

Isolation and characterization of a pseudogene related to human core 2 β -1,6-*N*-acetylglucosaminyltransferase

MARTI F.A. BIERHUIZEN*, KENTARO MAEMURA‡ and MINORU FUKUDA§

Glycobiology Program, La Jolla Cancer Research Foundation, La Jolla, California 92037, USA

Received 26 April 1995, revised 17 May 1995

In a previous study, we isolated genomic clones encoding core 2 β -1,6-*N*-acetylglucosaminyltransferase (C2GnT) and blood group IGnT and proposed that these two genes were produced from a common ancestral gene by duplication, diversion and intron insertion. In the present study, we have isolated a pseudogene which is highly related to the gene of C2GnT. The sequence analysis of this pseudogene indicated that the pseudogene was produced by duplication of a common precursor gene for C2GnT. These results taken together strongly suggest that the ancestral gene was first duplicated and one of the duplicated genes directly evolved into the IGnT gene. The other duplicated gene was further duplicated to produce the C2GnT gene and the pseudogene.

Keywords: core 2 β -1,6-*N*-acetylglucosaminyltransferase, pseudogene, gene evolution

Introduction

Cell surface carbohydrates are often characteristic of different cell lineage and different stages of differentiation [1–3]. Among these carbohydrates, those synthesized by β -1,6-*N*-acetylglucosaminyltransferases are particularly specific to cell-types. For example, a common *O*-glycan branch is formed by core 2 β 1,6-*N*-acetylglucosaminyltransferase (C2GnT) and the addition of the β -1,6-*N*-acetylglucosaminyl linkage to the Gal β 1 \rightarrow 3GalNAc backbone results in the formation of the hexasaccharide, NeuNAc α 2 \rightarrow 3Gal β 1 \rightarrow 3 (NeuNAc α 2 \rightarrow 3Gal β 1 \rightarrow 4GlcNAc β 1 \rightarrow 6)GalNAc [4, 5]. When β -1,3-*N*-acetylglucosaminyltransferase (extension enzyme) is also present, the β -1,6-*N*-acetylglucosaminyl branch is further extended to form poly-*N*-acetyllactosaminyl side chains in *O*-glycans [6, 7]. The poly-*N*-acetyllactosaminyl side chain is often modified to have a sialyl Le^x terminus, NeuNAc α 2 \rightarrow 3Gal β 1 \rightarrow 4(Fuc-

α 1 \rightarrow 3)GlcNAc \rightarrow R [4, 6]. A dramatic increase in core 2 branching was observed when human T-lymphocytes were activated from the resting state [5] and this increase leads into the expression of sialyl Le^x in activated T-lymphocytes [8]. The maturation of thymocytes from cortical to medullary thymus is associated with the turning off of C2GnT [9]. In pathological conditions such as leukaemia and immunodeficiency, leukocytes express an increased amount of C2GnT [10–13]. Moreover, highly metastatic tumour cells express much more branched oligosaccharides than low metastatic counterparts [14]. These results indicate that the increase of core 2 branches in hematolymphoid malignancy, and tumours in general, reflect the cell surface carbohydrates of immature cells. The results also suggest that core 2-based carbohydrates are involved in metastatic spreading of tumour cells.

During development of human erythrocytes, the blood group I-branching enzyme (IGnT), is substantially increased. While linear poly-*N*-acetyllactosamine (Gal β 1 \rightarrow 4GlcNAc β 1 \rightarrow 3)_n is expressed on fetal erythrocytes, it is replaced by branched, I-active poly-*N*-acetyllactosamine, Gal β 1 \rightarrow 4GlcNAc β 1 \rightarrow 3(Gal β 1 \rightarrow 4GlcNAc β 1 \rightarrow 6) on adult erythrocytes [15, 16]. Because of this importance, we have cloned cDNAs encoding the human C2GnT and IGnT genes [17, 18].

*Present address: Department of Hematology, Erasmus University, Rotterdam, The Netherlands.

‡Present address: Department of Internal Medicine, Osaka Medical College, Takatsuki, Japan.

§To whom correspondence should be addressed.

pseudogene C2GnT E1+E2	CTCGAGGRRCCCTTATCTCACCCATATACTAAAAATCAARACAAAAATGGATTAAACTGAAAACTTCTATCGGAAAAACATAGGGAAAAAGCTTATGACATT	100
pseudogene C2GnT E1+E2	AGTCTGGGCAATGATTTTTTTTCGGATATGACCCCCAAAGCACCGGCAAGCAAAAGRAAAAAATAGACAAACAGAAATACATCAAACATAAAAGCTTCTGTC	200
pseudogene C2GnT E1+E2	ACAGCAAAGGAAACAATGAAACADAATGTTAAG--AGACAACCTTALCAAAATGGA--AGAAAATAATGTT--CAAAACATACATC---TGAAAGGGGATTAATA -----TGGGGAAAGGATATGTTGGATTTTCTAAICATACTGGACTGTGTAACCTGTTCTCTGGTATACATTAAGGGAAAGCTGCATASA	289 84
pseudogene C2GnT E1+E2	TATTCAGTACATAAATA--GAACTCAAA---TA-----ACTTC--AATAGTARGAAA---ACAAATTAATBCCATT--AAAAATGGGAAA---AGGGAGCAG GTITTTCTTTTCAAAATAATTAACCTAAATTTTCAGGGTGTTTTCAAATTTTTCAGTAATCCACAAC--ATAGAAATTAGAAAATTAAGTAAGACTAGAGAGAG	369 183
pseudogene C2GnT E1+E2	---TCCCTGCATG--G---CTTTTCCTTA--GCA--AAAGCT---GATGGTA--GCGAA--CTGTGTGCTTTGGGGCCAGATCATGCTCT--AGGGTG--G---AATG GAAATTTTAAAGAAATAAAGTTAACCACTCAATTTTGGAAAAAAGCAATAATGTAGGATTTAGTTCAACCAACATATTTTCAGAGGTGTGTAATAAT	452 283
pseudogene C2GnT E1+E2	AAGTACCTGCAAAATCCAGCTGACCATTTGAAAGCCACCCAGGGAATGATTTTTCAAAAGCTTATGGGACAAATCATTATTCCTGGATGCCCTAGTACATTT AAGTGTCTATATAATCCCTTGTCAATTTTTCATTT--TCA--TAGATTTTTCAAAAGCTTGTGTCAGCAATCATTATTCCTGGATGCCCT--GACATGTT	552 380
pseudogene C2GnT E1+E2	AATTCCTGACAGCATGAAAGTCTTTCGAAATAGTTCAGGATGTCACCTGGAATCAGAGTCTTAAGTGATCTGACTTTCCTTAATTTTAAATGCGGC AAATACCTGACAGCATGTAAGTCTTTCGAAAT--GGG--CAGGATGTCACCTGGAATCAG--CACTAAGTGATTCAGACTTTCCTTAATTTTAAATGCGCT	652 477
pseudogene C2GnT E1+E2	AATTCCTCATTTCAAAGTGGCCTTTG--AGCTCTGATAAATGCAAACCTGACAACCTTCAAGGC--ACAAATGGAGGGTAAATAGTTGGTGTCTTAAGCCTAGAA GCTCTTCATTTCAAAGTGGCCTTTG--AGC--CTGATAAATGCAAACCTGACAACCTTCAAGGC--ACAAATGGAGGGTAAATAGTTGGTGTCTTGGAGCATAGAA	750 576
pseudogene C2GnT E1+E2	GACTGCCCTTAC--AAG--AAA--TCCCTGATTTGGCAATTTGAAATGCTGAGG--AATTTGCTGT----GGAGACNTTTTCTTATCCCCAC--AAATAC--ACTTTT GACTGCCCTTAC--AAG--AAA--TCCCTGATTTGCTTTGAAATGCTGAGG--ACCTTGTCTG--GAA--GGAGACNTTTTCTTATCCCCAC--AAATAC--ACTTTT	847 675
pseudogene C2GnT E1+E2	GGTCTTATTTTTCCTTATTCACC----TCCGTTTAAAGGATTCATCAAAAAGCTTCAATCTGTAAGCTCA--ATAATGTTGGAGCTTGTGGGAGAAATCCT GGTCTTATTTTTCCTTATTCACC----TCCGTTTAAAGGATTCATCAAAAAGCTTCAATCTGTAAGCTCA--ATAATGTTGGAGCTTGTGGGAGAAATCCT	944 775
pseudogene C2GnT E1+E2	AGTAGTATATTAATTCACCAAAGTTTAC--GGGG--GATGTA--ATGAAATCCAAAAGGTAAATCTTGAGATCTTAAACAGTGAAATTTACTAAGTGGCCCT AGTAGTATATTAATTCACCAAAGTTTAC--AGGGG--GATGTA--ATGAAATCCAAAAGGTAAATCTTGAGATCTTAAACAGTGAAATTTAAAGTGGCCCT	1044 874
pseudogene C2GnT E1+E2	TGGTGTATACCTGACGCTTTATAAACAATGACCAGTTATG--TAC--TCTTCAATGTA--CTTCTTTCATCAAGAGAG--GTA--ATAATATTTAGAACCCCTTAAGAA CGGTGTATACCTGACGCTTTATAAACAATGACCAG----TG--AC-----TGTCTTCTTTCATCAAGAGAG--GTA--ATAATATTTAGAACCCCTTAAGTAA	1144 963
pseudogene C2GnT E1+E2	GGAAGAGGTTAGTTTCCAATAGCATATTTCTATA--TGGTTCATTA--AAA--TGAAT--GCTTTGACAGGCTCTAGAGAGCCATCTATATGCTCAGAAATTT AGAAGAGGTTAGTTTCCAATAGCATATTTCTATA--TGGTTCATTA--AAA--TGAAT--GCTTTGACAGGCTCTAGAGAGCCATCTATATGCTCAGAAATTT	1244 1063
pseudogene C2GnT E1+E2	TATTGCATTTCATGTTGACAA--AAA--AAATCAGTATATTCCTTTT--TAGCTG--CAGTGA--TGGCATTG--GGTCA----TTT--AGTAA--TCTTTGTTGGCCCTGACAGTT TATTGCATTTCATGTTGACAA--AAA--AAATCAGTATATTCCTTTT--TAGCTG--CAGTGA--TGGCATTG--GGTCA----TTT--AGTAA--TCTTTGTTGGCCCTGACAGTT	1342 1161
pseudogene C2GnT E1+E2	GGAGAGTGTGGTTTATGCTTTTGGAGTGGGTTCTTGGCTGAGCTCAACTGCATGATGGA--CTCTTCCAA--TGAGTGCAT--ACTGGAAGTACTTAATAAT GGAGAGTGTGGTTTATGCTTTTGGAGTGGGTTCTTGGCTGAGCTCAACTGCATGATGGA--CTCTTCCAA--TGAGTGCAT--ACTGGAAGTACTTAATAAT	1442 1261
pseudogene C2GnT E1+E2	GTTTGTATGATGATTTTCCATATAAACAACCTA--AAATTTGTTAGGAAGCTCAAGTTGTTAATGGGGAA--ACA--TCTTAA--ACTTAAGAGGATGCCAT GTTTGTATGATGATTTTCCATATAAACAACCTA--AAATTTGTTAGGAAGCTCAAGTTGTTAATGGGGAA--ACA--TCTTAA--ACTTAAGAGGATGCCAT	1541 1361
pseudogene C2GnT E1+E2	CCATAAAGAAAGAAAGTGGAA--AAG--GGTATG--GATATA--TTAATGGAAGCTGACATA--TGTGGGACTGTCAA--GGG--CACTCTCCCTTGAAGCACCCTAT CCATAAAGAAAGAAAGTGGAA--AAG--GGTATG--GATATA--TTAATGGAAGCTGACATA--TGTGGGACTGTCAA--GGG--CACTCTCCCTTGAAGCACCCTAT	1641 1461
pseudogene C2GnT E1+E2	TTTTTC--GGCAGTGCCTATTTT--TGGTTCAGTAGGGAGTATGTTGGG--TATGTCTATAG--ATGAAAAA--CCAAAAGTTTATGGAGTGGGTTGAGGACACA CTTTTTC--GGCAGTGCCTATTTT--TGGTTCAGTAGGGAGTATGTTGGG--TATGTCTATAG--ATGAAAAA--CCAAAAGTTTATGGAGTGGGTTGAGGACACA	1741 1561
pseudogene C2GnT E1+E2	GACAGCCCAGATAGTATCTCTAGGCCATCATCC--AAGGATGCTGAAAGTCCC--GGCTCAATCTCTTTAGCCATAAGTAA--AA--TTCTGCTGAAATGCATG TACAGCCCAGATAGTATCTCTAGGCCATCATCC--AAGGATGCTGAAAGTCCC--GGCTCAATCTCTTTAGCCATAAGTAA--AA--TTCTGCTGAAATGCATG	1841 1661

and *IGnT* after labelling by random oligonucleotide priming [20] using [α - 32 P] dCTP and a kit from Boehringer Mannheim. The cDNA fragments were obtained by PCR amplification [21] of cDNA sequences which encode the beginning of the stem region of the protein to the 3'-untranslated region as described previously [19]. A probe for Southern hybridization was made by 32 P-labelling using random oligonucleotide priming [20].

Southern blot analysis

Phage DNAs were digested with various restriction enzymes and subjected to Southern blotting and hybridization as described previously [18]. Briefly, the blots were hybridized with 32 P-labelled cDNA inserts of *IGnT* or *C2GnT*. The hybridization was in $6 \times$ SSPE, pH 7.4, 0.5% SDS, $50 \mu\text{g ml}^{-1}$ of denatured, sheared salmon sperm DNA containing 50% formamide at 42 °C for 16 h [22]. The blot was then washed several times in $2 \times$ SSPE, pH 7.4, 0.5% SDS at room temperature for periods of 10 min and subsequently exposed to Kodak X-Omat AR film.

DNA sequencing

The DNA fragments of interest were subcloned into pcDNA1 (Invitrogen) and nucleotide sequences were determined by the dideoxy chain termination method [23] utilizing T7 DNA polymerase (US Biochemical Corp.) and [α - 35 S] dATP (DuPont - New England Nuclear). In order to judge if any portions of exon sequences were present, various primers used for cDNA sequencing [17, 18] were used as described [19]. Once the phage DNA yielded a sequence for a particular portion of *C2GnT* or *IGnT* exon sequence, sequencing was extended further using newly synthesized oligonucleotides based on the obtained sequence data. All of the oligonucleotides used were synthesized on an Applied Biosystems DNA synthesizer.

Results and discussion

Isolation and characterization of the human C2GnT and IGnT genes

As shown previously, eight genomic clones were initially isolated from about 1×10^6 plaques of the placental genomic DNA library. Among these, clone 20 was found to contain the *C2GnT* gene while clones 2, 3, 7, 9, 13 and 14 were found to contain various parts of the *IGnT* genes. The sequencing of those genomic clones and an additional clone containing *C2GnT* exon 1 showed the following results.

C2GnT is coded by two exons, of which the second exon encodes the whole translation product. In contrast, the complete coding sequence for *IGnT* is divided over

three exons. As shown previously, *C2GnT* and *IGnT* share three regions of extensive homology in their catalytic domains [18]. However, the high homologous region B is split between exons 1 and 2 in the *IGnT* gene while the same region is encoded entirely by exon 2 in the *C2GnT* gene [19].

Isolation and characterization of a pseudogene related to C2GnT

During these studies, we also isolated clone 6 which hybridized with *C2GnT*. The nucleotide sequence of clone 6 differs from that of *C2GnT* or *IGnT*. However, when the nucleotide sequence of clone 6 was tested for homology with known sequences, it became evident that clone 6 contains a nucleotide sequence which is highly related to *C2GnT* (Fig. 1).

The sequence corresponding to the initiation methionine of *C2GnT* lies in nucleotides 792–794. The nucleotide sequence of clone 6 can be aligned by regarding this codon as the initiation methionine. The resultant translated sequence, however, becomes out of frame starting from nucleotides 990–992 (Fig. 2). Moreover, the nucleotide sequence of clone 6 is very similar to the *C2GnT* cDNA sequence but no sequence corresponding to *C2GnT* intron 1 can be found (see Fig. 1). These results strongly indicate that this newly isolated gene is a pseudogene which is highly related to *C2GnT*.

The striking feature of the nucleotide sequence of this gene is, however, that it lacks a polyadenylation signal after nucleotide 2539. The corresponding sequence in *C2GnT* contains a polyadenylation signal and polyadenylation takes place after nucleotide 2539 (see Fig. 1). These results strongly argue against the hypothesis that this pseudogene was produced from *C2GnT* mRNA.

The pseudogene contains a 5'-upstream sequence homologous to one of the *Kpn*I repetitive sequences (*Kpn*13) and this sequence is present from nucleotide 1 to nucleotide 800 (Fig. 3). *Kpn*I repetitive sequence is the major human long interspersed repeated DNA sequence [24]. *Kpn*I repeats often have deletions and rearrangements at the 5'-end. The number of copies in the total genome is 5×10^4 to 1×10^4 for the 3'-end and 0.4×10^4 to 2.0×10^4 copies for the 5'-end. In this pseudogene, deletions took place apparently at the 3'-end of the *Kpn*I repeat. Interestingly, this *Kpn*I repeat is superimposed by an *Alu* repetitive sequence [25] from nucleotide 1 to 295 (Fig. 4). *Alu* repeat sequences represent the major human short interspersed DNA sequence and nearly 1×10^6 copies of the *Alu* repeat are present in the human genome [25]. We have analysed the 5'-sequence upstream from the exon 1 of *C2GnT*. However, there is no *Kpn*I or *Alu* sequence corresponding to the upstream sequence of the pseudogene. These results suggest that the *Alu* or *Kpn*I sequence was introduced after the ancestral gene was

```

848
ACCAAGCAAAGTCCCTGATTGGCATTGAAATGCTGAGGCAGTTGCTGT--GGAGACATTTTCTTATCCCACTAAATACCCTTGTG
      M L R Q L L - - R H F S Y P T K Y H F V
      M L R T L L R R R L F S Y P T K Y Y F M

935
GTTCTTATTTTCCCTAGTCACC---TCCGTTTAAGGATTTCATCAAAGTCCAAATCTGTAAGCGTCACATATGTGGAGCTTGTGGGA
V L I F S L V T - S V L R I H Q K S K S V S V T Y V E L V C
V L V L S L I T F S V L R I H Q K P E F V S V R H L E L A G

1025
GAGAATCCTAGTAGTCAATTAATGACCAAAAGTTTACGGGGGGGATGTAGATGAAATCCAAAAGGTAACACTTGAGATGCTAACAGT
E N P S S H I N C T K V L / G D V D E I Q K V K L E M L T V
E N P S S D I N C T K V L Q G D V N E I Q K V K L E I L T V

1114
GAAATTTACTAAGTCCCTTGGTGTATACCTGACGGCTTTATAACATGACCAGTTA-TGTACTTCTTCATGACTTCTTCATCAAGAG
K F T K C P W C I P D G F I N M T S / C T S S C T S F I K R
K F K K R P R W T P D D Y I N M T S - - - - D C S S F I K R

1204
ATGTAGATATATTTAGAACCCCTTAAGAAGGAAGAGGTGAGGTTTCCAATAGCATATCTATACTGGTTCATTAATAAACTGAAACGGT
C R Y I V E P L K K E E V R F P I A Y S I L V H Y K T E T L
R K Y I V E P L S K E E A E F P I A Y S I V V H H K I E M L

1294
TGACAGGCTCCAGAGGCCATCTATATGCCTCAGAATTTCTATGCAATTCATGTGGACAAAAAAATCAGCAGATTCCCTTTTGTAGCTGCA
D R L Q R A I Y M P Q N F Y C I H V D / K S A D S F L A A
D R L L R A I Y M P Q N F Y C V H V D T K S E D S Y L A A

1382
GTGATGGGCATGGGTCA--TTTCAGTAACATCTTTGTGGCCTGTCAGTTGGAGAGTCTGGTTTATGCCTTGTGGAGTCGGCTTCTGCCT
V M G I G S / F S N I F V A C Q L E S L V Y A L W S R V L A
V M G I A S C F S N V F V A S R L E S V V Y A S W S R V Q A

1472
GACCTCAACTGCATGAGGGACCTCTGCACAGTGCAGACTGGAAGTACTTAATACATGTTTGTAGTATGGATTTCCTATTAATAAAC
D L N C M R D L C T V S A D W K Y L I H V C S M D F P I K T
D L N C M K D L Y A M S A N W K Y L I N L C G M D F P I K T

1561
AACCTA-AAATGTTAGCAAGCTCAAGTTGTTAATGGGTGAAGCAGTCTCAAAGCCAAGGATGCCATCCAATAAAGAAAGAAAGGTTG
N L / I V R K L K L L M G E D S L K A K R M P S N K E E R W
N L E I V R K L K L L M G E N N L E T E R M P S H K E E R W

1651
AAAAAGTGGTATGCAGATATTAATGAAAGCTGACACATGTGGGGACTGTCAAAGGGCATCCTCCGCTGGAAGCACCCTATTTTTCAGGC
K K W Y A D I N G K L T H V G T V K G H P P L E A P I F S G
K K R Y E V V N G K L T N T G T V K M L P P L E T P L F S G

1741
AGTGCCTATTTGTGGTCACTAGCGAGTATGTGGGCCATGTGCTAGAGGATGAAAAACCCAAAAGTTTATGGACTGGGTGCGAGGCACA
S A Y F V V S R E Y V G H V L E D E K T Q K F M E W V R G T
S A Y F V V S R E Y V G Y V L Q N E K I Q K L M E W A Q D T

1831
GACAGCCAGATAAGTATCTTAGCCATCATCCGAAGGATCGCTGAAGTCCCTGGCTCATTCGCCCTTAAGCCATAAGTACAAGTTGCTC
D S P D K Y L * A I I R R I A E V P G S F A L S H K Y K L S
Y S P D E Y L W A T I Q R I P E V P G S L P A S H K Y D L S

1921
GGAATGCATGCCGTTGCTAGGTTTGTCAAGTGGCAGTACTGAGGATGACGTTTTCGAAGGATGCTCCCTACCCACCTGCAGTGGGGTC
G M H A V A R F V K W Q Y S E D D V F K D A P Y P P C S G V
D M Q A V A R F V K W Q Y F E G D V S K G A P Y P P C D G V

2011
TCCATGCACTCAGCATGCAATTTTCGGAGCCAGCAGCTTGAAGTGGATGCTGTGTAACACCTATGGGTGCAAGCTTATACCTTTGACATG
S M H S A C I F G A S S L N W M L C K H L W V Q A Y T F D M
H V R S V C I F G A G D L N W M L R K H H - L F A N K F D V

2101
GATGTTGACCTCCTTGCCACCTAGTGTGTTGGATGAGCATCTGAGGCATAAAGCTTTGGAGACTTTAAAACACTGACCATTATTAGCAATT
D V D L L A T * C L D E H L R H K A L E T L K H *
D V D L F A I Q C L D E H L R H K A L E T L K H *

```

Figure 2. Comparison of translated amino acid sequences of the pseudogene and *C2GnT*. Each set of three lines (from top to bottom) shows the nucleotide sequence of the pseudogene, its deduced amino acid sequence and *C2GnT* amino acid sequence. In order to maximize the homology between two predicted amino acid sequences, gaps are allowed. Stop codons and frame shifts are denoted by asterisks and slashes, respectively, in the pseudogene amino acid sequence.

pseudogene	-----	
Alu seq.	GGCCGGGCGCGGTGGCTCACGCCTGTAATCCCAAGCACTTTGGGAGGCCGA	50
pseudogene	-----	
Alu seq.	GGCGGGCGGATCACCTGAGGTCAGGAGTTTCGAGACCAGCCTGGCCAACAT	100
pseudogene	-----CTCGAGGRRCCCTTA	14
Alu seq.	GGTGAAACCCCGTCTCTACTAAAAATACAAAATTAGCCG-GGCGTGGT	149
pseudogene	TCTCACCCATATACTAAAATCAACACAAAATGGATTAACTGAAAACT	64
Alu seq.	GCSCCCCCCTGRTAGTCCCACTACTCCGGG-GGCTTAGCCCGGAAAACT	198
pseudogene	TCTATCGAAACATAGGGAAAAAGCT-TA-TGACATTAG-TCTGGSCAA	111
Alu seq.	GCT-T--GAACCCGGAGGEGGAEGTTCAGTGAAGCCGAGATCGCCGCAC	245
pseudogene	TG-ATTTTTTTTGGATATGAC-CCCAAGCCCGCAAGCAAAAATRAAA	159
Alu seq.	TGDAITCCAGCCTG-GGCGACAGAGCCAGACTCCGTC--TCAAAAAAAA	292
pseudogene	AAAAGACAAACAGAATTACATCAAACATAAAAGCTTCTGCACAGCAAG	209
Alu seq.	AAA-----AAAA-----AA--	301
pseudogene	GAAACAATGAACACAGTTAAGAGACAACCTACCAGAATGGAAGAAAATAT	259
Alu seq.	-----	301

Figure 4. Comparison of the nucleotide sequences of the pseudogene and *Alu* repetitive sequence. *Alu* repetitive sequence has homology with the upstream sequence of the pseudogene.

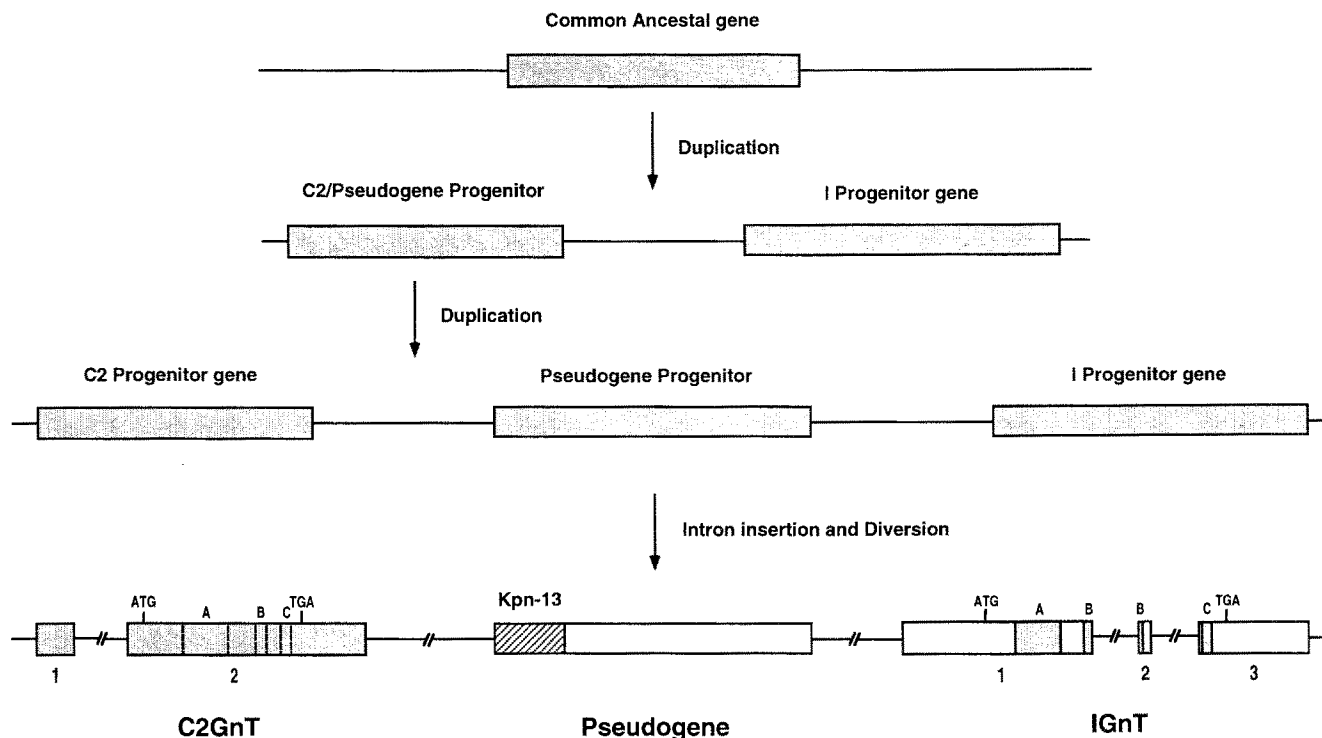


Figure 5. Evolutionary pathway of the *C2GnT* gene, the *IGnT* gene and the pseudogene. A common ancestral gene was duplicated and the resultant gene directly evolved into *IGnT* gene after intron insertion and diversion. The other gene was further duplicated to produce *C2GnT* progenitor gene and the pseudogene. The *C2GnT* progenitor gene evolved into *C2GnT* gene by intron insertion and diversion.

different sialyltransferases. One of them, sialyl motif L, corresponds to residues 178–225 in rat Gal β 1 \rightarrow 4GlcNAc α -2,6-sialyltransferase and was recently found to be the binding site for CMP-NeuNAc [27]. The residues 178 to 225 are, however, split between exon 2 and exon 3 [28, 29], supporting the conclusion that those homologous regions were not brought together by exon shuffling.

The studies on the genomic organization of glycosyltransferases revealed that there are two different types of genomic organization. One is represented by the *C2GnT* gene, in which the entire coding region is coded by one exon. These include *N*-acetylglucosaminyltransferase I [30], and fucosyltransferases III–VI [31]. The other is represented by the *IGnT* gene, in which the coding regions are split by several exons. These include β -1,4-galactosyltransferase [32], α -2,6-sialyltransferase [28, 29] and α -1,3-galactosyltransferase [33]. The present study strongly suggests that these seemingly different gene organizations among different glycosyltransferases were most likely produced from the common ancestral genes by gene duplication, diversion and intron insertion. Further studies are needed on the *C2GnT* and *IGnT* gene family to test this hypothesis by studying the chromosome localization of the pseudogene and the genomic organization of *C2GnT* and *IGnT* in the lower animal kingdom.

Acknowledgements

We thank Dr Kiyohiko Angata for useful discussion and Ms Bobbi Laubhan for secretarial assistance. This work was supported by Grants CA33000 and CA33895 from the National Cancer Institute.

References

1. Fukuda M (ed). (1992) *Cell Surface Carbohydrates and Cell Development* pp. 1–329. Boca Raton, FL: CRC Press.
2. Feizi T (1985) *Nature* **314**: 53–57.
3. Fukuda M (1985) *Biochim Biophys Acta* **780**: 119–50.
4. Fukuda M, Carlsson SR, Klock JC, Dell A (1986) *J Biol Chem* **261**: 12796–806.
5. Piller F, Piller V, Fox RI, Fukuda M (1988) *J Biol Chem* **263**: 15146–50.
6. Maemura K, Fukuda M (1992) *J Biol Chem* **267**: 24379–86.
7. Bierhuizen MFA, Maemura K, Fukuda M (1994) *J Biol Chem* **269**: 4471–79.
8. Ohmori K, Takada A, Yoneda T, Burna K, Harashima K, Tsuyuoka K, Hasegawa A, Kannagi R (1993) *Blood* **81**: 101–11.
9. Baum LG, Pang M, Perillo NL, Terry W, Uittenbogaart C, Fukuda M, Seilhammer JJ (1995) *J Exp Med* **181**: 877–87.
10. Piller F, Le Desit F, Weinberg KI, Parkman R, Fukuda M (1991) *J Exp Med* **173**: 1501–10.
11. Higgins EA, Siminovitch KA, Zhuang D, Brockhausen I, Dennis JW (1991) *J Biol Chem* **266**: 6280–90.
12. Saitoh O, Piller F, Fox RI, Fukuda M (1991) *Blood* **77**: 1491–99.
13. Brockhausen I, Kuhns W, Schacter H, Matta KL, Sutherland DR, Baker MA (1991) *Cancer Res* **51**: 1257–63.
14. Yousefi S, Higgins E, Daoling Z, Pollex-Kruger A, Hindsgaul O, Dennis JW (1991) *J Biol Chem* **266**: 1772–82.
15. Fukuda M, Fukuda MN, Hakomori S (1979) *J Biol Chem* **254**: 3700–3.
16. Fukuda M, Fukuda MN (1984) In *The Biology of Glycoproteins* (RJ Ivatt ed.) pp. 183–234. New York, Plenum Press.
17. Bierhuizen MFA, Fukuda M (1992) *Proc Natl Acad Sci USA* **89**: 9326–30.
18. Bierhuizen MFA, Mattei M-G, Fukuda M (1993) *Genes Dev* **7**: 468–78.
19. Bierhuizen MFA, Maemura K, Kudo S, Fukuda M (1995) *Glycobiology* **5**: 417–25.
20. Feinberg A, Vogelstein B (1983) *Anal Biochem* **132**: 6–13.
21. Saiki RK, Gelfand DH, Stoffel S, Scharf SJ, Higuchi R, Horn GT, Mullis KB, Erlich HA (1988) *Science* **239**: 487–91.
22. Sambrook J, Fritsch EF, Maniatis T (1989) *Molecular Cloning: A Laboratory Manual*. Cold Spring Harbor, NY: Cold Spring Harbor Laboratory.
23. Sanger F, Nicklen S, Coulson AF (1977) *Proc Natl Acad Sci USA* **74**: 5463–67.
24. Sun L, Paulson KE, Schmid CW, Kadyk L, Leinwand L (1984) *Nucleic Acids Res* **12**: 2669–90.
25. Schmid CW, Jelinek WR (1982) *Science* **216**: 1065–70.
26. Kudo S, Fukuda M (1989) *Proc Natl Acad Sci USA* **86**: 4619–23.
27. Datta AK, Paulson JC (1995) *J Biol Chem* **270**: 1497–1500.
28. Svensson EC, Soreghan B, Paulson JC (1990) *J Biol Chem* **265**: 20863–68.
29. Wang X, O'Hanlon TP, Young RF, Lau JTY (1990) *Glycobiology* **1**: 25–31.
30. Hull E, Sarkar M, Spruijt MPN, Hoppener JWM, Dunn R, Schacter H (1991) *Biochem Biophys Res Commun* **176**: 608–15.
31. Weston BW, Smith PL, Kelly RJ, Lowe JB (1992) *J Biol Chem* **267**: 24575–84.
32. Hollis GF, Douglas JG, Shaper NL, Shaper JH, Stafford-Hollis JM, Evans RJ, Kirsch IR (1989) *Biochem Biophys Res Commun* **162**: 1069–75.
33. Yamamoto F, McNeill PD, Hakomori S (1991) *Biochem Biophys Res Commun* **175**: 986–94.