

# Glycoproteogenomics characterizes the CD44 splicing code associate with bladder cancer invasion

## -Supplementary figures and tables-

Cristiana Gaitero<sup>1,2,3,4,5</sup>, Janine Soares<sup>1,2,3,4,6</sup>, Marta Relvas-Santos<sup>1,2,3,4,7,8,9</sup>, Andreia Peixoto<sup>1,2,3,7,8</sup>, Dylan Ferreira<sup>1,2,3,4,5,7,8</sup>, Paula Paulo<sup>2,3,10</sup>, Andreia Brandão<sup>2,3,10</sup>, Elisabete Fernandes<sup>1,2,3,11</sup>, Rita Azevedo<sup>12</sup>, Carlos Palmeira<sup>1,2,3,11,13</sup>, Rui Freitas<sup>1,2,3,4,7,8</sup>, Andreia Miranda<sup>1,2,3,14</sup>, Hugo Osório<sup>7,14,15</sup>, Jesús Prieto<sup>5</sup>, Luís Lima<sup>1,2,3</sup>, André M. N. Silva<sup>9,16</sup>, Lúcio Lara Santos<sup>1,2,3,4,11,16,17</sup>, José Alexandre Ferreira<sup>1,2,3,4,16,a</sup>

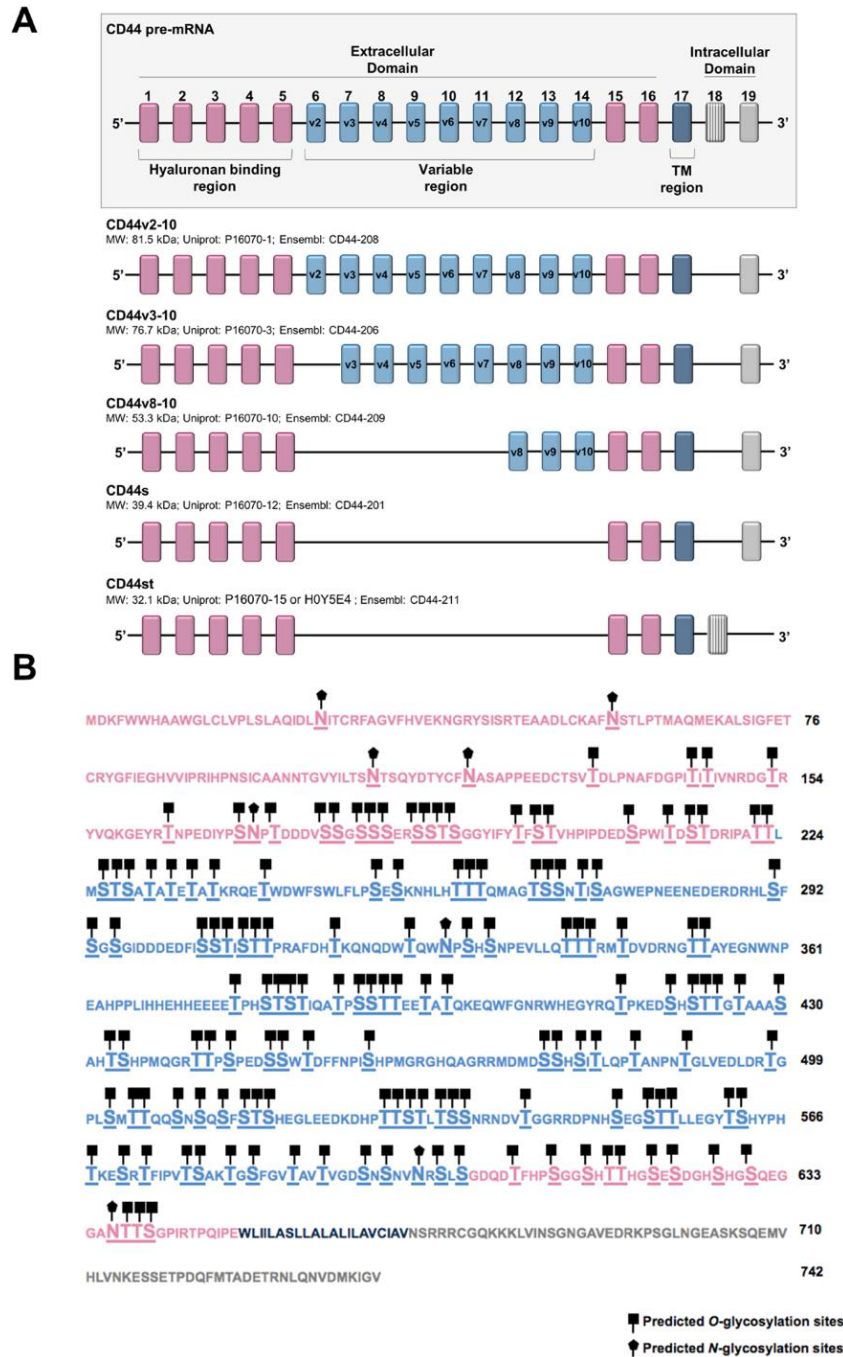
<sup>1</sup>Experimental Pathology and Therapeutics Group, IPO Porto Research Center (CI-IPOP), Portuguese Oncology Institute (IPO Porto), 4200-072 Porto, Portugal; <sup>2</sup>RISE@CI-IPOP (Health Research Network), Portuguese Oncology Institute of Porto (IPO Porto), 4200-072 Porto, Portugal; <sup>3</sup>Porto Comprehensive Cancer Center (P.ccc), 4200-072 Porto, Portugal; <sup>4</sup>Institute of Biomedical Sciences Abel Salazar (ICBAS), University of Porto, 4050-013 Porto, Portugal; <sup>5</sup>Center for Applied Medical Research (Centro de Investigación Médica Aplicada, CIMA), University of Navarra, 31008 Pamplona, Navarra, Spain; <sup>6</sup>REQUIMTE-LAQV, Department of Chemistry, University of Aveiro, 3810-193 Aveiro, Portugal; <sup>7</sup>Institute for Research and Innovation in Health (i3S), University of Porto, 4200-135 Porto, Portugal; <sup>8</sup>Institute for Biomedical Engineering (INEB), University of Porto, 4200-135 Porto, Portugal; <sup>9</sup>REQUIMTE-LAQV, Department of Chemistry and Biochemistry, Faculty of Sciences of the University of Porto, 4169-007 Porto, Portugal; <sup>10</sup>Cancer Genetics Group, IPO Porto Research Center (CI-IPOP), Portuguese Oncology Institute (IPO Porto), 4200-072 Porto, Portugal; <sup>11</sup>FP-I3ID, University Fernando Pessoa, 4249-004 Porto, Portugal; <sup>12</sup>Laboratoire d'Etude du Métabolisme des Médicaments (LEMM), CEA, INRA, Université Paris Saclay, F-91191, Gif-sur-Yvette cedex, France; <sup>13</sup>Immunology Department, Portuguese Oncology Institute (IPO Porto), 4200-072 Porto, Portugal; <sup>14</sup>Faculty of Medicine of the University of Porto, 4200-319 Porto, Portugal; <sup>15</sup>Ipatimup—Institute of Molecular Pathology and Immunology of the University of Porto, University of Porto, 4200-135 Porto, Portugal; <sup>16</sup>GlycoMatters Biotech, 4500-162 Espinho, Portugal; <sup>17</sup>Department of Surgical Oncology, Portuguese Oncology Institute (IPO Porto), 4200-072 Porto, Portugal

### **<sup>a</sup>Corresponding author:**

José Alexandre Ferreira (jose.a.ferreira@ipoporto.min-saude.pt)

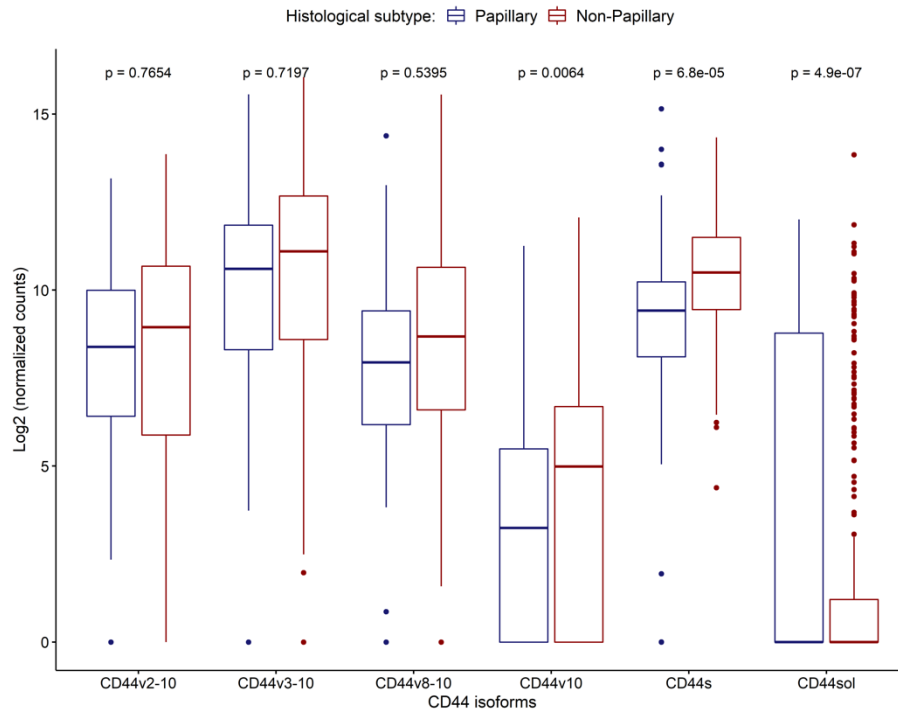
ORCID ID: <https://orcid.org/0000-0002-0097-6148>

Experimental Pathology and Therapeutics Group, Research Centre, Portuguese Oncology Institute of Porto, R. Dr. António Bernardino de Almeida 62, 4200-072 Porto, Portugal; Tel. +351 225084000 (ext. 5111).



**Figure S1. A) Schematic representation of human CD44 pre-mRNA and experimentally confirmed isoforms resulting from alternative splicing and extensively explored in this study.** Exons generally regarded as conserved are represented in pink (constitutive exons e1-e5 and e15-e16), variable exons are represented in blue (exons v6-v14), the transmembrane region is represented in dark blue (exon e17), and the intracellular tail is represented in grey (exons e18-e19). Exon 18, filled with gray stripes, contains an early 3'UTR and is only present in the CD44st isoform. The estimated molecular weight for each variant (without posttranslational modifications), Uniprot and Ensemble codifications are also highlighted. **B) Amino acid sequence evidencing potential O- and N-glycosites.** O-glycosylation sites were predicted using the NetOGlyc 4.0 server (<http://www.cbs.dtu.dk/services/NetOGlyc/>) and N-glycosylation sites with the NetNGlyc 1.0 server (<http://www.cbs.dtu.dk/services/NetNGlyc/>). The different domains

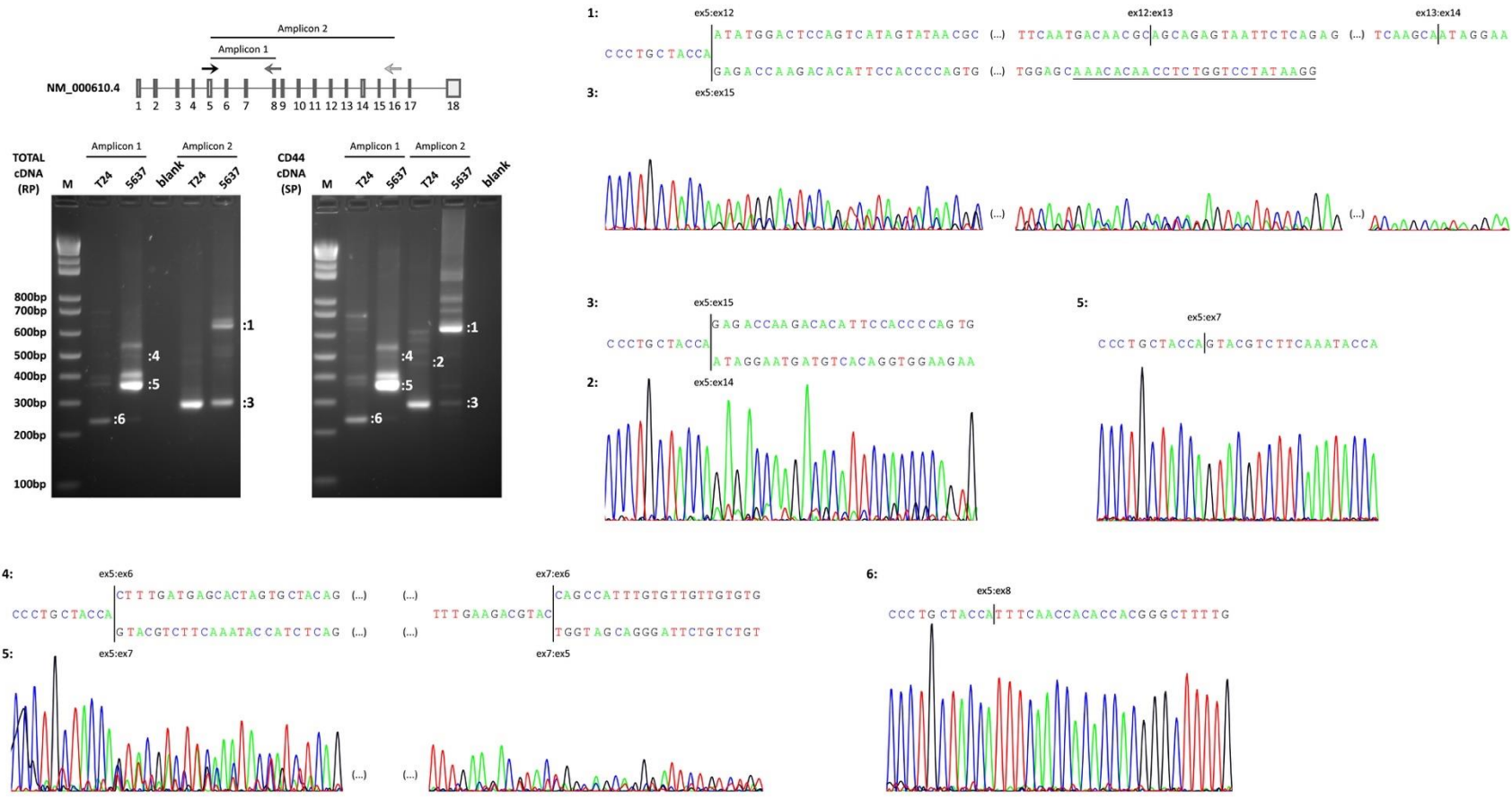
of the protein are highlighted using the same color code as in panel A and evidences the high density of *O*-glycosites in comparison to *N*-glycosites, which is particularly evident in the variable region.



**Figure S2. CD44 variants expression according to the histological subtypes of TCGA muscle invasive bladder tumors.** Wilcoxon test revealed a significant overexpression of CD44s in more aggressive non-papillary lesions.



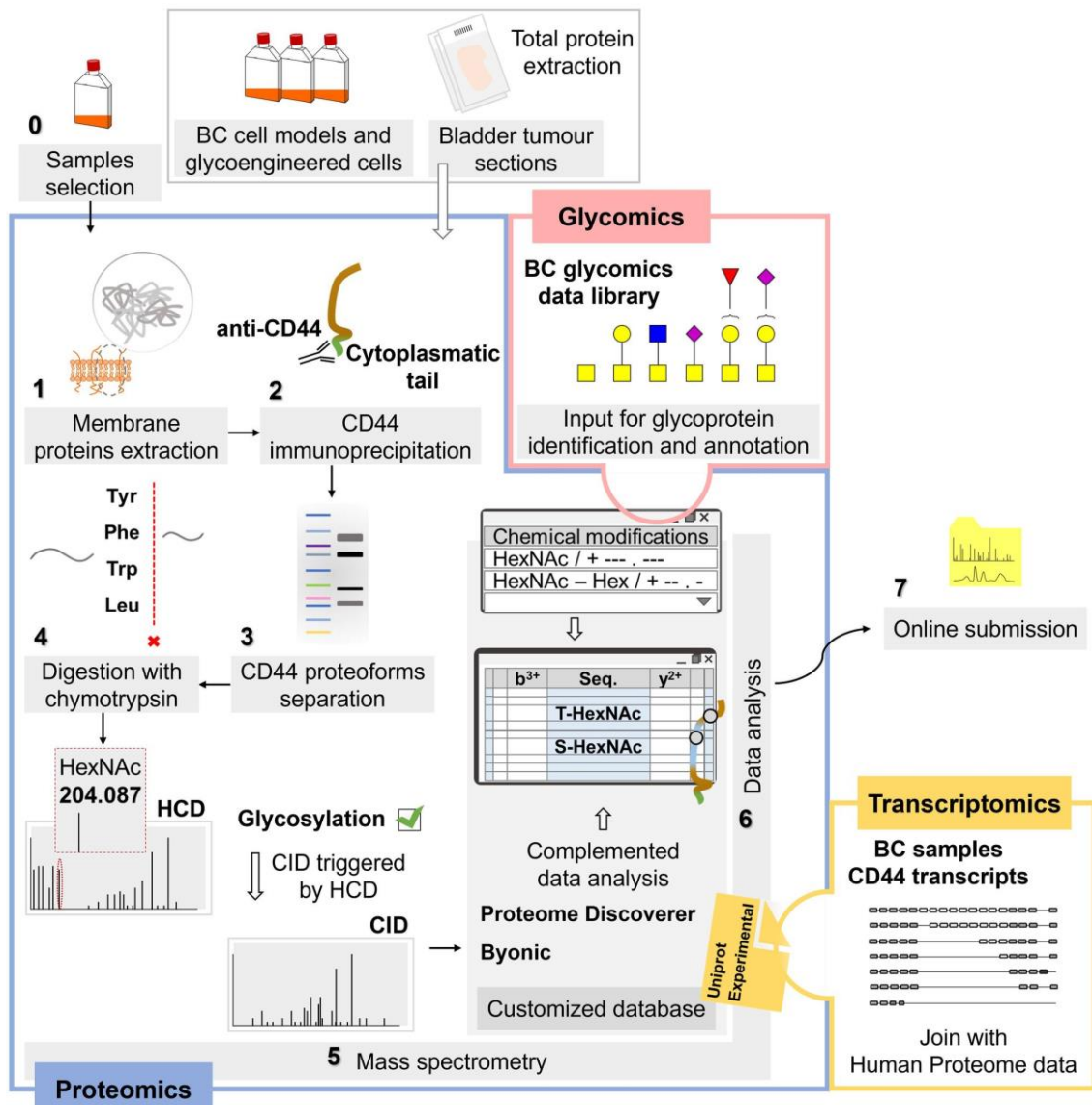
**Figure S3. CD44 splice variants expression in healthy human tissues.** Transcriptomics analysis of healthy tissues for CD44 isoforms were obtained and adapted from the GTEx Portal on 03/26/21 and/or dbGaP accession number phs000424.v8.p2. The Figure shows the complex mosaicism of CD44 mRNA in healthy tissues, including the co-expression of multiple isoforms in the same tissue. It highlights that CD44s is not expressed in healthy organs. Moreover, CD44st that differs from CD44s in terms of cytoplasmic tail but presents a similar extracellular domain was detected in low abundance in a restricted number of cells/organs, mostly from the skin, fibroblasts, subcutaneous adipocytes, the female reproductive system, the esophagus and breast. On the other hand, CD44s related structures lacking exon 15 (CD44s-exon 15) was present in high abundance in adipose tissues, fibroblasts, lungs, blood cells, across the female reproductive system and the gastrointestinal tract.



**Figure S4. Reverse Transcriptase-Polymerase Chain Reaction and Sanger Sequencing for the *CD44* gene using random primer (RP) and e16 specific primer amplification.** Two amplicons have been used: amplicon 1: e5-v4; amplicon 2: e5-e16. According to these experiments, 5637 were enriched longer isoforms in relation to T24 cells. However, we could not identify the precise nature of the long isoforms by sanger sequencing, most likely because of mixture of different isoforms. Supporting the distinct splicing codes presented by these cell lines, the following differences were observed: **e5-v8-v9-v10-(...)-e16 (band 1):** 5637 (main band)>T24; **e5-v10-(...)-e16 (band 2):** only in T24; **e5-e15-(...)-e16 (band 3):** T24 (main band; most likely corresponding to CD44s and/or st)>5637; **e5-v2-v3-(...)-e16 (band 4):** only in 5637; **e5-v3-v4-(...)-e16 (band 5):** 5637>T24; **e5-v4-(...)-e16 (band 6):** only in T24.

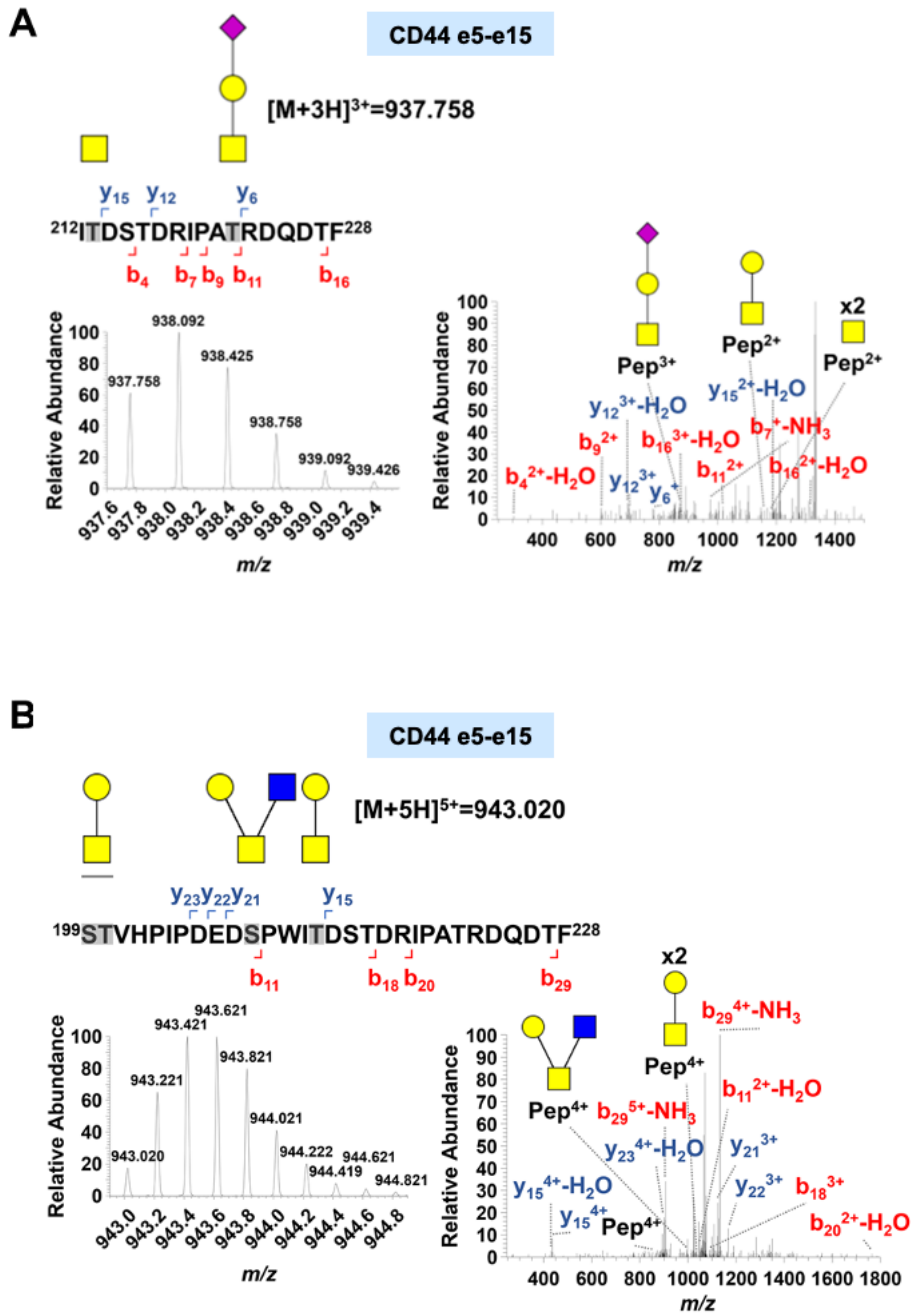




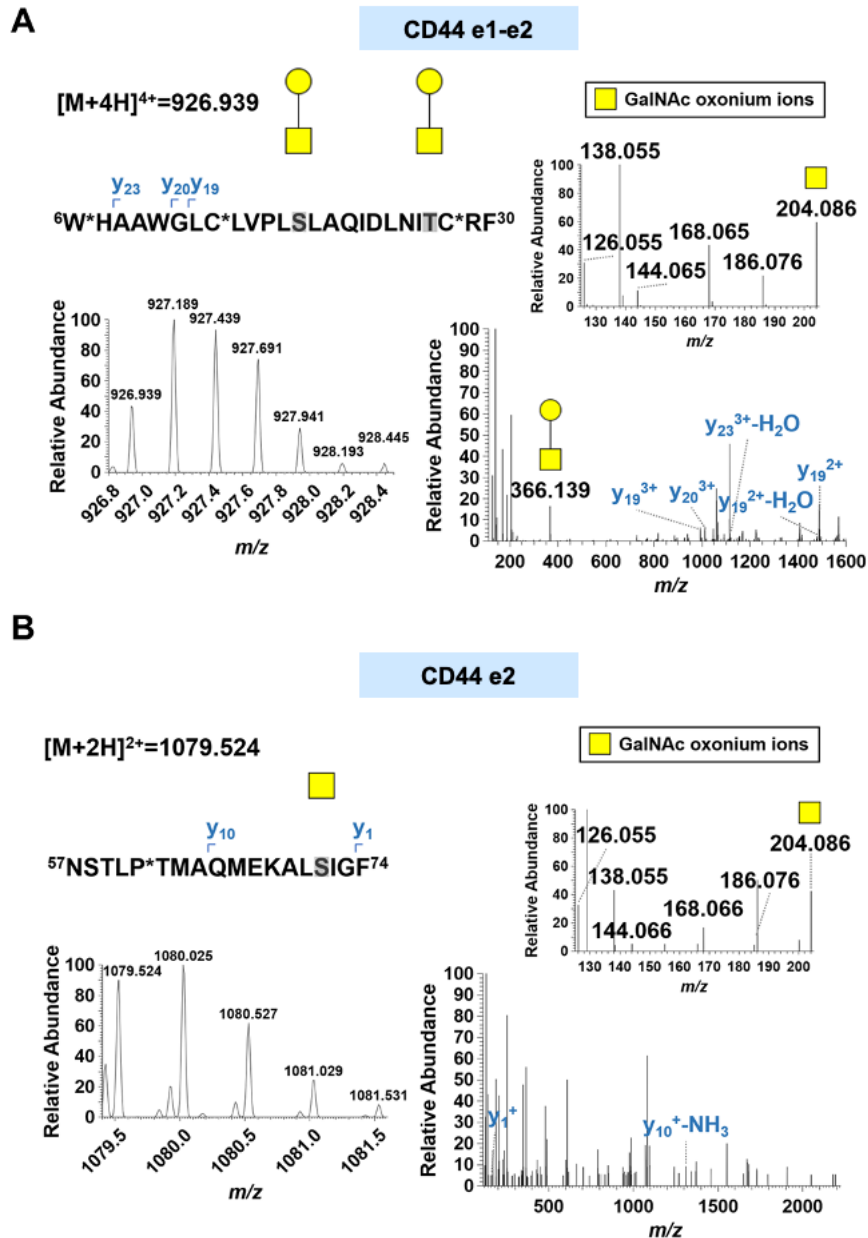


**Figure S5. Glycoproteogenomics workflow for identification of CD44 glycoproteoforms.** Briefly, bladder cancer cell models, glycoengineered cells and tumor tissue sections were used as starting material. **1.** Membrane glycoproteins from cell lines are isolated by differential ultracentrifugation. **2.** CD44 is immunoprecipitated using a monoclonal antibody targeting the cytoplasmatic tail domain expressed by most CD44 isoforms. **3.** CD44-IP enriched extracts were then digested with sialidase, separated by gradient SDS-PAGE electrophoresis. **4.** Bands are excised from the gels, proteins are reduced, alkylated, and digested with chymotrypsin *in gel*. **5.** Chymotrypsin digests are analyzed by nanoLC-MS/MS (HCD-MS2 and CID-MS2 with a dependent acquisition based on HCD detection of HexNAc oxonium ion at  $m/z$  204.087). **6.** Glycoproteoforms identification using a glycoproteogenomics approach. RNAseq was used to identify CD44 isoforms expressed by 5637 and T24 cells. Predicted protein sequences were included in a database together with the human proteome from Uniprot database and used for protein identification in the Proteome Discoverer bioinformatics platform. Glycomics analysis was performed in parallel to determine the glycan structures that may be present as variable CD44 post-translational modification. This information was included upon protein identification and glycopeptides annotation by MS/MS.

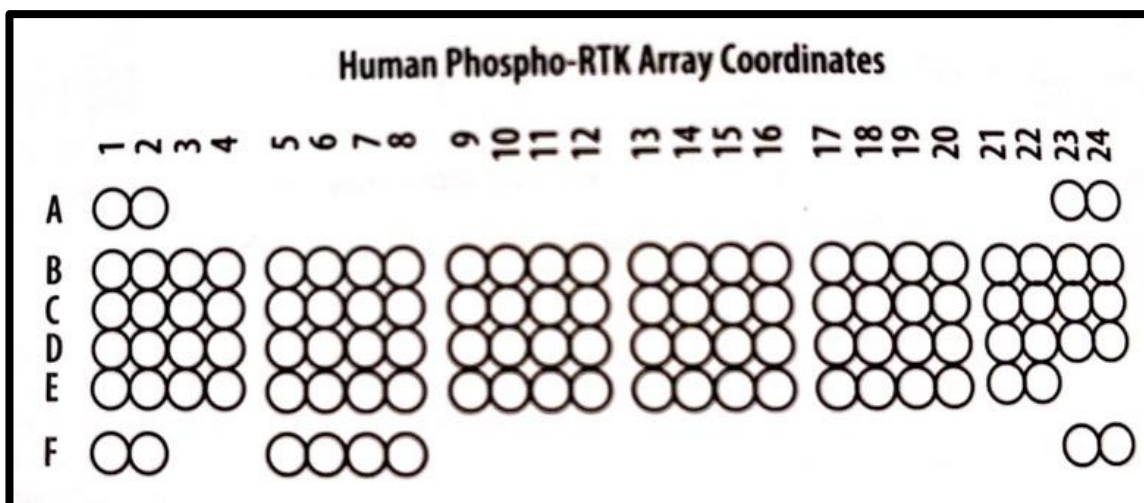




**Figure S6. CID spectra for CD44s glycopeptides carrying *O*-glycans identified in T24 cells. A) CID spectrum for CD44s-Tn/ST and B) CD44s-T/core 2 glycopeptides. CD44s glycopeptides carrying these *O*-glycans were identified by CID based on an HCD data dependent acquisition. Briefly, species presenting HCD product ion spectrum with an HexNAc fragment at  $m/z$  204.087 were elected for CID fragmentation. CID spectra confirmed the glycan structures and presented several *y*- and *b*-series peptide fragments that helped supporting glycosites assignment (highlighted in grey in the peptide sequence). Collectively, these spectra confirm the presence of the illustrated *O*-glycans in CD44s from T24 cells and the heterogeneous nature of glycosylation at the peptide level.**



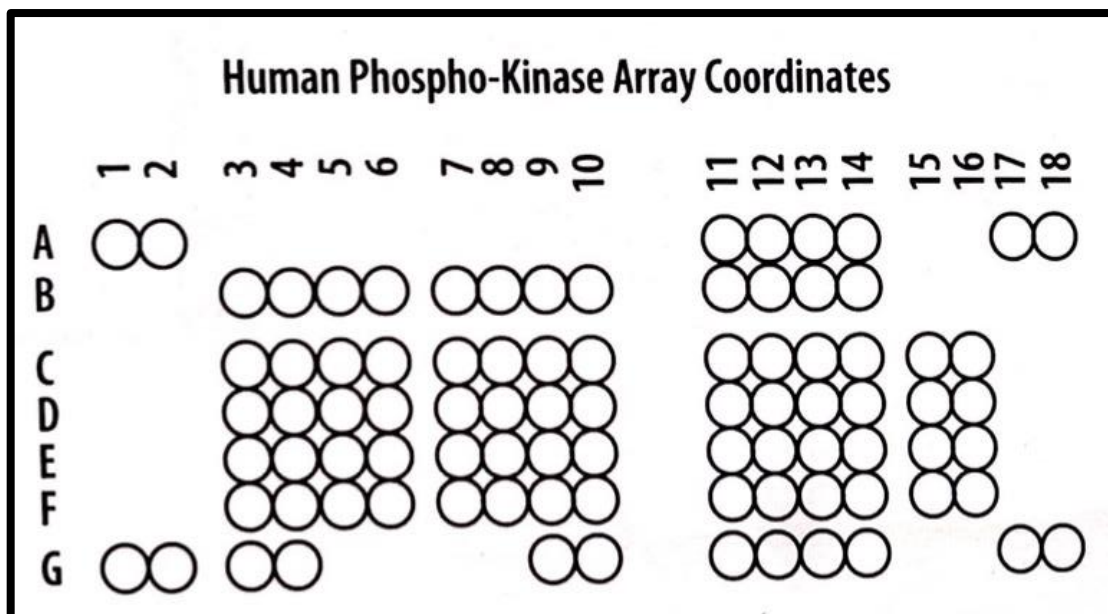
**Figure S7. A) CD44-T and B) Tn glycopeptides belonging to the conserved extracellular domain in different isoforms identified in muscle invasive bladder tumors by nanoLC-HCD-MS/MS. Both MS/MS spectra present the HexNAc oxonium ions, namely the ion at  $m/z$  204.086. The MS/MS spectrum of the CD44-T glycopeptide also presents an evident ion at  $m/z$  366.139, corresponding to the oxonium ion for the T antigen.  $y$ -series ions resulting from peptide backbone fragmentation were also highlighted. The sequence of the glycopeptide presents the glycosites identified in grey. The symbol \* corresponds to amino acid modifications: W – hydroxykynurenine; C – carbamidomethyl; P – glutamic semialdehyde.**



Coordinate	Receptor Family	RTK/Control
A1, A2	Reference Spots	-
B1, B2	Reference Spots	-
B3, B4	EGF R	EGF R
B5, B6	EGF R	ErbB2
B7, B8	EGF R	ErbB4
B9, B10	FGF R	FGF R1
B11, B12	FGF R	FGF R2 $\alpha$
B13, B14	FGF R	FGF R3
B15, B16	FGF R	FGF R4
B17, B18	Insulin R	Insulin R
B19, B20	Insulin R	IGF-I R
B21, B22	Axl	Axl
B23, B24	Axl	Dtk
C1, C2	Axl	Mer
C3, C4	HGF R	HGF R
C5, C6	HGF R	MSP R
C7, C8	PDGF R	PDGF R $\alpha$
C9, C10	PDGF R	PDGF R $\beta$
C11, C12	PDGF R	SCF R
C13, C14	PDGF R	Flt-3
C15, C16	PDGF R	M-CSF R
C17, C18	RET	c-Ret
C19, C20	ROR	ROR1
C21, C22	ROR	ROR2
C23, C24	Tie	Tie-1
D1, D2	Tie	Tie-2
D3, D4	NGF R	TrkA
D5, D6	NGF R	TrkB
D7, D8	NGF R	TrkC
D9, D10	VEGF R	VEGF R1
D11, D12	VEGF R	VEGF R2
D13, D14	VEGF R	VEGF R3
D15, D16	MuSK	MuSK

<b>D17, D18</b>	Eph R	EphA1
<b>D19, D20</b>	Eph R	EphA2
<b>D21, D22</b>	Eph R	EphA3
<b>D23, D24</b>	Eph R	EphA4
<b>E1, E2</b>	Eph R	EphA6
<b>E3, E4</b>	Eph R	EphA7
<b>E5, E6</b>	Eph R	EphB1
<b>E7, E8</b>	Eph R	EphB2
<b>E9, E10</b>	Eph R	EphB4
<b>E11, E12</b>	Eph R	EphB6
<b>E13, E14</b>	Insulin R	ALK
<b>E15, E16</b>	-	DDR1
<b>E17, E18</b>	-	DDR2
<b>E19, E20</b>	Eph R	EphA5
<b>E21, E22</b>	Eph R	EphA10
<b>F1, F2</b>	Reference spots	-
<b>F5, F6</b>	Eph R	Ephb3
<b>F7, F8</b>	-	RYK
<b>F23, F24</b>	Control (-)	PBS

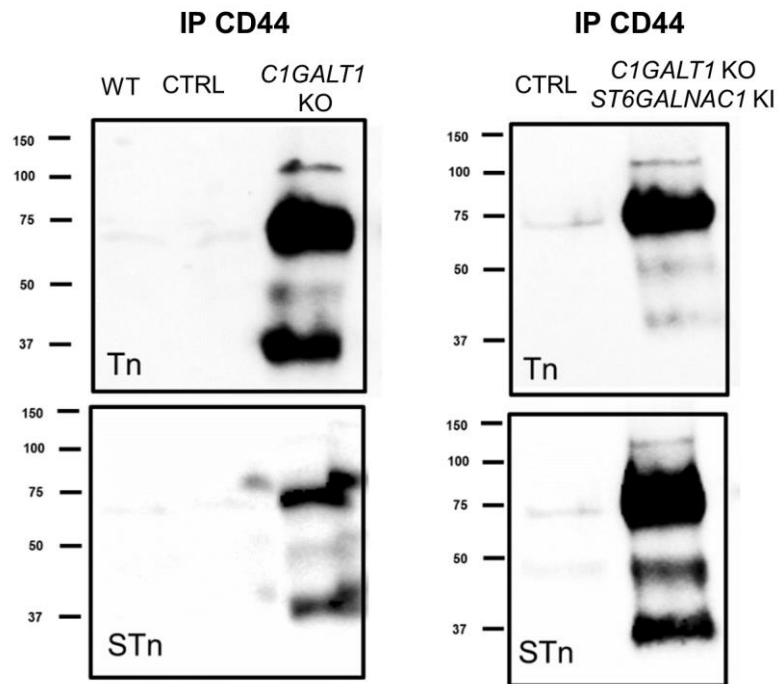
**Figure S8. The human phospho-receptor tyrosine kinase (Phospho-RTK) array coordinates.** The phospho-array presents capture and control antibodies spotted in duplicates for the 49 different phosphorylated human RTKs, which are represented by alpha-numeric coordinates.



Coordinate	Target/Control	Phosphorylation Site
A1, A2	Reference Spot	-
A11,A12	Akt 1/2/3	T308
A13, A14	Akt 1/2/3	S473
A17, A18	Reference Spot	-
B3, B4	CREB	S133
B5, B6	EGF R	Y1086
B7, B8	eNOS	S1177
B9, B10	ERK 1/2	T202/Y204, T185/Y187
B11, B12	Chk-2	T68
B13, B14	c-Jun	S63
C3, C4	Fgr	Y412
C5, C6	GSK-3 $\alpha/\beta$	S21/S9
C7, C8	GSK-3 $\beta$	S9
C9, C10	HSP27	S78/S82
C11, C12	p53	S15
C13, C14	p53	S46
C15, C16	p53	S392
D3, D4	JNK 1/2/3	T183/Y185, T221/Y223
D5, D6	Lck	Y394
D7, D8	Lyn	Y397
D9, D10	MSK 1/2	S376/S360
D11, D12	p70 S6 Kinase	T389
D13, D14	p70 S6 Kinase	T421/S424
D15, D16	PRAS40	T246
E3, E4	p38 $\alpha$	T180/Y182
E5, E6	PDGF R $\beta$	Y751
E7, E8	PLC- $\gamma$ 1	Y783
E9, E10	Src	Y419
E11, E12	PYK2	Y402
E13, E14	RSK 1/2	S221/S227

<b>E15, E16</b>	RSK 1/2/3	S380/S386/S377
<b>F3, F4</b>	STAT2	Y689
<b>F5, F6</b>	STAT5a/b	Y694/Y699
<b>F7, F8</b>	WNK1	T60
<b>F9, F10</b>	Yes	Y426
<b>F11, F12</b>	STAT1	Y701
<b>F13, F14</b>	STAT3	Y705
<b>F15, F16</b>	STAT3	S727
<b>G1, G2</b>	Reference Spot	-
<b>G3, G4</b>	$\beta$ -Catenin	-
<b>G9, G10</b>	PBS (Control -)	-
<b>G11, G12</b>	STAT6	Y641
<b>G13, G14</b>	HSP60	-
<b>G17, G18</b>	PBS (Control -)	-

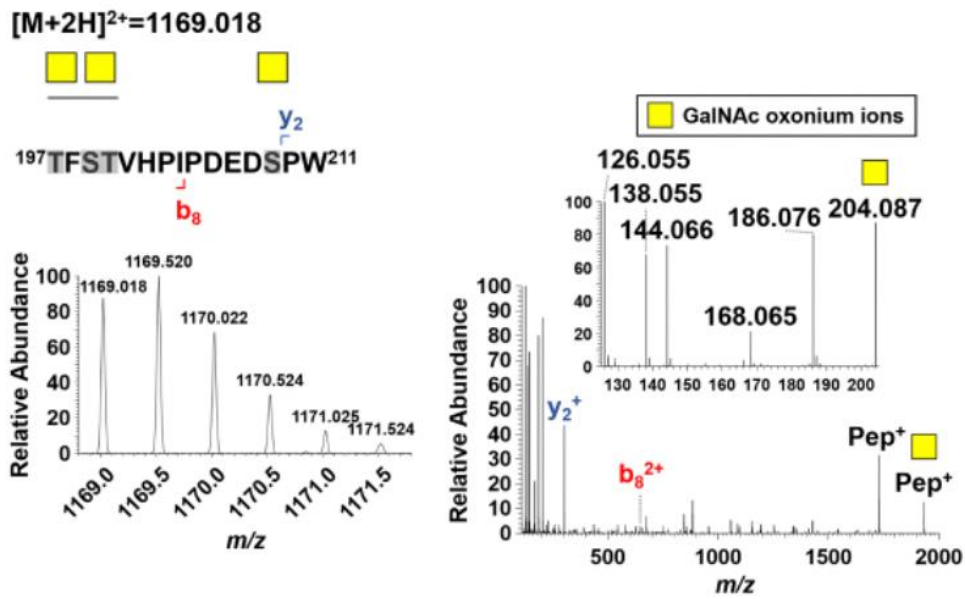
**Figure S9. The human phospho-kinase array coordinates.** The phospho-array presents capture and control antibodies spotted in duplicates for the detection of the relative levels of phosphorylation of 37 kinase phosphorylation sites and 2 related total proteins represented by alpha-numeric coordinates.



**Figure S10. Western Blots show co-expression of CD44 and short-chain *O*-glycans, namely STn and Tn, in glycoengineered T24 cells.** CD44 was immunoprecipitated from membrane protein extracts of T24 *C1GALT1* KO, T24 *C1GALT1* KO/*ST6GALNAC1* KI and corresponding controls and then blotted for STn, Tn, and CD44. A demarked expression of CD44-STn and CD44-Tn was detected for both cell model mainly in bands at 75 kDa and bellow.



CD44 e5



**Figure S11. CD44-Tn glycopeptides belonging to the conserved extracellular domain identified in T24 CIGALTIKO cells by nanoLC-HCD-MS/MS.** MS/MS spectrum highlights the typical HexNAC oxonium ion ( $m/z$  204.087) as well as other oxonium ion fragments consistent with GalNAc (highlighted at  $m/z$  138.055 and 144.066). The MS/MS spectrum of the CD44-Tn glycopeptide also highlights peptide ions with GalNAc losses. The predicted and possible glycosites are identified in grey.

**Table S1. Tailored-made Taqman Gene Expression Assay ID's to specifically detect mRNA encoding for total CD44 and 4 splicing variants.**

<b>Gene</b>	<b>TaqMan Gene Expression Assay ID</b>	<b>Exon junction Boundary</b>	<b>Interrogated sequences</b>
<b>CD44 Total</b>	Hs01075864_m1	e3-e4	NM_000610.3 (CD44 v2-10: ENST00000428726.8) NM_001001389.1 (CD44 v3-10: ENST00000415148.6) NM_001001390.1 (CD44 v8-10: ENST00000433892.6) NM_001001391.1(CD44s: ENST00000263398.11) NM_001202555.1 (CD44v10: ENST00000434472.6) NM_001202556.1(CD44s-exon15: ENST00000352818.8) NM_001202557.1 (CD44st: ENST00000442151.6)
<b>CD44 e5-v2</b>	Hs01075866_m1	e5-v6 (corresponds to exon v2 of canonical CD44)	NM_000610.3 (CD44 v2-10: ENST00000428726.8)
<b>CD44 e5-v3</b>	Hs01081480_m1	e5-v6 (corresponds to exon v3 of canonical CD44)	NM_001001389.1 (CD44 v3-10: ENST00000415148.6)
<b>CD44 e5-v8</b>	Hs01081475_m1	e5-v6 (corresponds to exon v8 of canonical CD44)	NM_001001390.1 (CD44 v8-10: ENST00000433892.6)
<b>CD44 e5-e15</b>	Hs01081473_m1	e5-v6 (corresponds to exon e15 of canonical CD44)	NM_001001391.1 (CD44s: ENST00000263398.11) NM_001202557.1 (CD44st: ENST00000442151.6)

<b>CD44sol</b>	Hs01081469_m1	v2(partial)-v3(partial)	NM_001001392.1 (CD44sol: ENST00000278386.10)
----------------	---------------	-------------------------	----------------------------------------------

**Table S2. List of antibodies and lectins used in this manuscript.**

Antigen	Antibody/Lectin	Source	Reference	Clone	Immunogens/targets
<b>CD44</b>	Rabbit polyclonal to CD44 antibody	Abcam	ab157107	Not provided	Synthetic peptide within Human CD44 a.a 650-750
<b>STn</b>	Mouse monoclonal to tag 72	Abcam	ab199002	B72.3+CC49	Membrane-enriched fraction of a human breast carcinoma liver metastasis [B72.3]; Purified TAG-72 protein [CC49]
<b>Tn</b>	Biotinylated and FITC-labeled VVA lectin	Vector Laboratories	B-1235-2 FL-1231-2	Not provided	$\alpha$ - or $\beta$ -linked terminal N-acetylgalactosamine structures
<b>T and ST after desialylation with <math>\alpha</math>-neuraminidase</b>	FITC-labeled PNA lectin	Vector Laboratories	FL-1071	Not provided	Gal ( $\beta$ -1,3) GalNAc structures
<b>B2M</b>	Rabbit monoclonal to beta 2 Microglobulin	Abcam	ab75853	EP2978Y	Not provided

**Table S3. CD44 isoforms identified in 5637 and T24 cells by RNAseq.**

<b>Ensembl</b>	<b>Uniprot</b>	<b>CD44 Isoforms</b>
ENST00000263398.11	P16070-12	<b>CD44s</b>
ENST00000526025.2	E9PKC6	<b>CD44splice E2</b>
ENST00000415148.6	P16070-4	<b>CD44v3-10</b>
ENST00000428726.8	P16070-1	<b>CD44v2-10</b>
ENST00000433892.6	P16070-10	<b>CD44v8-10</b>
ENST00000278386.10	P16070-19	<b>CD44sol</b>
ENST00000434472.6	P16070-11	<b>CD44v10</b>
ENST00000352818.8	P16070-18	<b>CD44s-exon15</b>
ENST00000442151.6	P16070-15 or H0Y5E4	<b>CD44st</b>
ENST00000526669.6	H0YD13	-
ENST00000425428.6	Q86UZ1	-
ENST00000528086.5	-	-
ENST00000526000.6	H0YDW7	-
ENST00000279452.10	H0Y2P0	-
ENST00000531118.5	-	-
ENST00000528455.5	H0YD17	-
ENST00000531873.5	H0YD90	-
ENST00000525209.5	-	-

**Table S4. Glycopeptides identified by nanoLC-HCD/CID-MS/MS in glycoproteogenomics settings for 5637 cells.**

Peptide sequences	5637									
	Glycan chains									
	HexNAc	HexNAcNeuAc	HexHexNAc	HexHexNAcNeuAc	HexHexNAcNeuAc(2)	HexNAc(2)	HexNAc(2)Hex	HexNAc(2)Hex(2)NeuAc	HexNAc(1)Hex(1)Fuc	
<b>CD44v2-10</b>										
LFLPSESKNHLHTTTQMAGTSnTISAGW			x				x			
IPsEsKNHLHTTTQMAGTSSNTISAGW	x		x							
LFLPSESKNHLHTTTQMAGTSnTIsAGW	x		x							
LFLPSESKNHLHttQMAGTSSnTISAGW			x							
DWFsW										x
FsWLF										x
LFLPSESKNHLHTTTQMAGTSSNTISAGW					x					
dWFsW					x					
FsWLF										
LFLPSESKnHLHTTTQmAGtSSnTISAGW	x					x				
fsWLF								x		
FsWLF								x		
LFLPSESKNHLHTTTQmAGtSSNTIsAGW	x						x			
DWFsW			x							
LFLPsEsKNHLHTTTQmAGTSSNTISAGW						x				
LPsEsKNHLHTTTQMAGTSSNTISAGW										
LPSESKNHLHTTTQMAGTSSNTIsAGW										
LFLPsEsKNHLHTTTQmAGTSSNTISAGW										
LFLPSESKnHLHttQMAGTSSNTISAGW										
LPSESKnHLHttQMAGTSSNTISAGW										
LPSESKnHLHttQMAGTSSNTISAGW										
LPSESKnHLHTTTQmAGTSSNTISAGW										
LFLPsEsKNHLHTTTQMAGTSSNTISAGW										
LFLPSEsKNHLHTTTQmAGTsNTISAGW										
LFLPsEsKNHLHTTTQmAGTSSNTISAGW										
LPSESKNHLHTTTQmAGTSSNTISAGW										
sWLF										
sWLF										
FsWLF										
LFLPsEsKnHLHTTTQMAGTSSNTISAGW										
STVHPiPEDESPWITDsTDRIPAtRDQDTf										
dWFsW		x								
sWLF										
FsWLF										
DWFsW										
FsWLF										
DWFsW										
<b>CD44v3-10</b>										
ItDSTDRIPATSTSSnTISAGW										
ItDSTDRIPAtStSSNTISAGW										
ITDSTDRIPATStSSNTISAGW										
ITDSTDRIPATSTsNTIsAGW		x								x
ITDsTDRIPATSTsNTIsAGW	x			x			x			
ITDsTDRIPATSTsNTIsAGW				x			x			
ITDSTDRIPATSTsNTISAGW				x						
ITDSTDRIPATSTsNTISAGW	x									
ITDSTDRIPAtStSSnTIsAGW	x			x						
ITDSTDRIPATSTsNTIsAGW	x									
ITDSTDRIPATSTsNTIsAGW						x				
ITDSTDRIPATSTsNTIsAGW										
ITDSTDRIPATSTsNTIsAGW										
<b>CD44s</b>										
StVHPiPEDEsPWITDSTDRIPATRDQDTf	x			x				x		
stVHPiPEDEsPWITDsTDRIPATRDQDTf	x			x						
ITDsTDRIPATRDQDTf										x
stVHPiPEDESPWITDSTDRIPAtRDQDTf										
STVHPiPEDESPWITDsTDRIPATRDQDTf										
StVHPiPEDESPWITDSTDRIPAtRDQDTf										
ItDsTDRIPATRDQDTf										
stVHPiPEDEsPWITDsTDRIPATRDQDTf										
ITDsTDRIPAtRDQDTf										
ItDsTDRIPAtRDQDTf										
STVHPiPEDESPWITDsTDRIPAtRDQDTf										
stVHPiPEDESPWITDsTDRIPATRDQDTf										
ItDsTDRIPAtRDQDTf										
ITDsTDRIPAtRDQDTf										
STVhPiPEDEsPWITDsTDRIPATRDQDTf										
STVHPiPEDEsPWITDsTDRIPATRDQDTf										
<b>CD44sol</b>										
qWScGGQKAKWtQRRGQVSGnGAF							x			
GEQGVVRNsRPVY							x			
GEQGVVRNsRPVY	x									
GEQGVVRNsRPVYDs	x			x						
QWScGGQKAKW	x									
tQRRGQVVsGNGAF				x						
TQRRGQVVsGnGAF	x									
etcSLHcSQSKVWAEKASDQQW							x		x	
scGGQKAKWTQRRGQVVsGnGAF	x									
ETcSLHcSQSKVW							x			
tQRRGQVVsGNGAFGEQGVVRNsRPVYDs							x		x	
TQRRGQVVsGNGAF				x						
gEQGVVRNsRPVY				x						
scGGQKAKWTQRRGQVVsGNGAF										
TQRRGQVVsGNGAF										
EtcSLHcSQSKVW									x	
AEEKASDQQWQW										x
ETcSLHcSQSKVWAEKASDQQWQW										
qWScGGQKAKWTQRRGQVVsGNGAF										
QWScGGQKAKWTQRRGQVVsGNGAF										

**Table S5. Glycopeptides identified by nanoLC-HCD/CID-MS/MS in glycoproteogenomics settings for T24 cell line.**

Peptide sequence	T24							
	Glycan chains							
	HexNAc	HexNAcNeuAc	HexHexNAc	HexHexNAcNeuAc	HexHexNAcNeuAc(2)	HexNAc(2)Hex	HexNAc(2)Hex(2)NeuAc	HexNAc(1)Hex(1)Fuc
<b>CD44v2-10</b>								
LFLPSESKNHLHTTTQMAGTSnTISAGW			x			x		
IPsEsKNHLHTTTQMAGTSSNTISAGW	x		x					
LFLPSESKNHLHTTTQMAGTSnTisAGW								
LFLPSESKNHLHtttQMAGTSSnTISAGW								
DWFsW								
FsWLF								x
LFLPSESKNHLHTTTQMAGTSSNTISAGW					x			
dWFsW								
FsWLF						x		
LFLPSESKnHLHTTTQmAGTSSnTISAGW							x	
fsWLF								
FsWLF						x		
LFLPSESKNHLHTTTQmAGTSSNTisAGW								
DWFsW								
LFLPsEsKNHLHTTTQmAGTSSNTISAGW						x		
LPsEsKNHLHTTTQMAGTSSNTISAGW								
LPSESKNHLHTTTQMAGTSSNTisAGW					x			x
LFLPsESKNHLHTTTQmAGTSSNTISAGW		x						
LFLPSESKnHLHtttQMAGTSSNTISAGW						x	x	
LPSESKnHLHtttQMAGTSSNTISAGW	x		x				x	
LPSESKnHLHTttQMAGTSSNTISAGW							x	
LPSESKnHLHTttQMAGTSSNTISAGW							x	
LPSESKNHLHTTTQmAGTSSNTISAGW								x
LFLPsESKNHLHTTTQMAGTSSNTISAGW		x						x
LFLPSEsKNHLHTTTQmAGTSSNTISAGW			x					
LFLPsEsKNHLHTTTQmAGTSSNTISAGW	x							
LPSESKNHLHTTTQmAGTSSNTISAGW		x						
sWLF								x
sWLF							x	
FsWLF				x				
LFLPsESKnHLHTTTQMAGTSSNTISAGW								
STVHPIPEDSPWITDsTDRIpAtRDQDTf								
dWFsW								
sWLF								
FsWLF								
DWFsW								
FsWLF								
DWFsW								
<b>CD44v3-10</b>								
ITDSTDRIpATsSSnTISAGW						x		
ITDSTDRIpATsTSSNTISAGW	x						x	
ITDSTDRIpATsTSSNTISAGW			x					
ITDSTDRIpATsSNTisAGW								
ITDStDRIPATsTssNTisAGW								
ITDStDRIPATsTSSNTisAGW								
iTDDTDRIpATsTSSNTISAGW								
iTDDTDRIpATsTSSNTisAGW								
ITDSTDRIpATsTSSnTisAGW								
ITDSTDRIpATsTssnTisAGW								
ITDSTDRIpATsTssnTisAGW								
ITDSTDRIpATsTssnTisAGW								
ITDSTDRIpATsTssnTisAGW								
<b>CD44s</b>								
StVHPIPEDsPWITDsTDRIpAtRDQDTf	x		x				x	
stVHPIPEDsPWITDsTDRIpAtRDQDTf	x		x					
ITDStDRIPAtRDQDTf								x
stVHPIPEDsPWITDsTDRIpAtRDQDTf			x				x	
STVHPIPEDsPWITDsTDRIpAtRDQDTf		x						x
StVHPIPEDsPWITDsTDRIpAtRDQDTf	x		x				x	
ItDsTDRIpAtRDQDTf						x		
stVHPIPEDsPWITDsTDRIpAtRDQDTf		x						x
ITDStDRIPAtRDQDTf		x						x
ITDStDRIPAtRDQDTf					x			
STVHPIPEDsPWITDsTDRIpAtRDQDTf								
stVHPIPEDsPWITDsTDRIpAtRDQDTf	x					x		
ItDStDRIPAtRDQDTf								
ITDStDRIPAtRDQDTf								
STVhPIPEDsPWITDsTDRIpAtRDQDTf								
STVHPIPEDsPWITDsTDRIpAtRDQDTf								
<b>CD44sol</b>								
qWScGGQKAKWtQRRGQVSGnGAF								
GEQGVVRNsRPVY								
GEQGVVRNsRPVY	x							
GEQGVVRNsRPVYDs								
QWScGGQKAKW								
tQRRGQVSGnGAF								
TQRRGQVSGnGAF								
etcSLHcSQSKKVVWAEKASDQQW						x	x	
scGGQKAKWTQRRGQVSGnGAF	x							
ETcSLHcSQSKKVV								
tQRRGQVSGnGAFGEQGVVRNsRPVYDs								
TQRRGQVSGnGAF								
gEQGVVRNsRPVY			x					
scGGQKAKWTQRRGQVSGnGAF								x
TQRRGQVSGnGAF								
ETcSLHcSQSKKVV				x				
AEEKASDQQWQW								
ETcSLHcSQSKKVVWAEKASDQQWQW		x						
ETcSLHcSQSKKVVWAEKASDQQWQW			x			x		
qWScGGQKAKWTQRRGQVSGnGAF								x
QWScGGQKAKWTQRRGQVSGnGAF		x						



**Table S6. Glycopeptides identified by nanoLC-HCD/CID-MS/MS in glycoproteogenomics settings for CD44s<sup>high</sup> MIBC showing areas of CD44 and STn co-localization.**

Peptide sequence	MIBC	
	Glycan chains	
	HexNAc	HexNAcNeuAc
<b>CD44v2-10</b>		
LFLPSESKNHLHTTTQMAGTSsNTISAGW		
IPsEsKNHLHTTTQMAGTSSNTISAGW		
LFLPSESKNHLHTTTQMAGTSsNTISAGW		
LFLPSESKNHLHttQMAGTSSnTISAGW		
DWFsW		
FsWLF		
LFLPSESKNHLHTTTQMAGTSSNTISAGW		
dWFsW		
FsWLF		
LFLPSESKnHLHTTTQmAGtSSnTISAGW		
fsWLF		
FsWLF		
LFLPSESKNHLHTTTQmAGtSSNTISAGW		
DWFsW		
LFLPsEsKNHLHTTTQmAGTSSNTISAGW		
LPsEsKNHLHTTTQMAGTSSNTISAGW		
LPSESKNHLHTTTQMAGTSSNTISAGW		
LFLPsESKNHLHTTTQmAGTSSNTISAGW		
LFLPSESKnHLHttQMAGTSSNTISAGW		
LPSESKnHLHttQMAGTSSNTISAGW		
LPSESKnHLHttQMAGTSSNTISAGW		
LPSESKNHLHTTTQmAGTSSNTISAGW		
LFLPsESKNHLHTTTQMAGTSSNTISAGW		
LFLPsEsKNHLHTTTQmAGTSSNTISAGW		
LPSESKNHLHTTTQmAGTSSNTISAGW	x	
sWLF	x	
sWLF		
FsWLF		
LFLPsESKnHLHTTTQMAGTSSNTISAGW	x	x
STVHPIPDEdSPWITDsTDRIpAtRDQDTf		
dWFsW		
sWLF		x
FsWLF	x	
DWFsW	x	
FsWLF		x
DWFsW		x
<b>CD44v3-10</b>		
ITDSTRIPATSTSSnTISAGW		
ITDSTRIPATStSSNTISAGW		
ITDSTRIPATStSSNTISAGW		
ITDSTRIPATSTSSnTISAGW		
ITDSTRIPATSTSSnTISAGW		
ITDSTRIPATSTSSnTISAGW		
ITDSTRIPATSTSSnTISAGW		
ITDSTRIPATSTSSnTISAGW		
ITDSTRIPATStSSnTISAGW		
ITDSTRIPATSTSSnTISAGW		
ITDSTRIPATSTSSnTISAGW	x	x
ITDSTRIPATSTSSnTISAGW		x
ITDSTRIPATSTSSnTISAGW	x	
<b>CD44s</b>		
StVHPIPDEdSPWITDsTDRIpAtRDQDTf		
stVHPIPDEdSPWITDsTDRIpAtRDQDTf		
ITDsTDRIpAtRDQDTf		
stVHPIPDEdSPWITDsTDRIpAtRDQDTf		
StVHPIPDEdSPWITDsTDRIpAtRDQDTf		
ItDsTDRIpAtRDQDTf		
stVHPIPDEdSPWITDsTDRIpAtRDQDTf		
ITDsTDRIpAtRDQDTf		
ITDsTDRIpAtRDQDTf	x	
STVHPIPDEdSPWITDsTDRIpAtRDQDTf	x	
stVHPIPDEdSPWITDsTDRIpAtRDQDTf		
ItDsTDRIpAtRDQDTf		x
ITDsTDRIpAtRDQDTf	x	x
STVhPIPDEdSPWITDsTDRIpAtRDQDTf		x
STVHPIPDEdSPWITDsTDRIpAtRDQDTf	x	x
<b>CD44sol</b>		
qWScGGQKAKWTQRRGQQVsGnGAF		
GEQGVVRNsRPVY		
GEQGVVRNsRPVY	x	
GEQGVVRNsRPVYDs		
QWScGGQKAKW		
tQRRGQQVsGnGAF		
TQRRGQQVsGnGAF		
etcSLHcSQSKVWAEKASDQQW		
scGGQKAKWTQRRGQQVsGnGAF	x	
ETcSLHcSQSKVW		
tQRRGQQVsGnGAFGEQGVVRNsRPVYDs		
TQRRGQQVsGnGAF		
gEQGVVRNsRPVY		x
scGGQKAKWTQRRGQQVsGnGAF		
TQRRGQQVsGnGAF		
EtcSLHcSQSKVW		
AEEKASDQQWQW		
ETcSLHcSQSKVWAEKASDQQWQW		
qWScGGQKAKWTQRRGQQVsGnGAF		
QWScGGQKAKWTQRRGQQVsGnGAF		x

**Table S7. Glycopeptides identified by nanoLC-HCD/CID-MS/MS in glycoproteogenomics settings for T24 CIGALT1 KO cell model.**

Peptide sequence	T24 C1GALT1 KO		
	Glycan chains		
	HexNAc	HexNAcNeuAc	HexNAc(2)
<b>CD44v2-10</b>			
LFLPSEKNHLHTTTQMAGTSsNTISAGW			
IPsEsKNHLHTTTQMAGTSSNTISAGW			
LFLPSEKNHLHttTQMAGTSsNTISAGW			
LFLPSEKNHLHttQMAGTSSnTISAGW			
DWFsW			
FsWLF			
LFLPSEKNHLHTTTQMAGTSSNTISAGW			
dWFsW			
FsWLF			
LFLPSEKnhLHTTTQmAGtSSnTISAGW			
fsWLF			
FsWLF			
LFLPSEKNHLHttTQmAGtSSNTISAGW			
DWFsW			
LFLPsEsKNHLHTTTQmAGTSSNTISAGW			
LPsEsKNHLHTTTQMAGTSSNTISAGW			
LPSEKNHLHttTQMAGTSSNTISAGW			
LFLPsESKNHLHTTTQmAGTSSNTISAGW			
LFLPSEKnhLHttQMAGTSSNTISAGW			
LPSEKnhLHttQMAGTSSNTISAGW			
LPSEKnhLHttQMAGTSSNTISAGW			
LPSEKnhLHttQMAGTSSNTISAGW			
LPSEKnhLHttQMAGTSSNTISAGW			
LFLPsESKNHLHTTTQMAGTSSNTISAGW			
LFLPsEsKNHLHTTTQmAGTSsNTISAGW			
LFLPsEsKNHLHTTTQmAGTSSNTISAGW	x		
LPSEKnhLHttTQmAGTSSNTISAGW			
sWLF			
sWLF			
FsWLF			
LFLPsESKnhLHttTQMAGTSSNTISAGW	x	x	
STVHPIPEDESPWITDsTRIPAtRDQDTf	x		
dWFsW			
sWLF			
FsWLF			
DWFsW			
FsWLF			
DWFsW			
<b>CD44v3-10</b>			
ItDSTRIPATSTSSnTISAGW			
ItDSTRIPAtStSSNTISAGW			
ITDSTRIPATStSSNTISAGW			
ITDSTRIPATSTsNTISAGW			
ITDStDRIPATSTsNTISAGW			
ITDStDRIPATSTsNTISAGW			
iTDSTRIPATSTsNTISAGW			
iTDSTRIPATSTSSnTISAGW			
ITDSTRIPATStSSnTISAGW			
ITDSTRIPATSTsNTISAGW			
ITDSTRIPATSTSSnTISAGW			
ITDSTRIPATStSSnTISAGW			
ITDSTRIPATSTsNTISAGW			
ITDSTRIPATSTSSnTISAGW			
ITDSTRIPATSTSSnTISAGW			
ITDSTRIPATSTSSnTISAGW			
ITDSTRIPATSTSSnTISAGW			
<b>CD44s</b>			
StVHPIPEDEsPWITDsTRIPATRDQDTf			
stVHPIPEDEsPWITDsTRIPATRDQDTf			
ITDStDRIPATRDQDTf			
sTVHPIPEDESPWITDsTRIPAtRDQDTf			
StVHPIPEDESPWITDsTRIPATRDQDTf			
StVHPIPEDESPWITDsTRIPAtRDQDTf			
ItDsTRIPATRDQDTf			
sTVHPIPEDEsPWITDsTRIPATRDQDTf			
ITDStDRIPAtRDQDTf			
ItDsTRIPAtRDQDTf	x		
STVHPIPEDESPWITDsTRIPAtRDQDTf	x		
sTVHPIPEDESPWITDsTRIPATRDQDTf			
ItDsTRIPAtRDQDTf			
ITDStDRIPAtRDQDTf			
STVhPIPEDESPWITDsTRIPATRDQDTf			
STVHPIPEDEsPWITDStDRIPAtRDQDTf			
<b>CD44sol</b>			
qWScGGQKAKWtQRRGQQVSGnGAF			
GEQGVVRNsRPVY			
GEQGVVRNsRPVY			
GEQGVVRNsRPVYDs			
QWScGGQKAKW			
tQRRGQQVsGNGAF			
TQRRGQQVsGnGAF	x		
etcSLHcSQSKKVVAAEKAQDQW			
scGGQKAKWtQRRGQQVsGnGAF	x		
ETcSLHcSQSKKVV			
tQRRGQQVSGNGAFGEQGVVRNsRPVYDs			
TQRRGQQVsGNGAF			x
gEQGVVRNsRPVY			
scGGQKAKWtQRRGQQVSGNGAF			
TQRRGQQVsGNGAF			
EtcSLHcSQSKKVV			
AEEKAsDQQWQW			
ETcSLHcSQSKKVVAAEKAQDQWQW			
qWScGGQKAKWtQRRGQQVsGNGAF			
QWScGGQKAKWtQRRGQQVsGNGAF			