

Identification, Mapping and Relative Quantitation of SARS-CoV-2 Spike Glycopeptides by Mass-Retention Time Fingerprinting

Chalk, R.¹, Greenland, W.², Moreira, T.¹, Coker, J.¹, Mukhopadhyay S.M.M¹, Williams, E.¹, Manning, C.¹, Bohstedt, T.¹, McCrorie, R.¹, Fernandez-Cid, A.¹, Burgess-Brown, N.A.¹

¹Centre for Medicines Discovery, ORCRB, Oxford University, OX3 7DQ, UK

²Agilent Technologies, Lakeside, Cheadle Royal Business Park, Cheadle, Cheshire, SK8 3GR, UK

Abstract

We describe a novel analytical method for rapid and robust identification, mapping and relative quantitation of glycopeptides from SARS-CoV-2 Spike protein. The method may be executed using any LC-TOF mass spectrometer, requires no specialised knowledge of glycan analysis and makes use of the differential resolving power of reversed phase HPLC. While this separation technique resolves peptides with high efficiency, glycans are resolved poorly, if at all. Consequently, glycopeptides consisting of the same peptide bearing different glycan structures will all possess very similar retention times and co-elute. While this has previously been viewed as a disadvantage, we show that shared retention time can be used to map multiple glycan species to the same peptide and location. In combination with MSMS and pseudo MS3, we have constructed a detailed mass-retention time database for Spike. This database allows any ESI-TOF equipped lab to reliably identify and quantify spike glycans from a single overnight elastase protein digest in less than 90 minutes.

Key words:

SARS-CoV-2, Spike, RBD, Glycoprotein, Glycopeptide, Glycan, Mass Spectrometry, HPLC, Database

Introduction

Glycosylation is known to play an important role in the efficacy and antigenicity of therapeutic proteins [1-3]. The current SARS-CoV-2 pandemic has spurred urgent research, much of it devoted to preparing vaccines, therapeutic antibodies or antibody tests based on Spike protein, the virus's primary surface antigen [4]. This 145 kDa protein forms a trimer [5] with each subunit bearing twenty-two potential N-linked glycosylation sites and two O-linked sites of which approximately seventeen are occupied [5]. The unusually heavy and complex glycosylation observed in Spike protein is believed to play an important role in the pathogenicity of SARS-CoV-2 by mimicking host cell glycans and allowing the virus to evade the normal immune response [6]. Analysis of expressed Spike protein by mass spectrometry presents unique challenges in terms of its size and the number and complexity of its glycans. These challenges have been commendably met to date by laboratories with wide experience in glycan analysis and access to very sensitive, high-end nano-LC-MSMS mass spectrometers [1, 7-9]. However, in our laboratory and in others a rapid and more robust methodology is needed for routine analysis of different batches of expressed Spike protein. In addition, any method which is reliant on LC-MSMS of glycopeptides may not necessarily detect specific glycans which fail to fragment under the conditions selected. LC-MS, by contrast, generates a mass, retention time and relative abundance for all ionizable species. We have developed a simple Mass-Retention Time Fingerprinting (MRTF) method for rapid and robust identification, mapping and relative quantitation of Spike glycans. Overnight digestion using a single enzyme followed by a 65-minute LC-MS run using any accurate mass instrument are the only experimental requirements. The resulting LC-MS data contains accurate mass, retention time and relative abundance values for each glycopeptide component. This dataset needs only to be matched against the pre-existing Spike glycopeptide database reported here, as shown in Figure 1. We describe this method as "analytical mode", which is both conceptually simple to understand, and straightforward to implement in a typical mass spectrometry laboratory. For scientific completeness, we also describe the "discovery mode" which we have used to generate the data for our Mass-Retention Time Fingerprinting

Table 1b. Spike elastase glycopeptide mass retention time database (PCDL) containing data for 140 observed glycopeptides and data for a further 306 inferred glycopeptides (RT 32-60 min and key)

Glycan posn.	RT (min)	Mass	Glycopeptide	Observed/ Inferred
N149	32.4	3333.3262	YYYYHKNKSWM Man9	Inf
	32.463	2847.1678	YYYYHKNKSWM Man6	Obs
	32.463	3171.2734	YYYYHKNKSWM Man8	Obs
	32.6	2767.1682	YYYYHKNKSWM GO	Inf
	32.6	2913.2261	YYYYHKNKSWM GOF	Inf
	32.6	2036.9038	YYYYHKNKSWM Man1	Inf
	32.6	2198.9566	YYYYHKNKSWM Man2	Inf
	32.6	2361.0094	YYYYHKNKSWM Man3	Inf
	32.6	2523.0622	YYYYHKNKSWM Man4	Inf
	32.6	2685.115	YYYYHKNKSWM Man5	Inf
	32.682	3009.2206	YYYYHKNKSWM Man7	Obs
	33.611	2694.1185	GP4 Man6	Obs
33.8	2760.1768	GP4 GOF	Inf	
33.819	2532.0657	GP4 Man5	Obs	
N17	32.463	2828.1148	CVNLTTRT Man9	Obs
	33.8	2261.9568	CVNLTTRT GO	Inf
	33.8	2408.0147	CVNLTTRT GOF	Inf
	33.8	1531.6924	CVNLTTRT Man1	Inf
	33.8	1693.7452	CVNLTTRT Man2	Inf
	33.8	1855.798	CVNLTTRT Man3	Inf
	33.8	2017.8508	CVNLTTRT Man4	Inf
	33.8	2179.9036	CVNLTTRT Man5	Inf
	33.8	2341.9564	CVNLTTRT Man6	Inf
	33.8	2504.0092	CVNLTTRT Man7	Inf
	33.841	2666.062	CVNLTTRT Man8	Obs
	N343	34.6	1848.7511	VFNAT GO
34.6		1994.809	VFNAT GOF	Inf
34.6		1118.4867	VFNAT Man1	Inf
34.6		1280.5395	VFNAT Man2	Inf
34.6		1442.5923	VFNAT Man3	Inf
34.6		2252.8563	VFNAT Man8	Inf
34.6		2414.9091	VFNAT Man9	Inf
34.656		1928.7507	VFNAT Man6	Obs
34.674		2090.8035	VFNAT Man7	Obs
34.682		1766.6979	VFNAT Man5	Obs
34.688		1604.6451	VFNAT Man4	Obs
N717		35.101	2195.8349	NFTI Man8
	35.353	1709.6765	NFTI Man5	Obs
	35.373	1385.5709	NFTI Man3	Obs
	35.374	2033.7821	NFTI Man7	Obs
	35.8	2357.8877	NFTI Man9	Inf
	35.86	1871.7293	NFTI Man6	Obs
	36.2	2406.9307	NFTI A1	Inf
	36.2	2552.9886	NFTI A1F	Inf
	36.2	2068.0261	NFTI A2	Inf
	36.2	2844.084	NFTI A2F	Inf
	36.2	1791.7297	NFTI GO	Inf
	36.2	1994.8091	NFTI GO +GlcNAc	Obs
36.2	2140.867	NFTI GOF +GlcNAc	Obs	
36.2	1953.7825	NFTI G1	Inf	
36.2	2115.8353	NFTI G2	Inf	
36.2	2261.8932	NFTI G2F	Inf	
36.2	1061.4653	NFTI Man1	Inf	
36.2	1223.5181	NFTI Man2	Inf	
36.295	1547.6237	NFTI Man4	Obs	
36.599	2099.8404	NFTI G1F	Obs	
36.77	2089.8308	NFTI Man5+380.2	Obs	
36.896	1937.7876	NFTI GOF	Obs	
N61	35.4	2765.11	PFFSNVTW G2F	Inf
	35.4	1564.6821	PFFSNVTW Man1	Inf
	35.4	1726.7349	PFFSNVTW Man2	Inf
	35.4	1888.7877	PFFSNVTW Man3	Inf
	35.4	2050.8405	PFFSNVTW Man4	Inf
	35.4	2212.8933	PFFSNVTW Man5	Inf
	35.4	2536.9989	PFFSNVTW Man7	Inf
	35.4	2699.0517	PFFSNVTW Man8	Inf
	35.4	2861.1045	PFFSNVTW Man9	Inf
	35.437	2619.0521	PFFSNVTW G2	Obs
	35.601	2456.9993	PFFSNVTW G1	Obs
	36	2498.0259	PFFSNVTW GO +GlcNAc	Inf
36	2441.0044	PFFSNVTW GOF	Inf	
36	2644.0838	PFFSNVTW GOF +GlcNAc	Inf	
36	2603.0572	PFFSNVTW G1F	Inf	
36.025	2294.9465	PFFSNVTW GO	Obs	
36.854	2374.9461	PFFSNVTW Man6	Obs	
N343	36.5	2181.8836	FGEVFNAT GO	Inf
	36.5	2327.9415	FGEVFNAT GOF	Inf
	36.5	1451.6192	FGEVFNAT Man1	Inf
	36.5	1613.672	FGEVFNAT Man2	Inf
	36.5	1775.7248	FGEVFNAT Man3	Inf
	36.5	1937.7776	FGEVFNAT Man4	Inf
	36.5	2261.8832	FGEVFNAT Man6	Inf
	36.5	2423.936	FGEVFNAT Man7	Inf
	36.5	2585.9888	FGEVFNAT Man8	Inf
	36.5	2748.0416	FGEVFNAT Man9	Inf
	36.528	2099.8304	FGEVFNAT Man5	Obs
	N343	36.8	2552.0511	LCPFGEVFNAT GO
36.8		2698.109	LCPFGEVFNAT GOF	Inf
36.8		1821.7867	LCPFGEVFNAT Man1	Inf
36.8		1983.8395	LCPFGEVFNAT Man2	Inf
36.8		2145.8923	LCPFGEVFNAT Man3	Inf
36.8		2307.9451	LCPFGEVFNAT Man4	Inf
36.8		2469.9979	LCPFGEVFNAT Man5	Inf
36.8		2794.1035	LCPFGEVFNAT Man7	Inf
36.8		2956.1563	LCPFGEVFNAT Man8	Inf
36.8		3118.2091	LCPFGEVFNAT Man9	Inf
36.864		2632.0507	LCPFGEVFNAT Man6	Obs
GP6		40.948	1570.7208	GP6
	41.645	1881.7363	GP6 +311.0	Obs
	42.6	1855.7569	PNITN GO	Inf
	42.6	2001.8148	PNITN GOF	Inf
	42.6	1125.4925	PNITN Man1	Inf
	42.6	1287.5453	PNITN Man2	Inf
	42.6	1611.6509	PNITN Man4	Inf
	42.6	1773.7037	PNITN Man5	Inf
	42.6	1935.7565	PNITN Man6	Inf
	42.6	2097.8093	PNITN Man7	Inf
	42.6	2259.8621	PNITN Man8	Inf
	42.6	2421.9149	PNITN Man9	Inf
42.631	1449.5981	PNITN Man3	Obs	
N801	46.236	3854.7283	KQIYKTPPKIDFGGFNFS G2F	Obs
	46.369	3788.67	KQIYKTPPKIDFGGFNFS Man8	Obs
	46.424	3626.6226	KQIYKTPPKIDFGGFNFS GOF +96.0	Obs
	46.5	3546.6176	KQIYKTPPKIDFGGFNFS G1	Inf
	46.5	3692.6755	KQIYKTPPKIDFGGFNFS G1F	Inf
	46.5	3708.6704	KQIYKTPPKIDFGGFNFS G2	Inf
	46.5	3950.7228	KQIYKTPPKIDFGGFNFS Man9	Inf
	46.506	3464.5644	KQIYKTPPKIDFGGFNFS Man6	Obs
	46.506	3626.6172	KQIYKTPPKIDFGGFNFS Man7	Obs
	46.509	3530.6227	KQIYKTPPKIDFGGFNFS GOF	Obs
	46.509	3733.7021	KQIYKTPPKIDFGGFNFS GOF +GlcNAc	Obs
	46.538	3140.4588	KQIYKTPPKIDFGGFNFS Man4	Obs
46.568	2978.406	KQIYKTPPKIDFGGFNFS Man3	Obs	
46.575	2654.3004	KQIYKTPPKIDFGGFNFS Man1	Obs	
46.581	2816.3532	KQIYKTPPKIDFGGFNFS Man2	Obs	
46.601	3302.5116	KQIYKTPPKIDFGGFNFS Man5	Obs	
46.668	3384.5648	KQIYKTPPKIDFGGFNFS GO	Obs	
46.966	3587.6442	KQIYKTPPKIDFGGFNFS GO +GlcNAc	Obs	
47.5	4145.8237	KQIYKTPPKIDFGGFNFS A1F	Inf	
47.5	4290.8612	KQIYKTPPKIDFGGFNFS A2	Inf	
47.5	4436.9191	KQIYKTPPKIDFGGFNFS A2F	Inf	
48.492	3999.7658	KQIYKTPPKIDFGGFNFS A1	Obs	
N801	48.1	3630.5282	YKTPPKIDFGGFNFS A1	Inf
	48.1	3776.5861	YKTPPKIDFGGFNFS A1F	Inf
	48.1	3921.6236	YKTPPKIDFGGFNFS A2	Inf
	48.1	4067.6815	YKTPPKIDFGGFNFS A2F	Inf
	48.121	3485.4907	YKTPPKIDFGGFNFS G2F	Obs
	48.27	4145.822	YKTPPKIDFGGFNFS Very complex	Obs
	48.332	3419.4324	YKTPPKIDFGGFNFS Man8	Obs
	48.339	3323.4379	YKTPPKIDFGGFNFS G1F	Obs
	48.408	3257.3796	YKTPPKIDFGGFNFS Man7	Obs
	48.453	3364.4645	YKTPPKIDFGGFNFS GOF+GlcNAc	Obs
	48.462	3161.3851	YKTPPKIDFGGFNFS GOF	Obs
	48.499	3983.7666	YKTPPKIDFGGFNFS Very complex	Obs
48.518	3095.3268	YKTPPKIDFGGFNFS Man6	Obs	
48.6	1716.8512	YKTPPKIDFGGFNFS	Inf	
48.6	2123.01	YKTPPKIDFGGFNFS (GlcNAc)2	Inf	
48.6	3015.3272	YKTPPKIDFGGFNFS GO	Inf	
48.6	3177.38	YKTPPKIDFGGFNFS G1	Inf	
48.6	3339.4328	YKTPPKIDFGGFNFS G2	Inf	
48.6	1919.9306	YKTPPKIDFGGFNFS GlcNAc stump	Inf	
48.6	2285.0628	YKTPPKIDFGGFNFS Man1	Inf	
48.6	2447.1156	YKTPPKIDFGGFNFS Man2	Inf	
48.6	2609.1684	YKTPPKIDFGGFNFS Man3	Inf	
48.6	2771.2212	YKTPPKIDFGGFNFS Man4	Inf	
48.6	3581.4852	YKTPPKIDFGGFNFS Man9	Inf	
48.623	3218.4066	YKTPPKIDFGGFNFS GO+GlcNAc	Obs	
48.647	2933.274	YKTPPKIDFGGFNFS Man5	Obs	
N343	50.406	3614.5291	NLCPFGEVFNAT Complex	Obs
	51.783	2908.1464	NLCPFGEVFNAT Man7	Obs
	51.942	3136.2575	NLCPFGEVFNAT G2F	Obs
	52.073	2131.0262	NLCPFGEVFNAT glyco	Obs
	52.167	2421.988	NLCPFGEVFNAT Man4	Obs
	52.185	2259.9352	NLCPFGEVFNAT Man3	Obs
	52.229	2666.094	NLCPFGEVFNAT GO	Obs
	52.265	2746.0936	NLCPFGEVFNAT Man6	Obs
	52.3	3281.295	NLCPFGEVFNAT A1	Inf
	52.3	3572.3904	NLCPFGEVFNAT A2	Inf
	52.3	2828.1468	NLCPFGEVFNAT G1	Inf
	52.3	2990.1996	NLCPFGEVFNAT G2	Inf
52.32	2974.2047	NLCPFGEVFNAT G1F	Obs	
52.445	3015.2313	NLCPFGEVFNAT GOF+GlcNAc	Obs	
52.5	1935.8296	NLCPFGEVFNAT Man1	Inf	
52.5	2097.8824	NLCPFGEVFNAT Man2	Inf	
52.5	3070.1992	NLCPFGEVFNAT Man8	Inf	
52.5	3232.252	NLCPFGEVFNAT Man9	Inf	
52.575	2584.0408	NLCPFGEVFNAT Man5	Obs	
52.632	2812.1519	NLCPFGEVFNAT GOF	Obs	
53.304	3427.3529	NLCPFGEVFNAT A1F	Obs	
53.961	3718.4483	NLCPFGEVFNAT A2F	Obs	
GP5	60.483	1861.7463	GP5	Obs

Key to Glycan Structures

Monosaccharide symbol

- Galactose
- Mannose
- ▲ Fucose
- N-Acetylglucosamine (GlcNAc)
- ◆ N-Acetylneuraminic acid (Neu5Ac)

The complete Spike PCDL database is available to download in .cdb or .xlsx format here:

https://zenodo.org/record/3958218#.Xxn_BChKhoY

Figure 3. Combined Extracted Ion Chromatogram (EIC) for 27 isoforms of glycopeptide GEVFNAT (N343) within +/- 2 min retention time window from RBD. Only three glycans are labelled, the remainder are listed in the accompanying table 2, below.

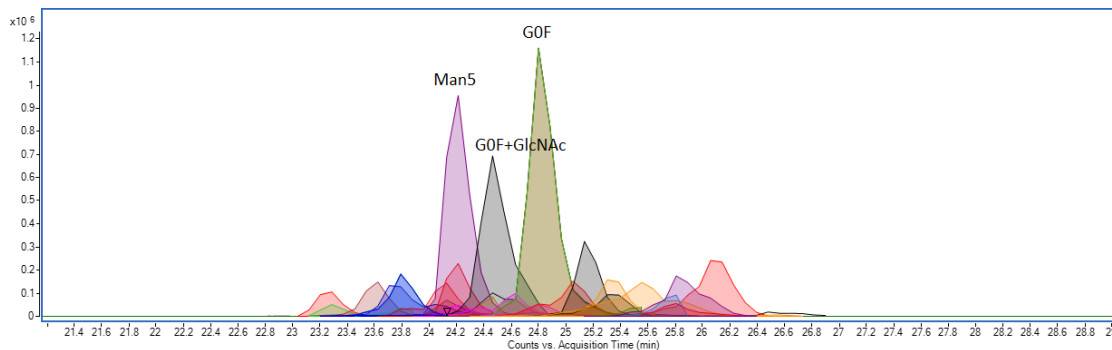


Table 2. GEVFNAT glycopeptide (N343) isoforms from RBD shown in Figure 3

Name	Mass	RT	Volume	ppm error	Name	Mass	RT	Volume	ppm error
GEVFNAT Complex NeuAc (F)2	2838.1078	23.26	733674	-0.2	GEVFNAT G1F	2342.9180	24.62	1090359	3.4
GEVFNAT Man8	2438.9113	23.30	744169	3.7	GEVFNAT G0	2034.8089	24.80	950868	3.1
GEVFNAT Man7	2276.8615	23.61	1867525	2.7	GEVFNAT G0F	2180.8688	24.82	12980259	2.0
GEVFNAT Complex NeuAc F	2692.0523	23.77	1266219	0.5	GEVFNAT A1(F)2-Gal+GlcNAc	2983.1427	24.91	732781	5.3
GEVFNAT Man6	2114.8091	23.80	2432950	2.7	(G)EVFNAT G0	1977.7887	25.05	2947447	2.5
GEVFNAT G1(F)2	2488.9768	24.09	1292274	2.8	GEVFNAT G0F+GlcNAc	2383.9464	25.16	3836871	2.6
GEVFNAT G2F	2504.9755	24.14	1007345	1.3	GEVFNAT A1F	2796.0641	25.33	1512501	3.6
GEVFNAT Man4	1790.7051	24.20	3040952	2.3	GEVFNAT A1(F)2-Gal+GlcNAc	2983.1444	25.34	1320994	4.8
GEVFNAT Man5	1952.7583	24.20	12492713	1.9	GEVFNAT A1F-Gal	2634.0105	25.51	608358	4.1
GEVFNAT Man3	1628.6504	24.20	666739	3.7	GEVFNAT A2F	3087.1629	25.77	805142	2.2
GEVFNAT G1F	2342.9183	24.37	1333546	3.3	GEVFNAT A1F	2796.0644	25.80	541989	3.5
GEVFNAT G0F	2180.8638	24.45	953034	4.3	GEVFNAT A1F-Gal	2634.0125	25.86	3692097	3.4
GEVFNAT G0F+GlcNAc	2383.9436	24.47	9038631	3.7					

Figure 4 illustrates a complete glycan fragmentation series for RBD glycopeptide GEVFNAT-Man5 showing the peptide stump (GEVFNAT-GlcNAc) and mannose ladders. Calculated mass errors are shown in table 3.

Figure 4. Complete glycan fragmentation series for RBD glycopeptide GEVFNAT-Man5 (N343)

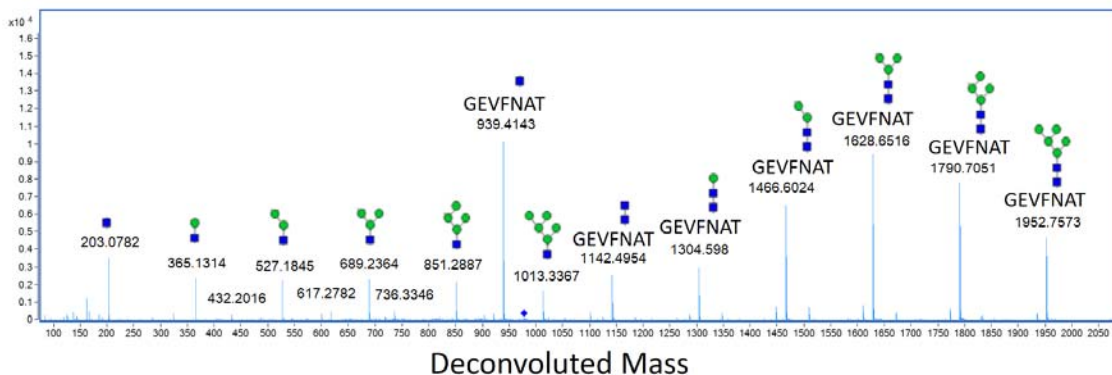


Table 3. Glycan assignment and mass errors (parts per million) for RDB glycopeptide GEVFNAT-Man5

Deconvoluted Mass Obs	Formula	Mass Calc	ppm	Assignment
203.0782	C ₈ H ₁₃ N O ₅	203.0794	-5.9	GlcNAc
365.1314	C ₈ H ₁₃ N O ₅ (C ₆ H ₁₀ O ₅) ₁	365.1322	-2.2	GlcNAc(Man) ₁
527.1845	C ₈ H ₁₃ N O ₅ (C ₆ H ₁₀ O ₅) ₂	527.1850	-0.9	GlcNAc(Man) ₂
689.2364	C ₈ H ₁₃ N O ₅ (C ₆ H ₁₀ O ₅) ₃	689.2364	0.0	GlcNAc(Man) ₃
851.2887	C ₈ H ₁₃ N O ₅ (C ₆ H ₁₀ O ₅) ₄	851.2907	-2.3	GlcNAc(Man) ₄
1013.3367	C ₈ H ₁₃ N O ₅ (C ₆ H ₁₀ O ₅) ₅	1013.3435	-6.7	GlcNAc(Man) ₅
939.4143	C ₃₂ H ₄₈ N ₈ O ₁₂ (C ₈ H ₁₃ N O ₅)	939.4185	-4.5	GEVFNAT (GlcNAc)
1142.4954	C ₃₂ H ₄₈ N ₈ O ₁₂ (C ₈ H ₁₃ N O ₅) ₂	1142.4979	-2.2	GEVFNAT (GlcNAc) ₂
1304.5498	C ₃₂ H ₄₈ N ₈ O ₁₂ (C ₈ H ₁₃ N O ₅) ₂ (C ₆ H ₁₀ O ₅) ₁	1304.5507	-0.7	GEVFNAT (GlcNAc) ₂ (Man) ₁
1466.6024	C ₃₂ H ₄₈ N ₈ O ₁₂ (C ₈ H ₁₃ N O ₅) ₂ (C ₆ H ₁₀ O ₅) ₂	1466.6036	-0.8	GEVFNAT (GlcNAc) ₂ (Man) ₂
1628.6516	C ₃₂ H ₄₈ N ₈ O ₁₂ (C ₈ H ₁₃ N O ₅) ₂ (C ₆ H ₁₀ O ₅) ₃	1628.6564	-2.9	GEVFNAT (GlcNAc) ₂ (Man) ₃
1790.7051	C ₃₂ H ₄₈ N ₈ O ₁₂ (C ₈ H ₁₃ N O ₅) ₂ (C ₆ H ₁₀ O ₅) ₄	1790.7092	-2.3	GEVFNAT (GlcNAc) ₂ (Man) ₄
1952.7573	C ₃₂ H ₄₈ N ₈ O ₁₂ (C ₈ H ₁₃ N O ₅) ₂ (C ₆ H ₁₀ O ₅) ₅	1952.7620	-2.4	GEVFNAT (GlcNAc) ₂ (Man) ₅

In the pseudo MS3 experiment glycans were lost by in-source decay. GEVFNAT-GlcNAc was isolated in the quadrupole and fragmented in the collision cell. Sequence confirmation for the peptide stump GEVFNAT-GlcNAc is shown in Figure 5 with mass errors calculated in Table 4.

Figure 5. Pseudo MS3 fragmentation analysis of RBD glycopeptide stump GEVFNAT-GlcNAc (N343)

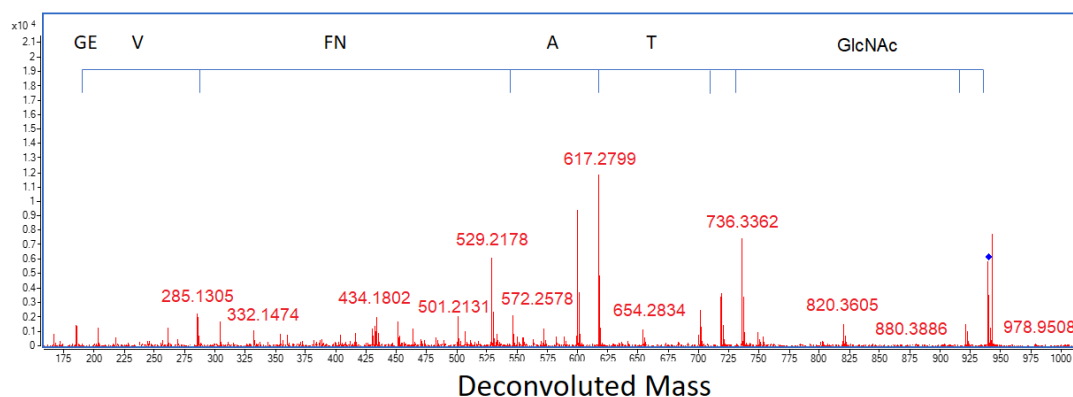
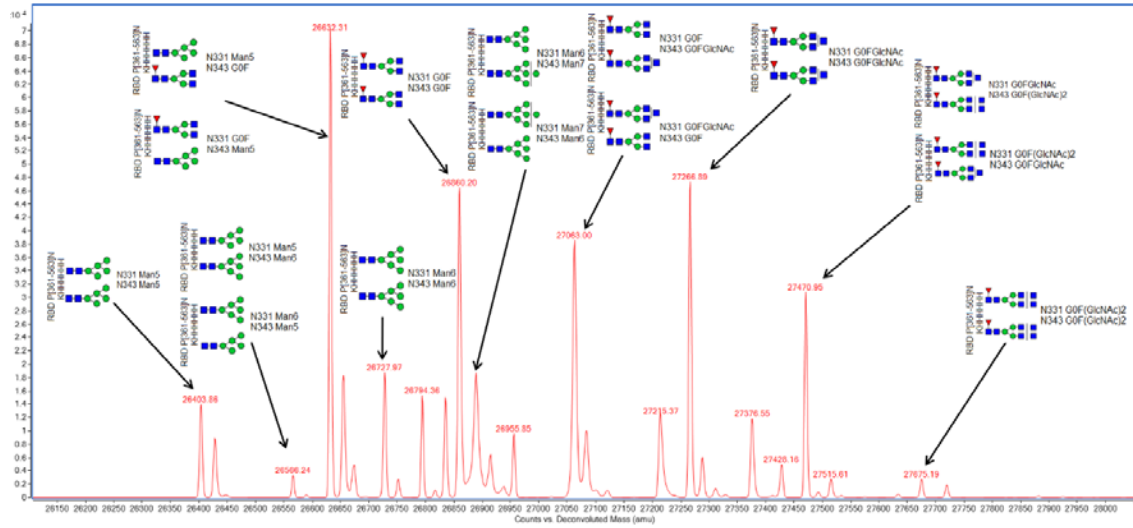


Table 4. Pseudo MS3 fragment ion assignment and mass errors (parts per million) for RBD glycopeptide stump GEVFNAT-GlcNAc (N343)

Deconvoluted Observed Mass	Formula	Calculated Mass	Error (ppm)	Fragment ion assignment	Peptide Sequence
186.0645	C ₇ H ₁₀ N ₂ O ₄	186.0641	2.1	b2	GE
285.1305	C ₁₂ H ₁₉ N ₃ O ₅	285.1325	-7.0	b3	GEV
546.2410	C ₂₅ H ₃₄ N ₆ O ₈	546.2438	-5.1	b5	GEVFN
617.2799	C ₂₈ H ₃₉ N ₇ O ₉	617.2809	-1.6	b6	GEVFNAT
718.3257	C ₃₂ H ₄₆ N ₈ O ₁₁	718.3286	-4.0	b7	GEVFNAT
736.3362	C ₃₂ H ₄₈ N ₈ O ₁₂	736.3392	-4.1	y7	GEVFNAT
939.4174	C ₃₂ H ₄₈ N ₈ O ₁₂ C ₈ H ₁₃ N O ₅	939.4185	-1.2	M GlcNAc	GEVFNAT GlcNAc
921.4064	C ₃₂ H ₄₆ N ₈ O ₁₁ C ₈ H ₁₃ N O ₅	921.408	-1.7	M GlcNAc - H ₂ O	GEVFNAT GlcNAc
820.3605	C ₂₈ H ₃₉ N ₇ O ₉ C ₈ H ₁₃ N O ₅	820.3603	0.2	b6 GlcNAc - H ₂ O	GEVFNAT GlcNAc
749.3213	C ₂₅ H ₃₄ N ₆ O ₈ C ₈ H ₁₃ N O ₅	749.3232	-2.5	b5 GlcNAc - H ₂ O	GEVFNAT GlcNAc
701.3021	C ₃₂ H ₄₃ N ₇ O ₁₁	701.3021	0.0	b7 - NH ₃	GEVFNAT
600.2529	C ₂₈ H ₃₆ N ₆ O ₉	600.2544	-2.5	b6 - NH ₃	GEVFNAT
529.2178	C ₂₅ H ₃₁ N ₅ O ₈	529.2173	0.9	b5 - NH ₃	GEVFNAT

Intact mass measurement of fully glycosylated Spike was unsuccessful due to the polydispersity of its innumerable glycoforms and the resulting dilution of ion signal. However, the smaller receptor binding domain, bearing only two glycosylation sites did prove amenable to intact mass analysis. Figure 6 shows twenty-one glycoforms for intact RBD, of which ten major glycoforms could be assigned. This showed that the principal glycan species were Man5, G0F and G0F+GlcNAc which was in agreement with the glycopeptide analysis.

Figure 6. Intact mass analysis of RBD showing the principal glycan species Man5, G0F and G0F+GlcNAc in agreement with glycopeptide analysis. (This method cannot differentiate individual glycosylation sites, hence when two structures are possible, both are shown)



Elastase was chosen as a single digestion enzyme because it was judged to give the best chance of generating glycopeptides with a single NXS/T motif, essential for unambiguous glycan mapping. For non-glycosylated Spike peptides, elastase generated 63 high quality MSMS hits and 26% coverage allowing for five missed cleavages. The same data searched for non-specific cleavage gave 135 high quality MSMS hits and 48% coverage allowing for twenty missed cleavages. Elastase itself contains 2 NXS/T motifs. We therefore prepared elastase only, at x10 the usual concentration, searched the resulting LC-MS data using the PCDL as a control, and no hits were found. The Spike protein LC-MS data did contain a small number of elastase autodigestion peptides.

Methods

Cloning, expression and purification of Spike

The gene encoding amino acids 1-1208 of the SARS-CoV-2 Spike glycoprotein ectodomain (S), with mutations of RRAR > GSAS at residues 682-685 (to remove the furin cleavage site) and KV > PP at residues 986-987 (to stabilise the protein), was synthesised with a C-terminal T4 fibrin trimerization domain, HRV 3C cleavage site, 8xHis tag, and Twin-Strep-tag [5]. The construct was sub-cloned into pHL-sec [10] using the AgeI and XhoI restriction sites and the sequence was confirmed by sequencing. Recombinant Spike was produced in *Expi293F*TM cells by transient transfection with purified DNA (0.5 mg/L cells) using a 1:6 DNA:L-PEI ratio, mixed in minimal medium, and sodium butyrate as an additive. Cells were grown in suspension in *FreeStyle293*TM medium with shaking at 150 rpm in 2 L smooth roller bottles, filled with 0.5 L cells at 2 e⁶/mL per bottle at 30°C with 8% CO₂ and 75% humidity. Supernatants from transfected cells were harvested 3-days post-transfection by centrifugation. Clarified supernatant was mixed with Ni²⁺ IMAC *Sepharose*[®] 6 Fast Flow (*GE*; 2 mL bed volume per L of supernatant) at room temperature for 2 h. Using a gravity flow column, resin

was collected and washed stringently with 50 CV each of base buffer (1X PBS), WB25 (BB + 25 mM imidazole), and WB40 (BB+ 40 mM imidazole), followed by elution with EB (0.30 M imidazole in 1X PBS). Protein was dialyzed into 1X PBS using *SnakeSkin*[™] 3,500 MWCO dialysis tubing, concentrated to 1 mg/mL using a 100,000 MWCO *VivaSpin* centrifugal concentrator (GE), and centrifuged at 21,000 x g for 30 min to remove aggregates. The trimeric Spike protein was flash frozen in LN₂ and stored at -80°C until use. Final purified yield was 1 mg of Spike protein per L of transfected cells.

Cloning, expression and purification of Receptor Binding Domain

The receptor binding domain (RBD; aa 330-532) of SARS-CoV-2 Spike (Genbank MN908947) was inserted into the pOPINTTNeo expression vector fused to an N-terminal signal peptide and a C-terminal 6xHis tag [11]. RBD was produced by transient transfection in *Expi293F*[™] cells (*ThermoFisher Scientific*, UK) using purified DNA (1.0 mg/L cells), a 1:3 DNA:L-PEI ratio, and sodium butyrate as an additive. Cells were grown in suspension in *FreeStyle293*[™] expression medium at 37°C with 8% CO₂ and 75% humidity. Supernatants from transfected cells were harvested 3-days post-transfection and the supernatant was collected by centrifugation. Clarified supernatant was incubated with 5 mL of Ni²⁺ IMAC *Sepharose*[®] 6 *Fast Flow* (GE) at room temperature for 2 h. Using gravity flow, resin was washed with 50 CV of base buffer (1X PBS) and 50 CV of WB (1X PBS + 25 mM imidazole) before elution with EB (0.5 M imidazole in 1X PBS). Protein was concentrated using a 10,000 MWCO *Amicon Ultra-15* before application to a *Superdex 75 16/600* column pre-equilibrated with 1X PBS pH 7.4. Peak monomeric fractions were pooled and concentrated to 2 mg/mL, flash frozen in LN₂, and stored at -80°C until use. Final purified yield was >15 mg RBD per L of transfected cells.

Sample preparation

SARS-CoV-2 Spike or RBD-6H at 1 mg/mL in PBS were prepared in aliquots of either 20 µL or 80 µL and diluted 1 in 3 in 100 mM ammonium bicarbonate, pH 8.0, followed by reduction by addition of 1, 4 Dithiothreitol (DTT) to 5 mM and incubation 37°C for 1 h. Next, the protein was alkylated by addition of iodoacetamide (IAA) to 15 mM and incubation in the dark for 30 min. This was followed by overnight digestion using elastase (*Promega*) at a ratio of 1:20 (w/w). The following day, the supernatant was dried using a rotary evaporator, and re-suspended in 60 µL of 0.1% formic acid for injection into the LC-MS.

'Analytical mode' LC-MS glycopeptide data acquisition

LC-MS 'analytical mode' was performed using a *1290 Infinity* UHPLC coupled to a *G6530A* ESI QTOF mass spectrometer (*Agilent Technologies*). TOF and quadrupole were calibrated prior to analysis and the reference ion 922.0098m/z was used for continuous mass correction. Sample was introduced using a 50 µL full-loop injection. Reversed phase chromatographic separation was achieved using an *AdvancedBio Peptide* reversed phase 2.7 µm particle, 2.1 mm x 100 mm column 655750-902 (*Agilent Technologies*). Mobile phase A was 0.1% formic acid in water and mobile phase B 0.1% formic acid in methanol (*Optima* LC-MS grade, *Fisher*). Initial conditions were 5% B and 0.200 mL/min flow rate. A linear gradient from 5% B - 60% B was applied over 60 min, followed by isocratic elution at 100% B for 2 min returning to initial conditions for a further 2 min. Post time was 10 min. MS source parameters were drying gas temperature 350°C, drying gas 8 L/min, nebulizer 30 psi, capillary 4000 V, fragmentor 150 V. MS spectrum range was 100 – 3200 m/z (centroid only), 2 GHz Extended Dynamic range, with the instrument in positive ion mode.

LC-MSMS glycopeptide data acquisition 'discovery mode'

LC-MSMS 'discovery mode' was performed as described above, with the following changes: Soft CID collision energy parameters for MSMS were slope 1.0, intercept 0 using argon as the collision gas (if using nitrogen slope 2.0, intercept 0) were used to favour glycan fragmentation over peptide

and charge states +2 to +5. The results were filtered to remove compounds <1000 Da (too small to be glycopeptides). Compound MSMS spectra were screened manually for the following oxonium reporter ions: Hex m/z 163.0601, HexNAc m/z 204.0866, HexHexNAc m/z 366.1395, Neu5Ac m/z 274.0921/ m/z 291.0949 and/or a Hexose ladder m/z 162.0528 Da. High quality m/z spectra were deconvoluted to neutral mass spectra with glycan *de novo* interpretation performed manually. Once a glycopeptide had been identified, it was entered into a personal compound data library database (PCDL, *Agilent Technologies*) as a mass and retention time. In addition, the database made use of known mammalian N-linked glycan processing. After the initial glycopeptide identification, other processed glycopeptides, which were considered likely to also be present, were added to the database at the same retention time and with a calculated mass. For example, if a glycopeptide with Man5 was identified by MSMS, Man1-9 and G0/F were added at the same retention time. If these glycans were subsequently found in the data, their actual retention times were updated, and the next round of processing to more complex glycans was added, in order to produce the most comprehensive PCDL possible, while still being manageable. Processing order:

Man(n) → G0/F → G1/F → G2/F → A1/F → A2/F → Very Complex

Valid glycan identifications resulted in a calculated peptide mass that could be matched to the sequence. Where high quality spectra were present, a peptide-GlcNAc stump was observed (Figure 4). This was used in a pseudo MS3 experiment with manual peptide *de novo* interpretation to confirm the peptide sequence (Figure 5). Mass data adjacent to the glycopeptide retention time was then searched for neutral differences corresponding to glycans, for example, Man5 → G0F or Man7 → G2F has a neutral delta mass of 228.1111 Da.

As expected, not all species could be matched to the sequence, presumably due to unexpected modifications. In this case, they were added to the database as 'GP' with an identifying number and as much information as could be extracted. Data for the most likely glycan was added to the PCDL, including a deconvoluted mass MSMS spectra were available, using nomenclature generating the most easily readable format.

A second round of glycopeptide discovery used *Bioconfirm* v10.0 data analysis software (*Agilent Technologies*). Sequences were matched by peptide accurate mass using the following parameters: peptide cleavage nonspecific, number of missed cleavages 20, N-linked modifications Man3, Man5-9, G0, G0F, G0F GlcNAc, G1, G1F, G2, G2F. Any peptide bearing the glycosylation motif NXS/T with two or more glycan hits within a retention time window +/-2 min was added to the PCDL, excepting missed cysteine alkylations.

In-source fragmentation due to glycopeptide ions absorbing excess energy could be identified in the MS by searching extracted ion chromatograms (EICs) of the oxonium reporter ions and also by related glycopeptides appearing with exactly the same retention times. Both were observed infrequently and at manageable levels.

Intact mass analysis

Concentrated protein samples were diluted to 0.02 mg/mL in 0.1% formic acid and 50 μ L was injected on to a 2.1 mm x 12.5 mm *Zorbax* 5 μ m 300SB-C3 guard column (*Agilent Technologies*) housed in a column oven set at 40°C. The solvent system used consisted of 0.1% formic acid (solvent A) and 0.1% formic acid in methanol (solvent B). Chromatography was performed as follows: Initial conditions were 90% A and 10% B and a flow rate of 1.0 mL/min. A linear gradient from 10% B to 80% B was applied over 35 seconds. Elution then proceeded isocratically at 95% B for 40 seconds followed by equilibration at initial conditions for a further 15 seconds. The mass spectrometer was configured with the standard ESI source and operated in positive ion mode. The ion source was operated with the capillary voltage at 4000 V, nebulizer pressure at 60 psig, drying gas at 350°C and

drying gas flow rate at 12 L/min. The instrument ion optic voltages were as follows: fragmentor 250 V, skimmer 60 V and octopole RF 250 V.

Discussion

Glycoprotein analysis is difficult. It is either performed in biopharmaceutical laboratories with proprietary expertise of glycan analysis on simple glycoproteins, such as immunoglobulins, or performed by a handful of academic labs with experience of glycan discovery from complex glycoproteins. Many protein researchers choose to ignore it, manipulating cell lines such that they cannot process beyond Man5, or to remove glycans entirely by mutation at the glycosylation motif or enzymatically [12]. While this approach has its merits, it has exposed a serious weakness in analytical capability when faced with a pathogen such as SARS-CoV-2 whose ability to evade the immune system is dependent upon heavy and complex glycosylation.

We have chosen an approach relying on elastase digestion to generate glycopeptides bearing a single glycan but with a sufficient number of amino acid residues to enable chromatographic separation by reversed-phase HPLC, as well as confident identification by accurate mass or *de novo* sequencing. Our choice of reversed phase HPLC has excellent discrimination for short elastase peptides, whereas glycans show little or no interaction with the column. Thus, species originating from a single glycosylation site with the same peptide sequence but several different glycans, eluted with the same retention time and could be discriminated by mass spectrometry. We used reversed phase HPLC and MSMS to characterise as many glycopeptides as possible. Although this required complex and time-consuming data analysis, it needed only be performed once, with the goal of building an accurate mass-retention time database for all observed Spike glycopeptides. Provided the same HPLC column and mobile phase conditions are used, retention times should not vary significantly. Thus, working in the analytical mode we describe, glycan structure and peptide sequence is assigned confidently, by accurate mass and retention time alone. LC-MS data need only to be searched against the mass-retention time database, and peak areas recorded, to generate a complete characterisation of Spike glycans.

We believe the MRTF method described here has advantages over other approaches to Spike glycan analysis. Previous studies relied upon very expensive equipment and software unavailable in most analytical laboratories. Working in 'analytical' mode, all that is necessary is to reproduce the chromatography, hence our method is a generic one, which can be run using any HPLC coupled to any accurate mass instrument and is not restricted to specific proprietary data analysis software. We used PCDL and *Masshunter*, but MRTF analysis can be performed on any vendor software or manually. Moreover, it demands no specialised expertise in glycobiology, and is thus accessible to many more researchers. Some published methods require multiple specific endoproteases, some of which cannot be readily sourced. Our method uses a single enzyme, elastase, which is inexpensive and widely available. Nor does it rely on glycosidases, which may not work efficiently and do not cleave O-linked glycans.

Our data contains an excess of glycopeptides with the motif (y)nNxS/T. This appears to be a very convenient function of elastase on glycopeptides, because the presence of the motif at the C-terminus facilitates *de novo* sequencing. We would be interested to know if this cleavage bias towards the C-terminus of the glycan motif is reproducible in other labs and whether it indicates steric hindrance within the elastase enzyme structure. If such bias is real, then these peptides are less likely to be a false positive result.

Receptor binding domain (RBD) from Spike protein is of interest in many labs for development of serological tests or neutralising antibodies. Because the yield of RBD was five times higher than Spike and more was initially available, we used it for method optimisation, and since it bears only two glycosylation sites which are also present on Spike, it functioned as a useful model. Consequently, N343 on glycopeptide GEVFNAT is over-represented in our demonstration PCDL. We consistently

observed the same three major glycans (Man5, G0F and G0F+GluNAc) on this peptide and these were also in agreement with intact mass analysis of RBD protein as shown in Figure 6. On closer inspection, glycans up to A2F could also be observed at lower levels. We suspect that sufficiently detailed analysis may reveal all possible glycan structures with low abundance at all available sites. The most important would therefore be the top three to five glycans. If the complete complement of Spike protein glycoforms proves too challenging for a single analysis, this site, which is the most complete, would make a good proxy for total glycosylation.

We acknowledge that the mass-retention time fingerprinting method described, like all database searching methods, is dependent on the reproducibility of the enzyme digestion and both the quality and the completeness of database being searched. The example PCDL database reported here is provided as a demonstration. Due to glycan complexity and the likely absence of specific glycans within the Spike batches prepared by us, it will always be incomplete. Moreover, individual glycopeptides were identified with variable degrees of certainty, and we recommend that they should be validated by the user. As with all glycan analysis methods, there is a bias towards glycopeptides that are easiest to identify by the techniques used, and such bias will also be reflected within the database. Once the PDCL has been created, it must be refined and extended over time to improve data quality, and it is our intention to do so.

References

1. Zhang, Y., et al., *Site-specific N-glycosylation Characterization of Recombinant SARS-CoV-2 Spike Proteins using High-Resolution Mass Spectrometry*. bioRxiv, 2020.
2. Solá, R.J. and K. Griebenow, *Glycosylation of therapeutic proteins*. BioDrugs, 2010. **24**(1): p. 9-21.
3. Vugmeyster, Y., et al., *Pharmacokinetics and toxicology of therapeutic proteins: advances and challenges*. World journal of biological chemistry, 2012. **3**(4): p. 73.
4. Tortorici, M.A. and D. Veessler, *Structural insights into coronavirus entry*. Advances in virus research, 2019. **105**: p. 93-116.
5. Wrapp, D., et al., *Cryo-EM structure of the 2019-nCoV spike in the prefusion conformation*. Science, 2020. **367**(6483): p. 1260-1263.
6. Yang, T.-J., et al., *Cryo-EM analysis of a feline coronavirus spike protein reveals a unique structure and camouflaging glycans*. Proceedings of the National Academy of Sciences, 2020.
7. Shajahan, A., et al., *Deducing the N-and O-glycosylation profile of the spike protein of novel coronavirus SARS-CoV-2*. bioRxiv, 2020.
8. Watanabe, Y., et al., *Site-specific analysis of the SARS-CoV-2 glycan shield*. BioRxiv, 2020.
9. Yao, H., et al., *Molecular architecture of the SARS-CoV-2 virus*. bioRxiv, 2020: p. 2020.07.08.192104.
10. Aricescu, A.R., W. Lu, and E.Y. Jones, *A time-and cost-efficient system for high-level protein production in mammalian cells*. Acta Crystallographica Section D: Biological Crystallography, 2006. **62**(10): p. 1243-1250.
11. Nettleship, J.E., et al., *Transient expression in HEK 293 cells: an alternative to E. coli for the production of secreted and intracellular mammalian proteins*, in *Insoluble proteins*. 2015, Springer. p. 209-222.
12. Chang, V.T., et al., *Glycoprotein structural genomics: solving the glycosylation problem*. Structure, 2007. **15**(3): p. 267-273.

Acknowledgements

The authors wish to thank Professor David Harvey for critical reading of the manuscript and for helpful comments. We thank Professor Ray Owens for kindly providing the RBD-6H construct and Professor Gavin Screaton and Dr. Juthathip Mongkolsapaya for kindly providing the Spike construct. AdvancedBio Peptide HPLC column was a gift from Agilent Technologies. This project has received

funding from the Innovative Medicines Initiative 2 Joint Undertaking (JU) under grant agreement No 875510. The JU receives support from the European Union's Horizon 2020 research and innovation programme and EFPIA and Ontario Institute for Cancer Research, Royal Institution for the Advancement of Learning McGill University, Kungliga Tekniska Hoegskolan, Diamond Light Source Limited. The SGC is a registered charity (number 1097737) that receives funds from AbbVie, Bayer Pharma AG, Boehringer Ingelheim, Canada Foundation for Innovation, Eshelman Institute for Innovation, Genentech, Janssen, Merck KGaA, Darmstadt, Germany, MSD, Ontario Ministry of Research, Innovation and Science (MRIS), Pfizer, São Paulo Research Foundation-FAPESP, Takeda, and Wellcome [106169/ZZ14/Z].