

Supplemental Data

Enhancing Mass spectrometry-based tumor immunopeptide identification: machine learning filter leveraging HLA binding affinity, aliphatic index and retention time deviation

Feifei Wei, Taku Kouro, Yuko Nakamura, Hiroki Ueda, Susumu Iizumi, Kyoko Hasegawa, Yuki Asahina, Takeshi Kishida, Soichiro Morinaga, Hidetomo Himuro, Shun Horaguchi, Kayoko Tsuji, Yasunobu Mano, Norihiro Nakamura, Takeshi Kawamura, and Tetsuro Sasada

Contents

Table S1. Patient list

Table S2. CandiSeq list of training set

Table S3. CandiSeq list of test set

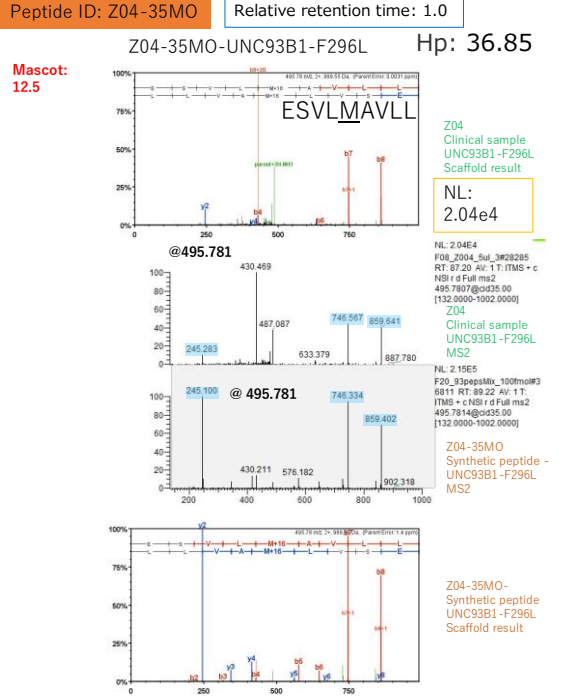
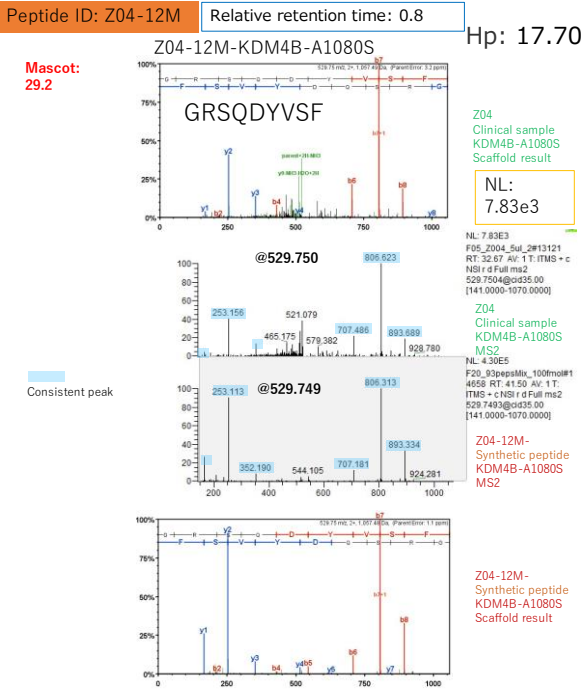
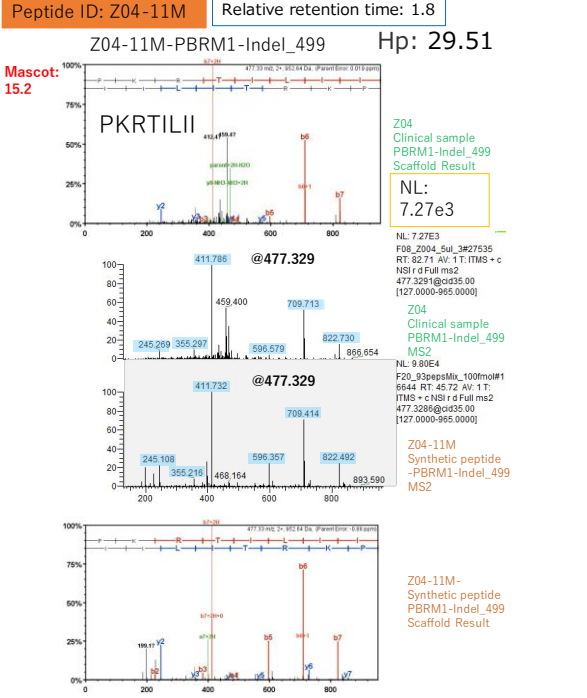
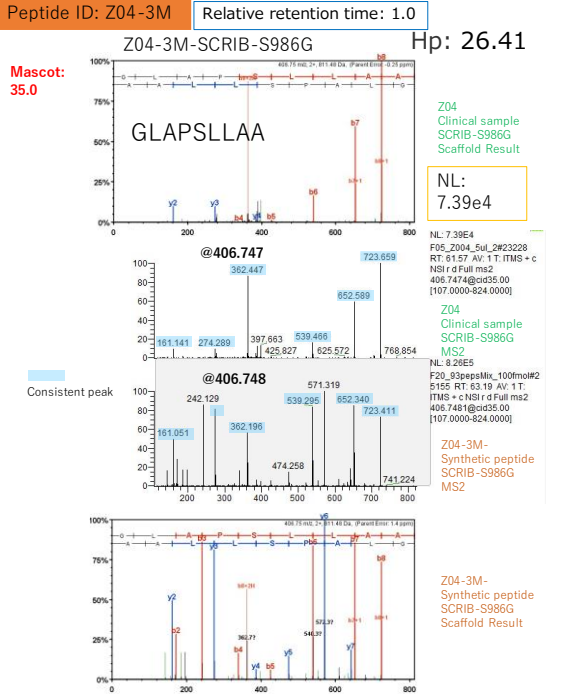
Table S4. List of overlapping peptides with more than 8 amino acid residues in the current study dataset (combination of training and test sets)

Table S5. Results of cut-off analysis

Figure S1. MS/MS spectra of HCS

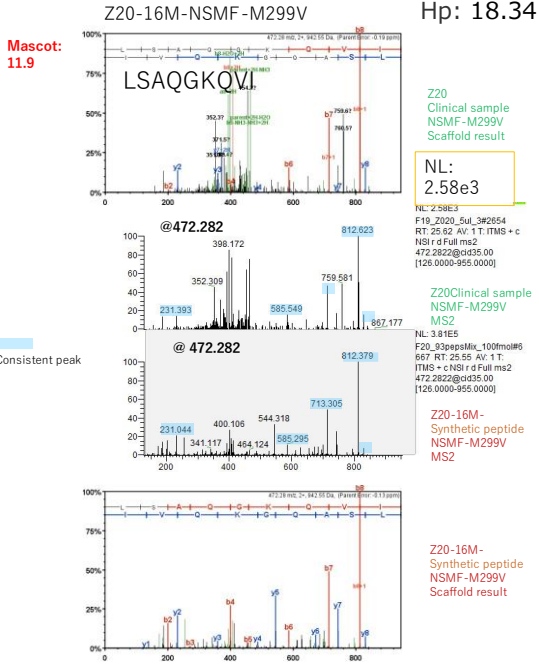
Figure S2. Selection of Hydrophobicity for predicting standard retention time

Figure S3. Count of overall detected peptides from patients with pancreas and kidney cancer



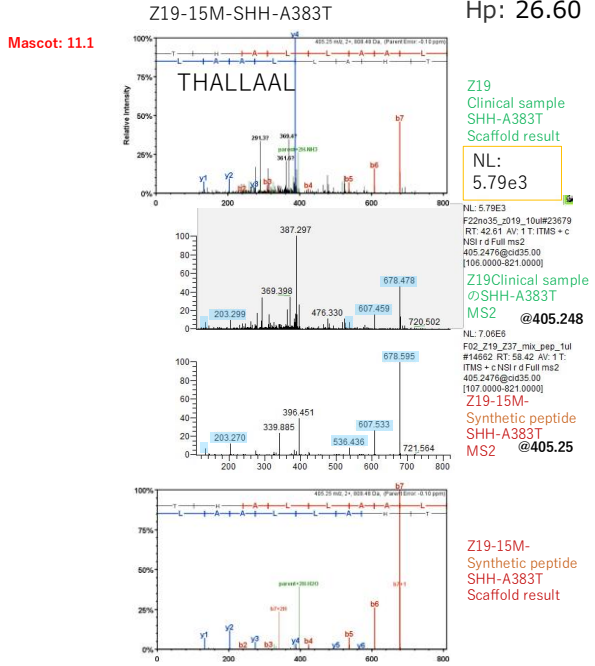
Peptide ID: Z20-16M

Relative retention time: 1.0



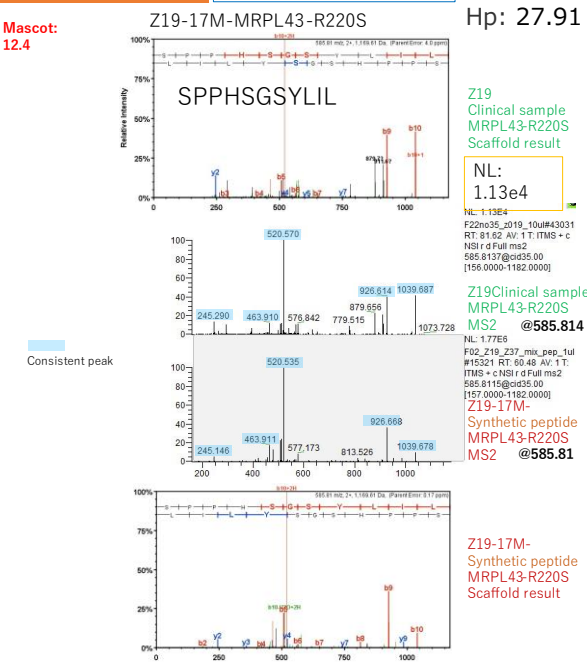
Peptide ID: Z19-15M

Relative retention time: 0.7



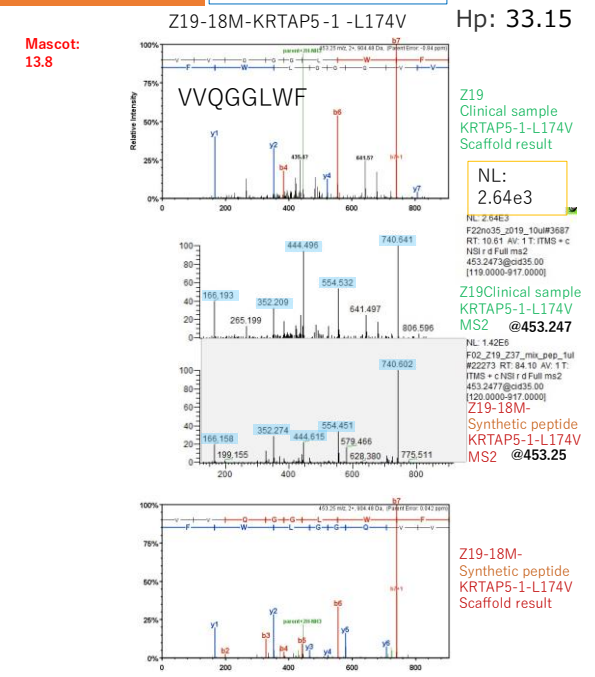
Peptide ID: Z19-17M

Relative retention time: 1.3



Peptide ID: Z19-18M

Relative retention time: 0.1



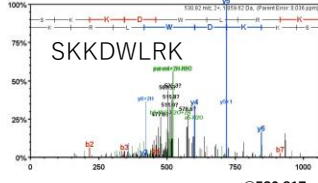
Peptide ID: Z34-2M

Relative retention time: 1.3

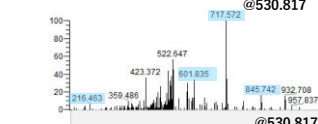
Z34-2M-FAM221A-Indel_167

Hp: 16.68

Mascot: 14.9

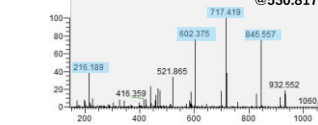


Z34
Clinical sample
FAM221A-
Indel_167
Scaffold result
NL:
1.79e3

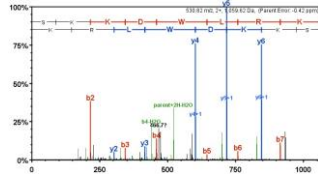


Z34Clinical sample
のFAM221A-
Indel_167 MS2

Consistent peak



Z34-2M-
Synthetic peptide
FAM221A-Indel_167
MS2



Z34-2M-
Synthetic peptide
FAM221A-
Indel_167
Scaffold result

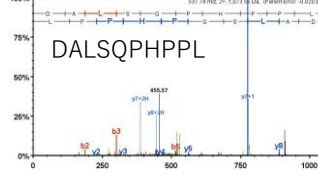
Peptide ID: Z35-13M

Relative retention time: 1.3

Z35-13M-NYAP1-A761P

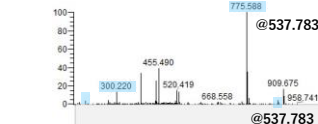
Hp: 21.58

Mascot: 14.5

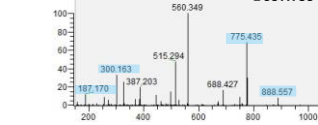


Z35
Clinical sample
NYAP1-A761P
Scaffold result

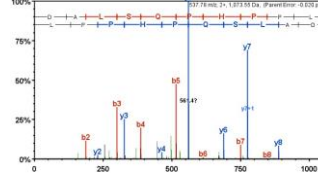
NL:
2.15e4



Z35
Clinical sample
NYAP1-A761P
MS2



Z35-13M-
Synthetic peptide
NYAP1-A761P
MS2



Z35-13M-
Synthetic peptide
NYAP1-A761P
Scaffold result

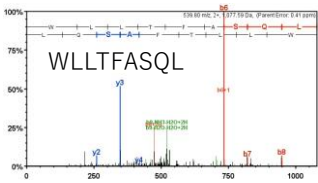
Peptide ID: Z35-16M

Relative retention time: 0.9

Z35-16M-SLC12A4-R712W

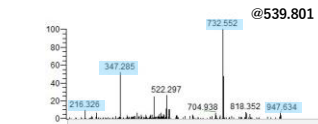
Hp: 41.82

Mascot: 11.5



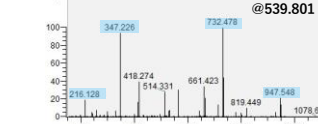
Z35
Clinical sample
SLC12A4-R712W
Scaffold result

NL:
1.02e4

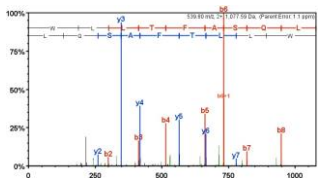


Z35Clinical sample
のSLC12A4R712W
MS2

Consistent peak



Z35-16M-
Synthetic peptide
SLC12A4-R712W
MS2



Z35-16M-
Synthetic peptide
SLC12A4-R712W
Scaffold result

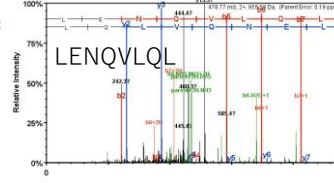
Peptide ID: Z37-2M

Relative retention time: 1.2

Z37-2M-GPR64-M360L

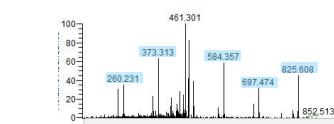
Hp: 26.64

Mascot: 19.3

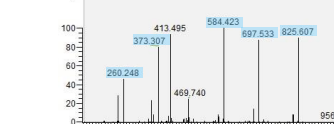


Z37
Clinical sample
GPR64-M360L
Scaffold result

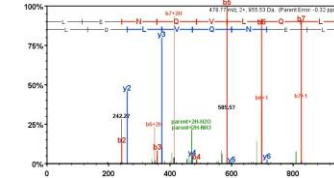
NL:
2.58e3



Z37
Clinical sample
GPR64-M360L MS2



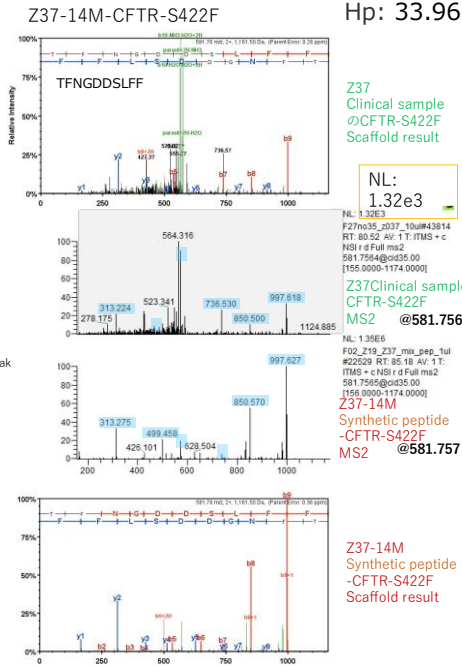
Z37-2M
Synthetic peptide -
GPR64-M360L
MS2 @478.774



Z37-2M-
Synthetic peptide
GPR64-M360L
Scaffold result

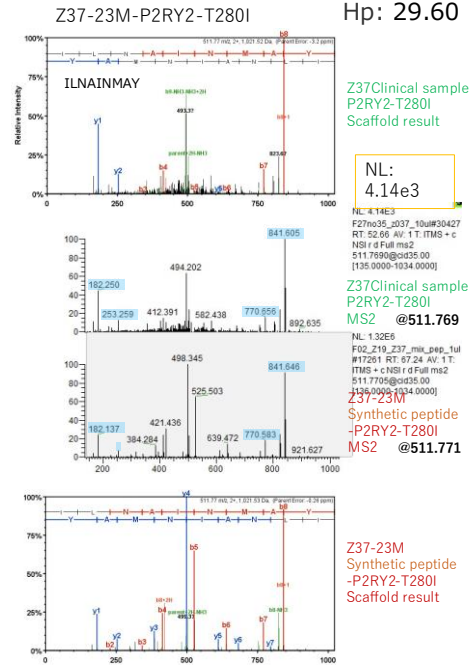
Peptide ID: Z37-14M Relative retention time: 0.9

Mascot: 11.7



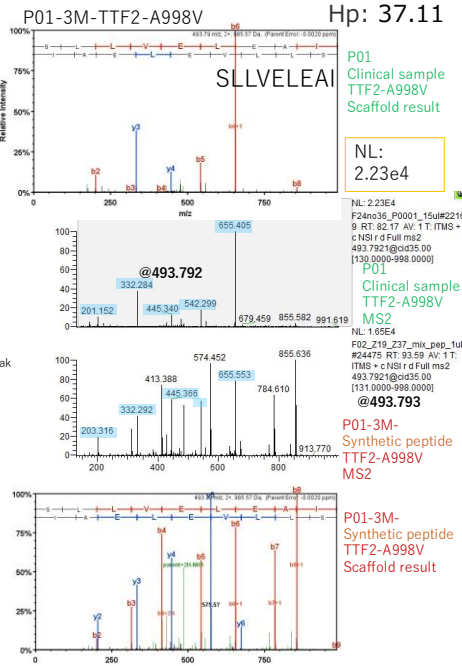
Peptide ID: Z37-23M Relative retention time: 0.8

Mascot: 11.2



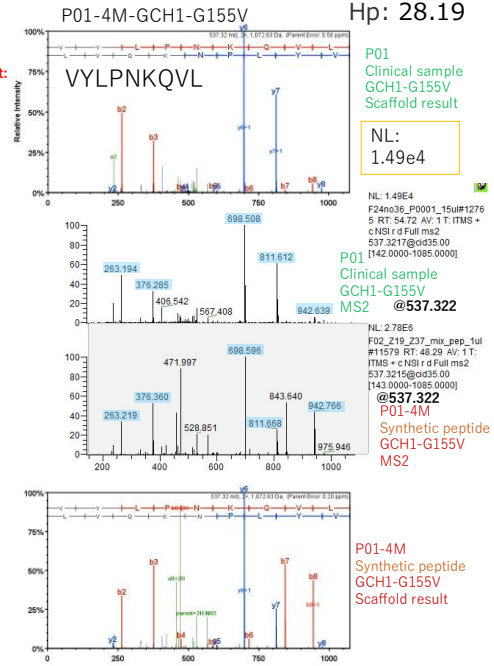
Peptide ID: P01-3M Relative retention time: 0.9

Mascot: 28.6



Peptide ID: P01-4M Relative retention time: 1.1

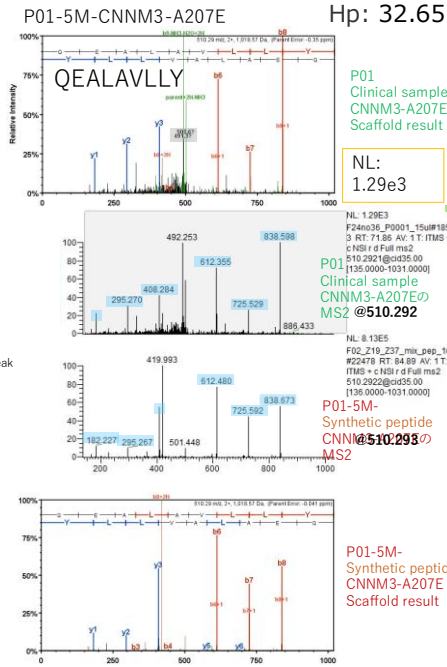
Mascot: 25.6



Peptide ID: P01-5M

Relative retention time: 0.8

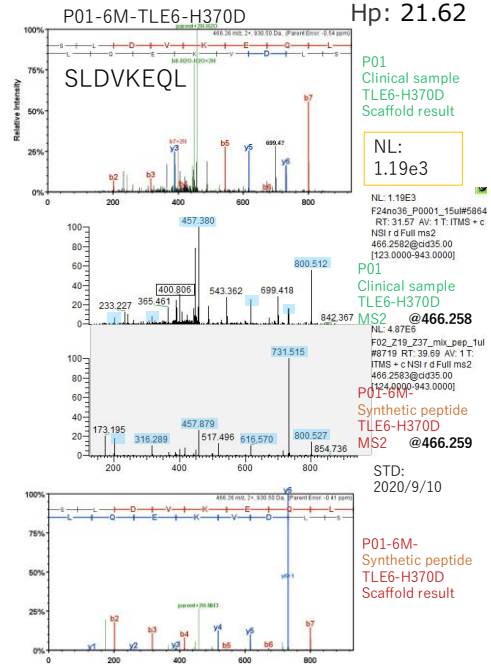
Mascot: 21.3



Peptide ID: P01-6M

Relative retention time: 0.8

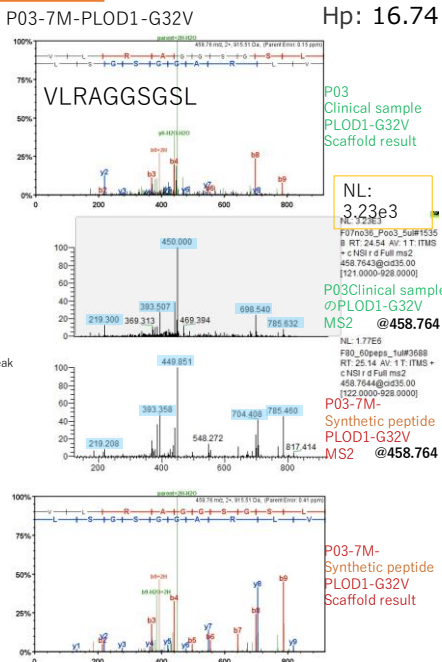
Mascot: 24.1



Peptide ID: P03-7M

Relative retention time: 1.0

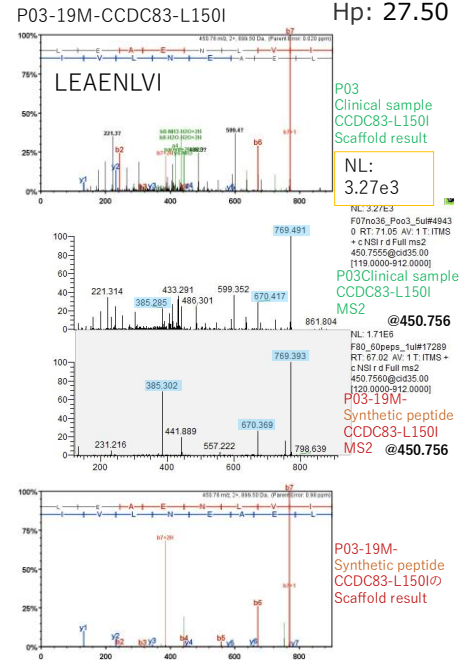
Mascot: 19.1



Peptide ID: P03-19M

Relative retention time: 1.1

Mascot: 15.5



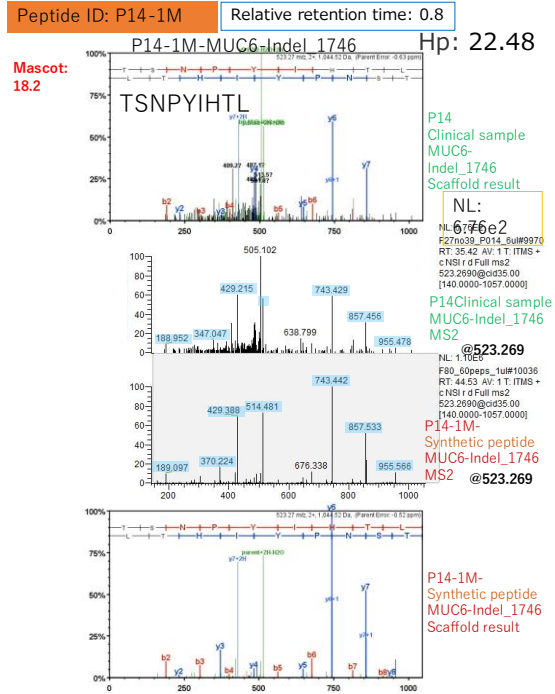
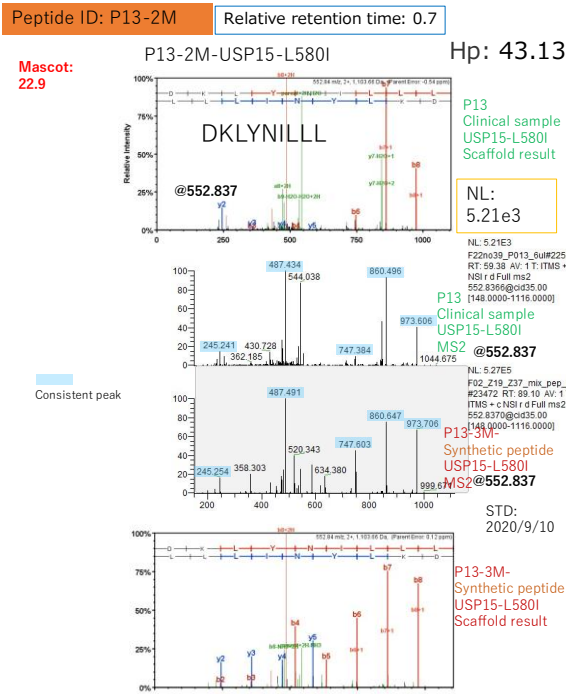
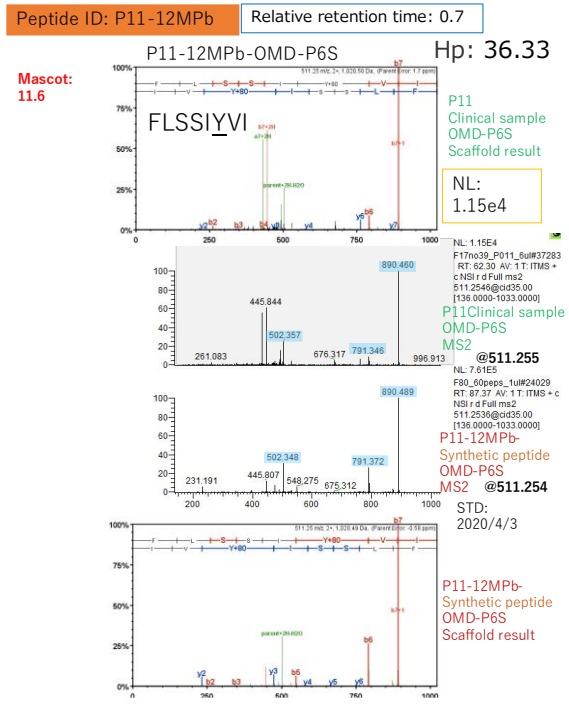
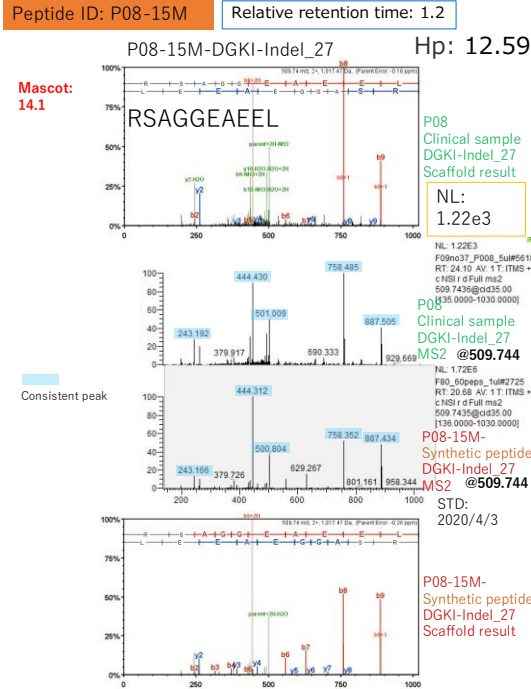


Figure S1. MS/MS spectra of HCS. HCS, reliable identifications with highly consistent MS/MS spectra.

a

x	y	R-squared	F-statistic	p-value
Hydrophobicity	SynRT_HSC	0.88	164.9	1.07E-11***
Gravy	SynRT_HSC	0.55	26.89	3.37E-05***
pl	SynRT_HSC	0.35	11.61	0.002523**
AliphaticIndex	SynRT_HSC	0.30	9.521	0.005403**
Entrp	SynRT_HSC	0.05	1.236	0.2783
InstabilityIndex	SynRT_HSC	0.05	1.175	0.2901
MinRank	SynRT_HSC	0.02	0.4739	0.4984

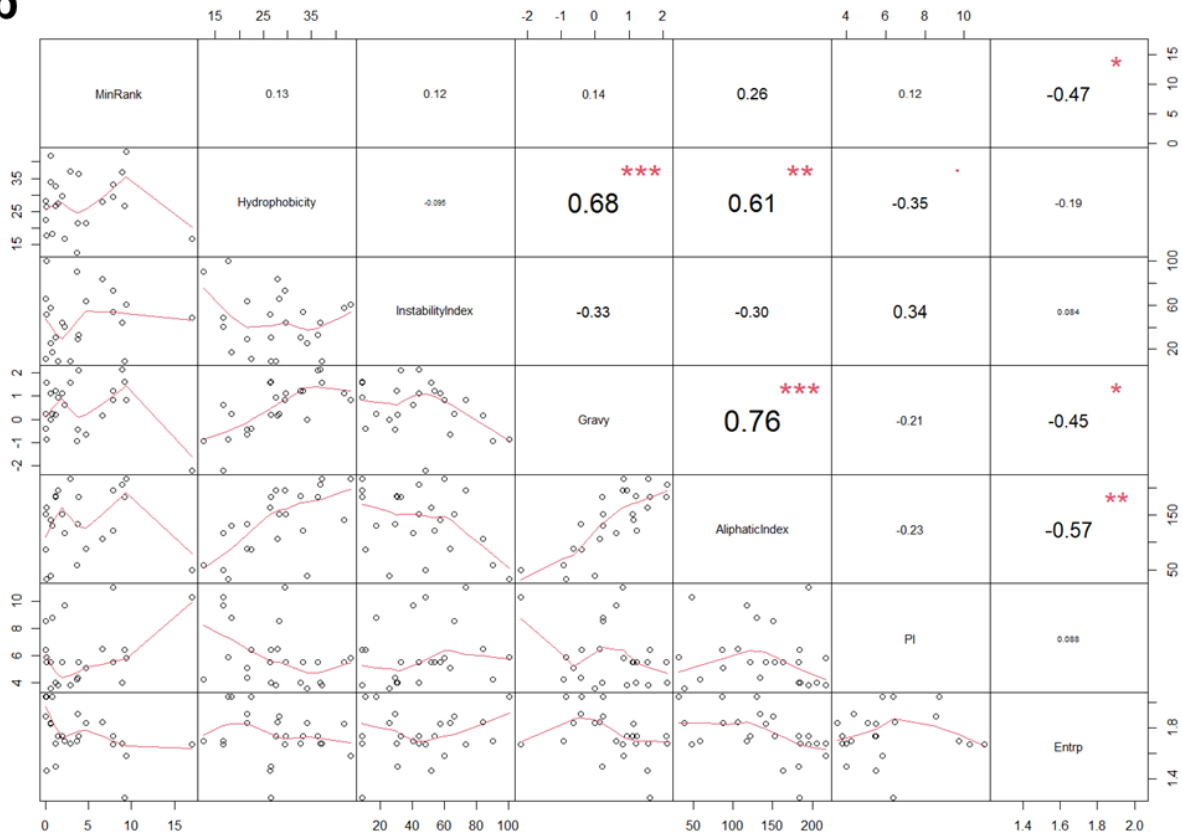
b

Figure S2. Selection of Hydrophobicity for predicting standard retention time: **a** outcome of univariate linear regression with the standard retention time (SynRT_HSC, y) for features predicted based on sequence information (x) across the 24 sequences exhibiting highly consistent MS/MS spectra (HCS). The features are ranked based on the R-squared value of the regression

line, and the linearity between Hydrophobicity and SynRT_HCS is optimal; **b** results of Spearman correlation analysis between features. Spearman coefficients plus the significance level were shown in the upper triangular, and the bivariate scatter plots with a fitted line were displayed in the lower part below the diagonal. Features exhibiting considerable linearity with SynRT_HCS (Gravy, pI and AliphaticIndex) display significant or strong correlations with Hydrophobicity. Considering the limited sample size of HCS and the risk of overfitting caused by multicollinearity in the multivariable regression model, the univariable linear regression model (as shown in **Figure 2**) using Hydrophobicity was selected to give a relatively robust performance. Significance level: *** P-value 0 - 0.001; ** P-value 0.001 - 0.01; * P-value 0.01 - 0.05; • P-value 0.05 - 0.10.

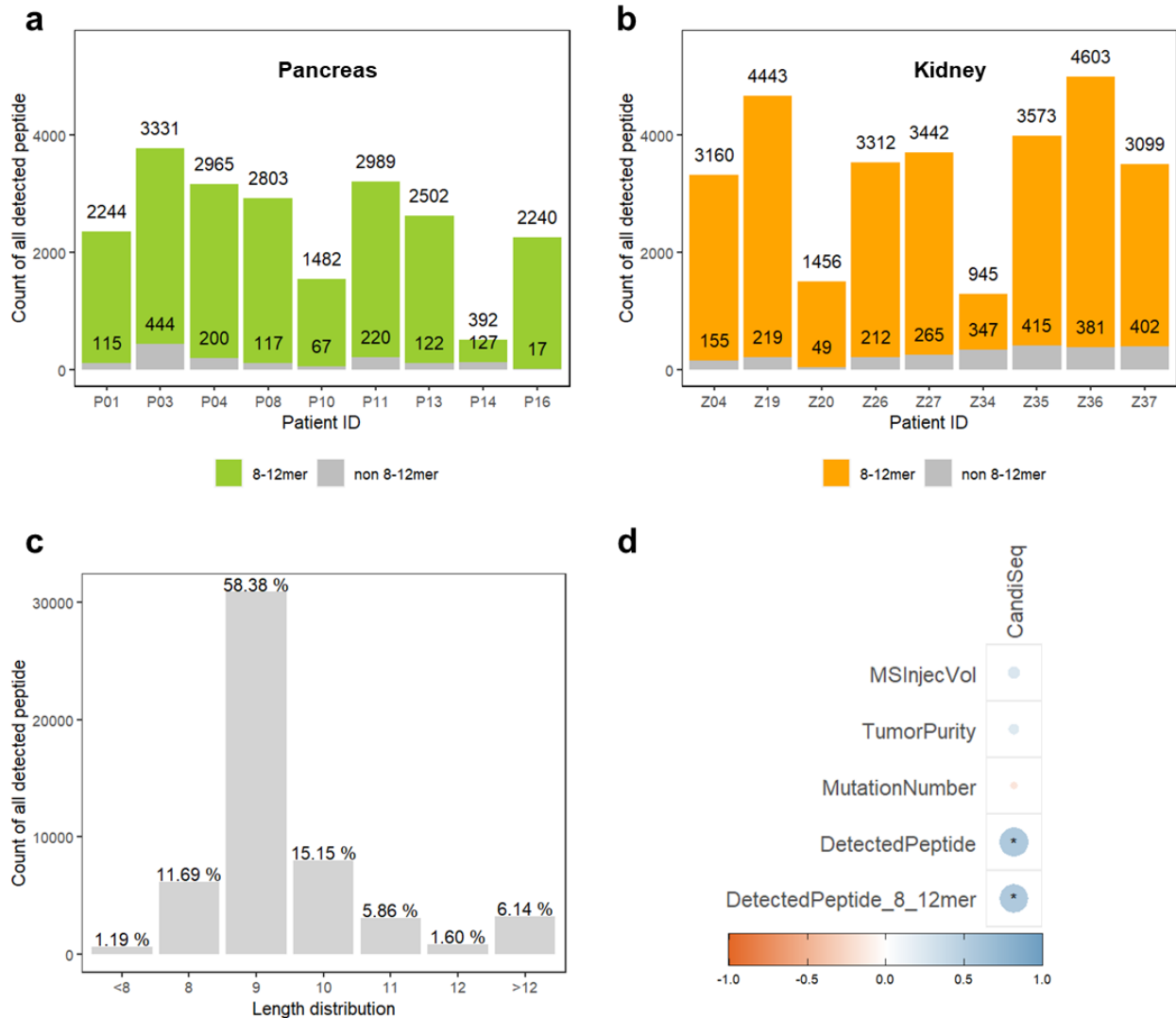


Figure S3. Count of overall detected peptides (Mascot ≥ 10) from patients with **a** pancreas and **b** kidney cancer; **c** the length distribution of all detected peptides; **d** Spearman's correlation coefficients between count of CandiSeqs (CandiSeq), tumor tissue input (MSInjecVol), tumor purity (TumorPurity), mutation number (MutationNumber), count of overall detected peptides (DetectPeptide), and count of 8-12mers in overall detected peptides (DetectedPeptide_8_12mer). Number labels in **a** and **b**, count number; number labels in **c**, percentage in total detected peptides; circle size in **d**, |Spearman's correlation coefficient|; circle color, red: negative, blue: positive value; * p-value < 0.05 , ** p-value < 0.01 , *** p-value < 0.001 .