (1989) *Nucleic Acids Res.* **14**, 3301-3302

18. Proia, R. L. (1988) *Proc. Natl. Acad. Sci. USA* **85**, 1883-1887
19. Breathnach, R. and Chambon, P. (1981) *Ann. Rev. Biochem.* **50**, 349-383
20. Mount, S. M. (1982) *Nucleic Acids Res.* **10**, 459-472
21. Reed, R. and Maniatis, T. (1988) *Genes Dev.* **2**, 1268-1276
22. Dynan, W. S. and Tijan, R. (1985) *Nature* **316**, 774-778
23. Bird, A. P. (1986) *Nature* **321**, 209-213
24. Lindsay, S. and Bird, A. P. (1987) *Nature* **327**, 336-338
25. Dynan, W. S. (1986) *Trends Genet.* **2**, 196-197
26. Deininger, P. L., Jolly, D. J., Rubin, C. M., Friedmann, T. and Schmid, C. W. (1981) *J. Mol. Biol.* **151**, 17-33
27. Maeda, N. and Smithies, O. (1986) *Ann. Rev. Genet.* **20**, 81-108
28. Unpublished results
29. Gilbert, W. (1978) *Nature* **271**, 501
30. Michelson, A. M., Blake, C. C., Evans, S. T. and Orkin, S. H. (1985) *Proc. Natl. Acad. Sci. USA* **82**, 6965-6969

## CONCLUDING REMARKS

$\alpha$ -N-acetylgalactosaminidase ( $\alpha$ -GalNAc) is representative of a class of enzymes present in lysosomes. Genetic deficiencies of specific lysosomal hydrolases result in different lysosomal storage diseases. To date, over 30 lysosomal storage disorders have been described in man, and have sparked wide interest in lysosomal cell biology and the biochemistry of lysosomal enzymes. Biochemical and processing studies of  $\alpha$ -GalNAc and other lysosomal hydrolases have now well defined the proteolytic, carbohydrate, and phosphorylation modifications necessary for the appropriate trafficking of these enzymes to the lysosomes. The deficient activity of  $\alpha$ -GalNAc results in Schindler disease. Two distinct clinical subtypes of this enzymatic deficiency have been described: Type I is characterized by a clinical course identical to infantile neuroaxonal dystrophy, whereas Type II is characterized by disseminated angiokeratoma corporis diffusum with no neurological involvement.

The revolution in molecular biology has permitted the study of lysosomal enzymes and their deficiencies at the molecular level. Approximately 20 cDNA sequences encoding lysosomal proteins have been isolated. As a result, a myriad of mutations causing specific lysosomal disorders have been identified and characterized and determined to be heterogeneous. The studies described here provide the molecular basis for the further study of human  $\alpha$ -GalNAc. The similarities of  $\alpha$ -GalNAc with  $\alpha$ -galactosidase A ( $\alpha$ -Gal A) should provide insight into the potential structural and functional domains as well as into their evolutionary relationships. As the sequences encoding  $\alpha$ -GalNAc and  $\alpha$ -Gal A from other species are isolated, it should be possible to determine the divergence time of these enzymes. Furthermore, X-ray crystallographic studies on recombinant  $\alpha$ -GalNAc and  $\alpha$ -Gal A should provide information on the conserved structural and functional domains of these proteins, as well as potential differences in domains responsible for functional specificity. Schindler disease has

provided an interesting study of  $\alpha$ -GalNAc deficiency. While the mutations and their resulting amino acid substitutions in each subtype have been identified and expressed, the physiological basis for the remarkable clinical heterogeneity between the two subtypes is not understood. However, the characterization of the  $\alpha$ -GalNAc cDNA sequence and knowledge of the specific base substitution in each Schindler disease subtype should permit further delineation of the enzyme's function, stability, and potential specificity and protein interactions.

In conclusion, these studies have provided the framework for the further study and understanding of the properties of human  $\alpha$ -GalNAc and the pathophysiology of its deficiency.