

## Supplementary Information

### **High-throughput identification of FLT3 wild-type and mutant kinase substrate preferences and application to design of sensitive *in vitro* kinase assay substrates**

Minervo Perez<sup>1,2</sup>, John Blankenhorn<sup>1</sup>, Kevin J. Murray<sup>3</sup>, Andy W. Tao<sup>2</sup>, Laurie L. Parker<sup>1</sup>

Institutional Affiliations:

<sup>1</sup>University of Minnesota, Department of Biochemistry, Molecular Biology and Biophysics, 420 Washington Avenue SE, Minneapolis, MN 55455

<sup>2</sup>Purdue University, Department of Medicinal Chemistry and Molecular Pharmacology, 201 S. University Street, West Lafayette, IN 47907

<sup>3</sup>University of Minnesota, Department of Veterinary Population Medicine, 319 15<sup>th</sup> Avenue South East, Minneapolis, MN 55455

Figures:

- Figure S1. Schematic representation of raw mass spectrometer file combination for ProteinPilot database searches
- Figure S2. Site selectivity matrices for 16H-KALIP FLT3 kinase variants
- Figure S3. Heat map representation of FLT3-WT time course KALIP experiment, Site Selectivity Matrix and artificial substrate library sequence scoring comparison

Tables

- Table S1. A summary of the FLT3 Artificial Substrate (FAS) candidate sequences synthesized and assayed *in vitro* with recombinant FLT3 variants.
- Table S2. Performance metrics and comparison of PSM models.

## Methods

*Extraction and reformatting of phosphopeptide sequences from peptide ID results*—The KinaMINE data formatter (Kinamine.jar) uses the Distinct Peptide Report and the FASTA file that was used in the proteomics search engine as input, filters the peptides from the report with a threshold of 1% FDR, to consolidate the sequences of all peptides that were phosphorylated in the experiment. It then outputs a .csv table (the “Positive Substrates” file, which is named by the user at the time of running the script) of those tyrosine-phosphorylated sequences, with each amino acid separated into an individual column and the phosphotyrosine aligned. This table also contains the accession number of the protein each peptide was from, which is used to extract the sequences of those proteins from the inputted FASTA file and calculate the “Substrate Background Frequency” (frequency of the 20 canonical amino acids found in each of the proteins individually; SBF), also output as a .csv. This .csv file also reports the total number of tyrosine residues within those protein sequences and the number of those tyrosine residues that were observed as phosphorylated in the experiment for subsequent use in determining FLT3’s “normalization score” in the Screener module of KINATEST-ID (described below).

*Phosphopeptide list comparison filtering*—To select the sequences that were phosphorylated in common between the WT and the two mutant forms of FLT3, we developed a filtering script in R (“Similarity and Difference Finder.R”) to extract sequence lists and generate corresponding Substrate Background Frequency tables for the proteins corresponding to the selected peptides. This script provides either the intersection or symmetric difference between those sets as two new output tables containing only the information relevant to the sequences desired.

*Approximating most likely “true negative” sequence list from substrate dataset*—The accession numbers for proteins that remain in the Substrate Background Frequency list after the previous filter are submitted to the reviewed human Uniprot/SwissProt database (<http://uniprot.org/uploadlist/>) to generate a FASTA file containing the sequences of those

proteins. The FASTA file is converted separately to .csv format using a script obtained from ([https://www.researchgate.net/post/Converting\\_a\\_fasta\\_file\\_to\\_a\\_tab-delimited\\_file10](https://www.researchgate.net/post/Converting_a_fasta_file_to_a_tab-delimited_file10)). This file and the filtered Positive Substrates list file (generated as described in the previous section) are used as input for the “NegativeMotifFinder.R” to extract additional tyrosine-containing sequences from those proteins that could in principle have been phosphorylated but were not detected (outputting a “Negative Motifs” .csv file that is named by the user upon running the script). “Negative Motifs” files and corresponding “Positive Substrates” files are later used by the Kinatestpart1.R script to calculate Matthews Correlation Coefficient (MCC) values that give a general threshold for which peptides will or will not be phosphorylated by the kinase of interest.

## Tables

Substrate	Substrate Sequence	Molecular Weights (g/mol)	[M+(1)H]	[M+(2)H]	[M+(3)H]	[M+(4)H]	[M+(5)H]	[M+(6)H]	[M+(7)H]	[M+(8)H]
FL-ABLtide	<u>EAIYAAPFAKKBGGG</u> APTYSPPPPPGGRKKRRQRLL	4346.15	N/A	2174.4	1449.8	1087.8	870.3	725.4	622.0	544.5
FLT3tide	FTDRLQQYISTR <u>GGBGG</u>	2109.37	2109.9	1055.7	703.8	N/A	N/A	N/A	N/A	N/A
A	GGDE <u>DNDNYCNPNEE</u> GGBGG	2265.26	2264.2	1132.1	N/A	N/A	N/A	N/A	N/A	N/A
B	GGDE <u>SDDYFNPNEE</u> GGBGG	2283.24	2284.9	1143.1	N/A	N/A	N/A	N/A	N/A	N/A
C	GGDE <u>SDIYANPNEE</u> GGBGG	2205.22	2206.6	1114.9	N/A	N/A	N/A	N/A	N/A	N/A
D	GGDE <u>SDNYFNPNEE</u> GGBGG	2282.26	2281.9	1140.6	N/A	N/A	N/A	N/A	N/A	N/A
E	GGDE <u>SDIYFNPNEE</u> GGBGG	2281.31	2282.9	1152.7 (M+Na)	N/A	N/A	N/A	N/A	N/A	N/A
F	GGDE <u>SDNYFNFNEE</u> GGBGG	2332.32	2334.7	1167; 1179 (M+Na); 1186.6 (M+K)	N/A	N/A	N/A	N/A	N/A	N/A
G	GGDE <u>SNDYFNTNEE</u> GGBGG	2286.25	2287.9	1155.1 (M+Na)	N/A	N/A	N/A	N/A	N/A	N/A
H	GGDE <u>HNQYEQPNEE</u> GGBGG	2341.33	2343.1	1172.1; 1182.6 (M+Na)	781.8	N/A	N/A	N/A	N/A	N/A

Table S1. A summary of the FLT3 Artificial Substrate (FAS) candidate sequences synthesized and assayed in vitro with recombinant FLT3 variants.

Abltide (EAIYAAPFAK; the substrate has been incorporated with an SH3 recognition and cell penetrating sequence and termed FL-Abltide) is a previously known FLT3 peptide substrate and has been used a reference substrate to monitor kinase activity.<sup>1</sup> The substrate sequences derived from the KINATEST-ID pipeline are underlined and were synthesized within the terbium binding motif shell (amino acids not underlined; sequence generated using the Aligner module of KINATEST-ID) with a biotinylated lysine (B) as an enrichment tag. The “Molecular Weights” column summarizes the theoretical weight (molecular weight) of the synthesized peptide sequences. The additional right-hand columns summarize the

major observed mass (M) to charge (m/z) signals  $[M+(n)H]$  for each peptide's LC-MS analysis. Sodium (Na) or potassium (K) adducts were present in the second charge state  $[M+2H]$  for FAStides-E-H. Charge states not observed are denoted as N/A.

Scoring Model	Database size	MCC	Sensitivity	Specificity	Accuracy	Precision	EER	AROC	Threshold
WT-2H	888	0.39	0.73	0.81	80.28	0.31	0.20	0.86	17
WT-OVLP	559	0.38	0.79	0.76	76.44	0.29	0.24	0.85	17
WT-16H	1559	0.34	0.91	0.55	60.95	0.28	0.39	0.78	12
D835Y-16H	2010	0.35	0.92	0.56	60.82	0.30	0.39	0.78	13
ITD-16H	344	0.43	0.54	0.92	88.58	0.45	0.11	0.89	32
SHARED-16H	244	0.45	0.66	0.89	86.50	0.42	0.59	0.89	45

Table S2. Performance metrics and comparison of PSM models. *Scoring model* gives the substrate lists used to develop the scoring model. *Database size* represents the number of substrates in the list. The Matthew's correlation coefficient (*MCC*) is a performance metric for binary classifiers with values within a -1 to +1 scale.<sup>2,3</sup> A value close to 0 indicates a model's prediction is random while a value close to 1 indicates a perfect prediction. A value close to -1 indicates the model is making inverse predictions. *Sensitivity* is the true positive rate (recall) of each scoring model. *Specificity* is the true negative rate (relative to the input dataset). *Accuracy* is the measure of a prediction model's ability to correctly predict an outcome's true classification (i.e. positives vs. negatives from the input dataset). *Precision* is the rate of a model's ability for predicting true results from all its predictions. Equal error rate (*EER*) is the rate where the acceptance and rejection errors are the same. Area under the receiver operator curve (*AROC*) describes the number of correct predictions (of true positives or true negatives) at each given score. *Threshold* is the chosen value for binary classification for each predictive model. The lack of balance in the dataset (i.e. more true positives than true negatives or the inverse) is a potential caveat of these metrics. MCC-based performance metrics are shown to be compatible with imbalanced datasets. However, optimized classifying metrics for imbalanced datasets have been developed<sup>2,4</sup> but require implementation of Bayesian statistics or the development of a support vector machines.<sup>2,5</sup> Based on the satisfactory performance of KINATEST-ID for the applications pursued so far, these

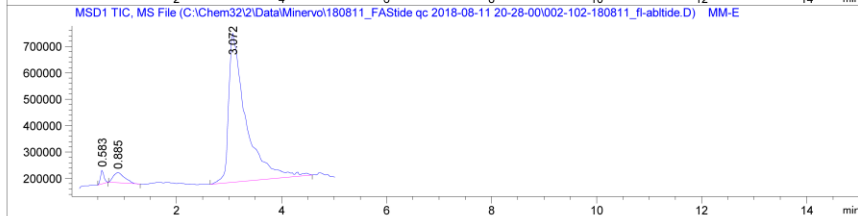
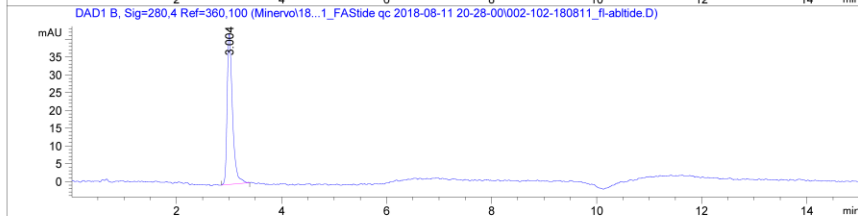
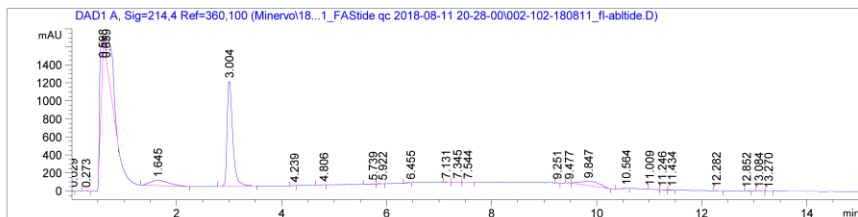
have not yet been examined. However, advanced performance metrics should be considered in future updates of the KINATEST-ID predictive models.

# FL-ABLtide

MS Peak Purity Range Report

Data File C:\Chem32\...\180811\_FAS tide qc 2018-08-11 20-28-00\002-102-180811\_fl-abl tide.D

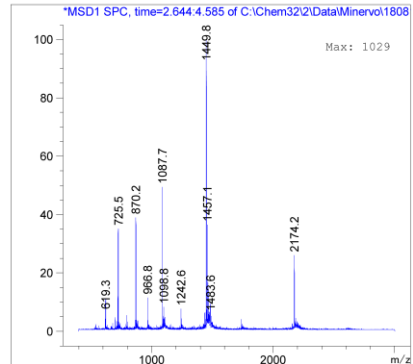
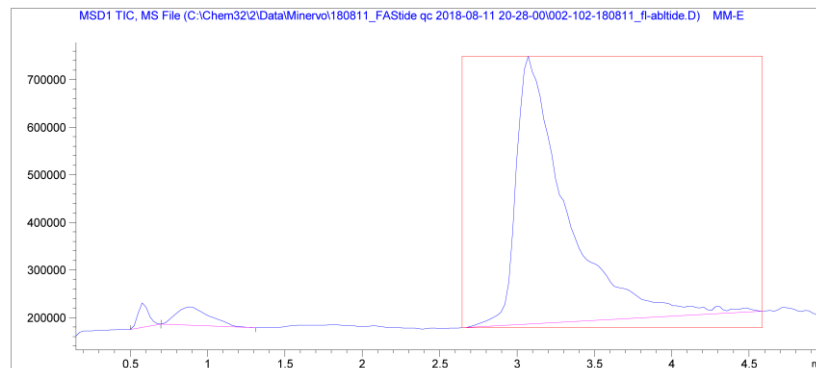
Sample Name: 180811\_fl-abl tide



MS Peak Purity Range Report

Data File C:\Chem32\...\180811\_FAS tide qc 2018-08-11 20-28-00\002-102-180811\_fl-abl tide.D

Sample Name: 180811\_fl-abl tide



MSCalcPurity Error # 407!!!

Peak #3 at 3.072 min ( 2.644 to 4.585 min)  
-> No purity results available. <-

```

=====
Fraction Information
=====
Fraction collection off
=====
No Fractions found.
=====
    
```

\*\*\* End of Report \*\*\*

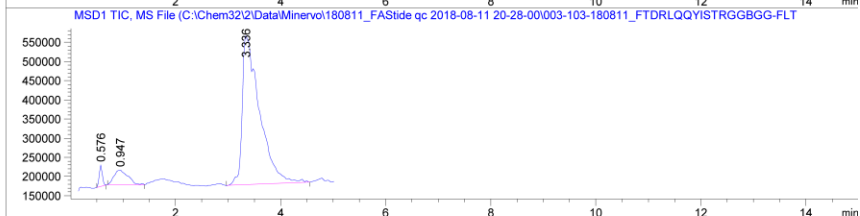
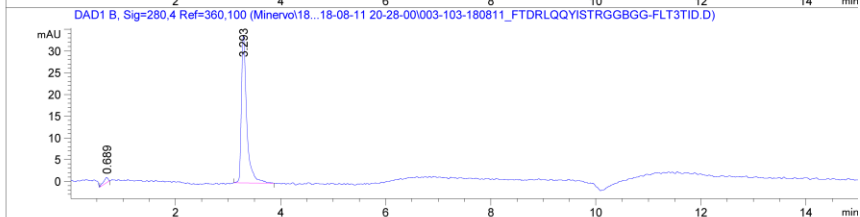
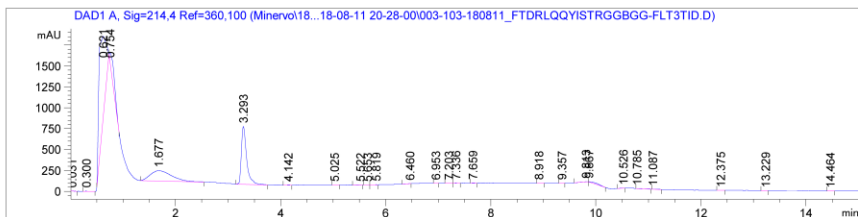


# FLT3tide

## MS Peak Purity Range Report

Data File C:\Chem32\...e qc 2018-08-11 20-28-00\003-103-180811\_FTDLRQQYISTRGGBGG-FLT3TID.D

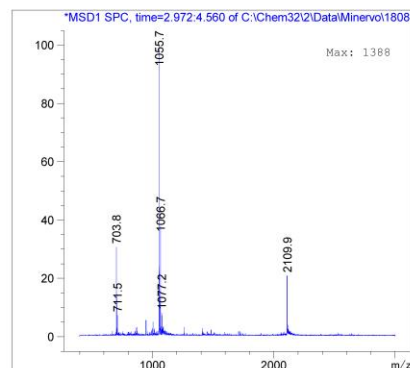
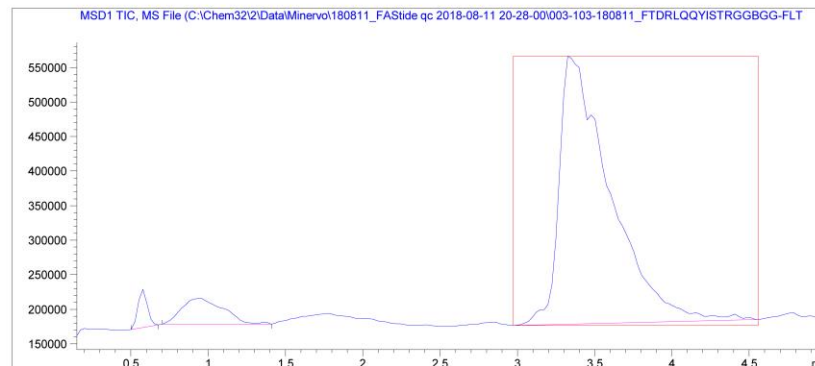
Sample Name: 180811\_FTDLRQQYISTRGGBGG-FLT3TIDE



## MS Peak Purity Range Report

Data File C:\Chem32\...e qc 2018-08-11 20-28-00\003-103-180811\_FTDLRQQYISTRGGBGG-FLT3TID.D

Sample Name: 180811\_FTDLRQQYISTRGGBGG-FLT3TIDE



MSCalcPurity Error # 1!!!

Peak #3 at 3.336 min ( 2.972 to 4.557 min)  
-> No purity results available. <-

\*\*\* End of Report \*\*\*

=====  
Fraction Information  
=====  
Fraction collection off  
=====  
No Fractions found.  
=====

# HPLC-MS Analytical Blank

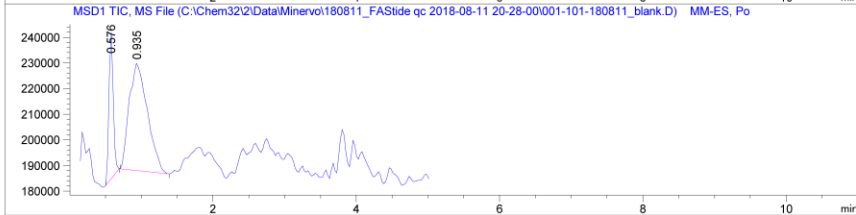
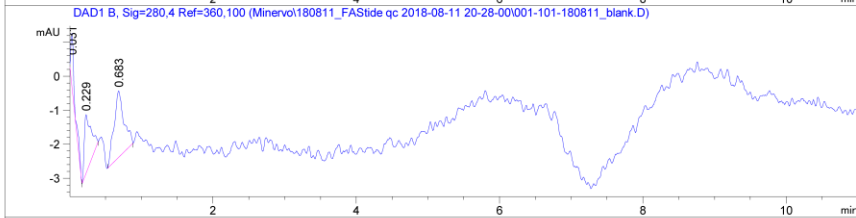
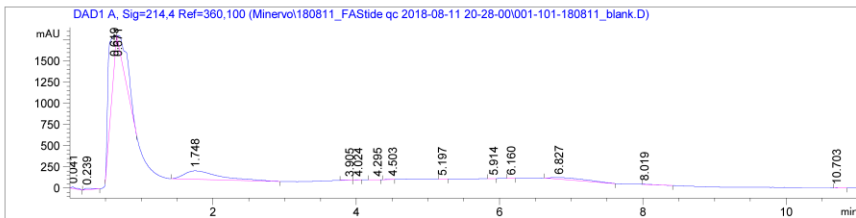
## MS Peak Purity Range Report

Data File C:\Chem32\...\inervo\180811\_FASTide qc 2018-08-11 20-28-00\001-101-180811\_blank.D  
 Sample Name: 180811\_blank

DMSO IN PBS diluted with hplc water with 0.1% fa

```

=====
Acq. Operator   : SYSTEM                      Seq. Line :    1
Acq. Instrument : ChemStation Online LC-MS     Location  :   101
Injection Date  : 8/11/2018 8:29:36 PM        Inj       :    1
                                           Inj Volume: 50.000 µl
Acq. Method     : C:\Chem32\2\Data\Minervo\180811_FASTide qc 2018-08-11 20-28-00\180223_
                  ANALYTICA BLANK-HPLC-MS.M
Last changed    : 8/11/2018 8:28:00 PM by SYSTEM
Analysis Method : C:\Chem32\2\Data\Minervo\180811_FASTide qc 2018-08-11 20-28-00\180223_
                  ANALYTICA BLANK-HPLC-MS.M (Sequence Method)
Last changed    : 8/12/2018 12:30:52 AM by SYSTEM
Sample Info     : DMSO IN PBS diluted with hplc water with 0.1% fa
  
```



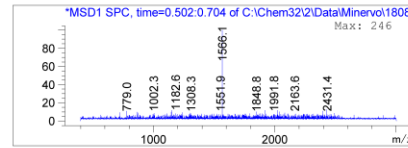
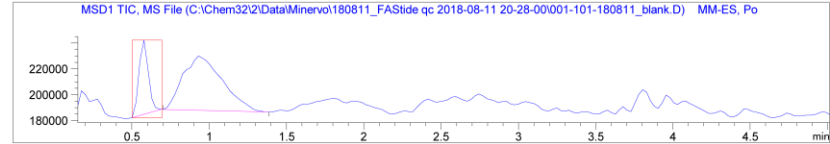
### Fraction Information

Fraction collection off

No Fractions found.

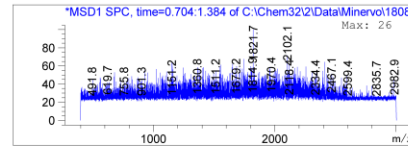
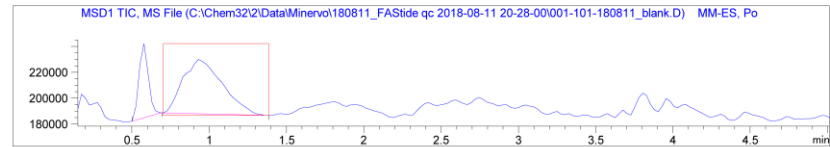
## MS Peak Purity Range Report

Data File C:\Chem32\...\inervo\180811\_FASTide qc 2018-08-11 20-28-00\001-101-180811\_blank.D  
 Sample Name: 180811\_blank



MSCalcPurity Error # 407!!!

Peak #1 at 0.576 min ( 0.502 to 0.694 min)  
 -> No purity results available. <-



MSCalcPurity Error # 407!!!

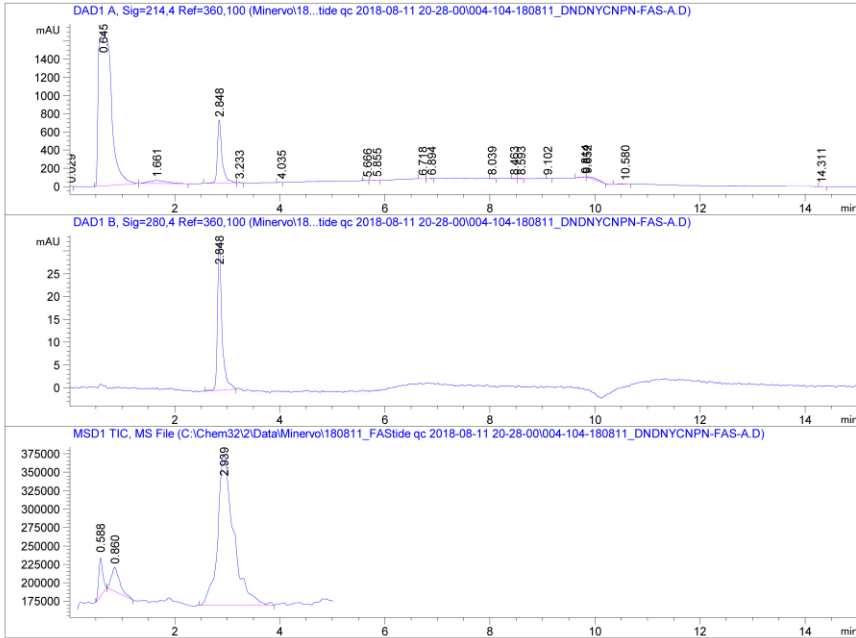
Peak #2 at 0.935 min ( 0.704 to 1.384 min)  
 -> No purity results available. <-

\*\*\* End of Report \*\*\*

# FASTide-A

## MS Peak Purity Range Report

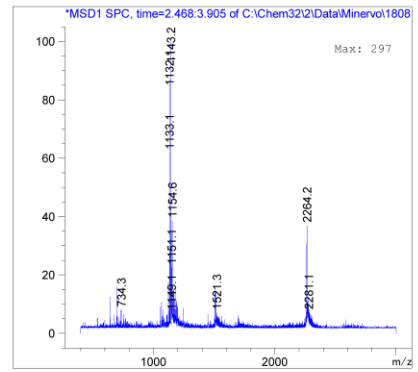
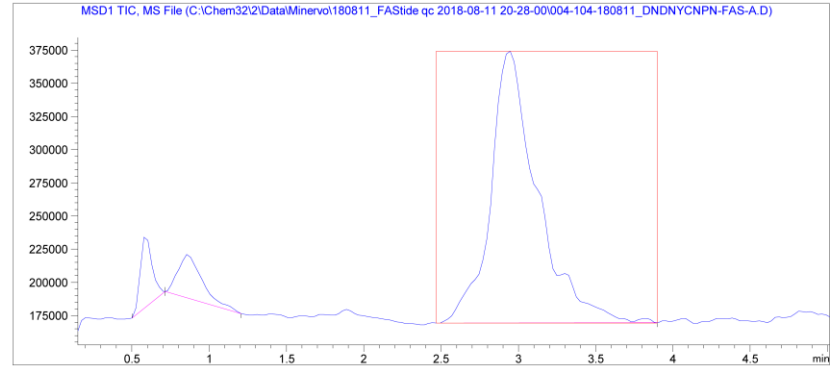
Data File C:\Chem32\...811\_FASTide qc 2018-08-11 20-28-00\004-104-180811\_DNDNYCNPN-FAS-A.D  
Sample Name: 180811\_DNDNYCNPN-FAS-A



=====  
Fraction Information  
=====  
Fraction collection off  
=====  
No Fractions found.  
=====

## MS Peak Purity Range Report

Data File C:\Chem32\...811\_FASTide qc 2018-08-11 20-28-00\004-104-180811\_DNDNYCNPN-FAS-A.D  
Sample Name: 180811\_DNDNYCNPN-FAS-A



MSCalcPurity Error # 407!!!

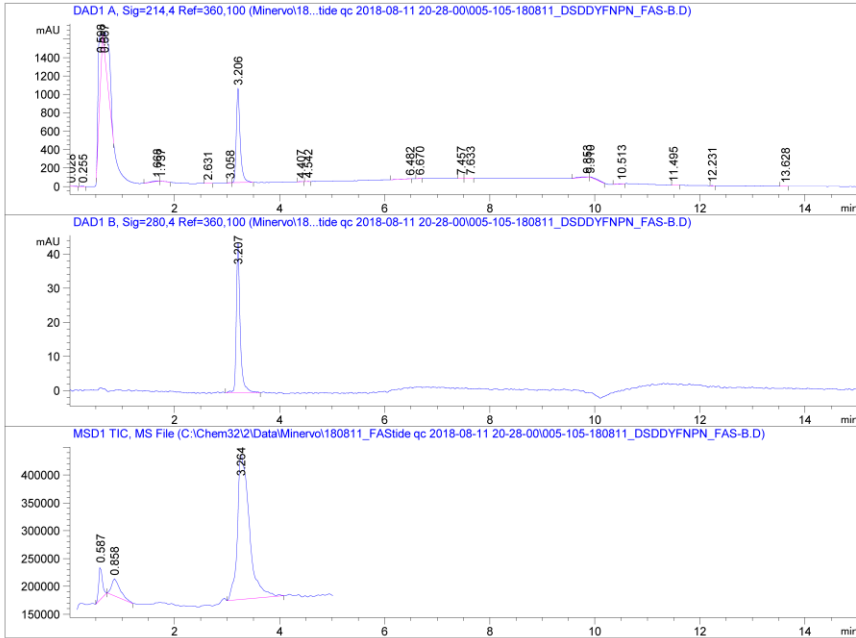
Peak #3 at 2.939 min ( 2.468 to 3.897 min)  
-> No purity results available. <-

\*\*\* End of Report \*\*\*

# FASTide-B

## MS Peak Purity Range Report

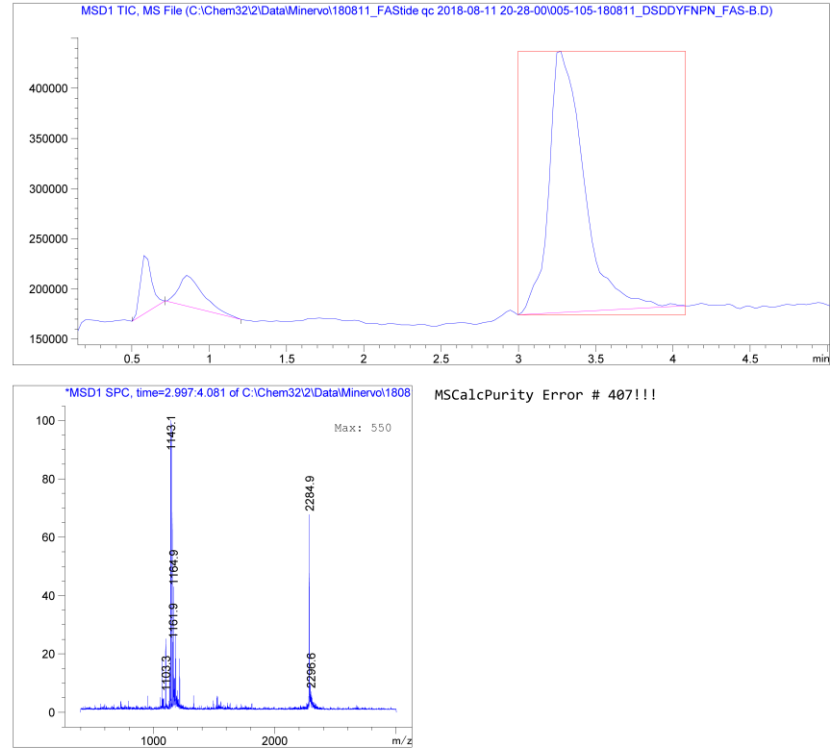
Data File C:\Chem32\...811\_FASTide qc 2018-08-11 20-28-00\005-105-180811\_DSDDYFNP\_N\_FAS-B.D  
Sample Name: 180811\_DSDDYFNP\_N\_FAS-B



=====  
Fraction Information  
=====  
Fraction collection off  
=====  
No Fractions found.  
=====

## MS Peak Purity Range Report

Data File C:\Chem32\...811\_FASTide qc 2018-08-11 20-28-00\005-105-180811\_DSDDYFNP\_N\_FAS-B.D  
Sample Name: 180811\_DSDDYFNP\_N\_FAS-B



Peak #3 at 3.264 min ( 2.997 to 4.081 min)  
-> No purity results available. <-

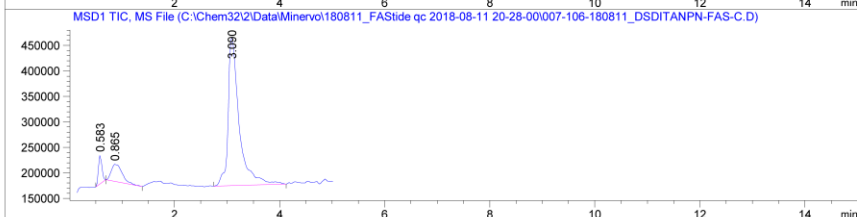
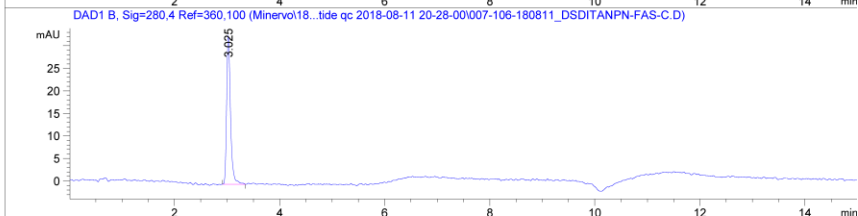
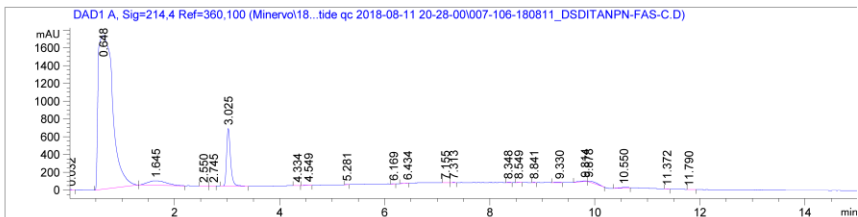
\*\*\* End of Report \*\*\*

# FASTide-C

MS Peak Purity Range Report

Data File C:\Chem32\...811\_FASTide qc 2018-08-11 20-28-00\007-106-180811\_DSDITANPN-FAS-C.D

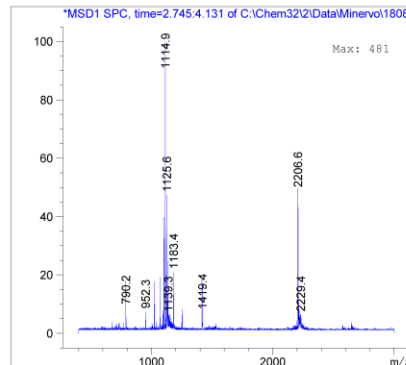
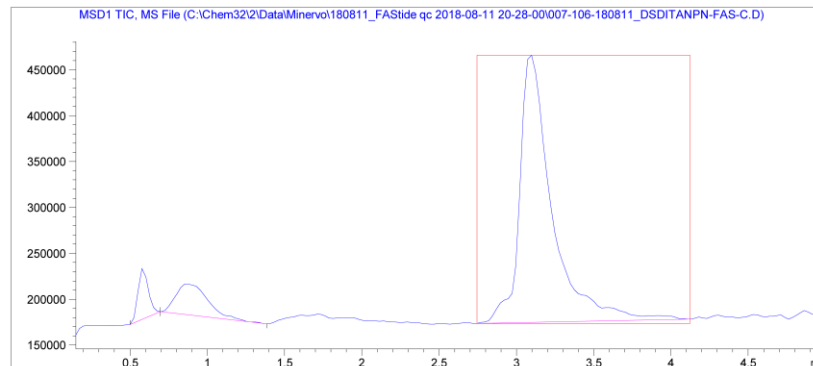
Sample Name: 180811\_DSDITANPN-FAS-C



MS Peak Purity Range Report

Data File C:\Chem32\...811\_FASTide qc 2018-08-11 20-28-00\007-106-180811\_DSDITANPN-FAS-C.D

Sample Name: 180811\_DSDITANPN-FAS-C



MSCalcPurity Error # 407!!!

Peak #3 at 3.090 min ( 2.745 to 4.123 min)  
-> No purity results available. <-

```

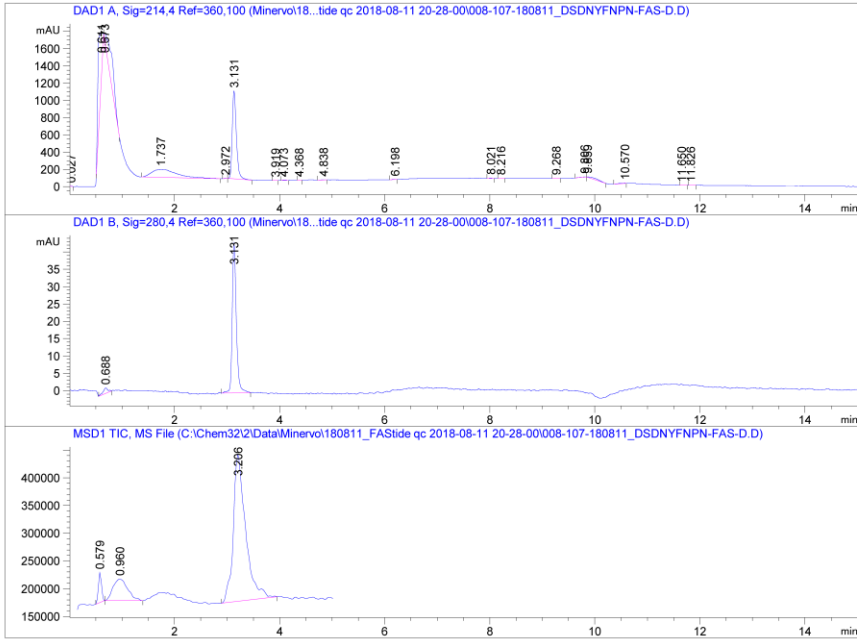
=====
Fraction Information
=====
Fraction collection off
=====
No Fractions found.
=====
    
```

\*\*\* End of Report \*\*\*

# FASTide-D

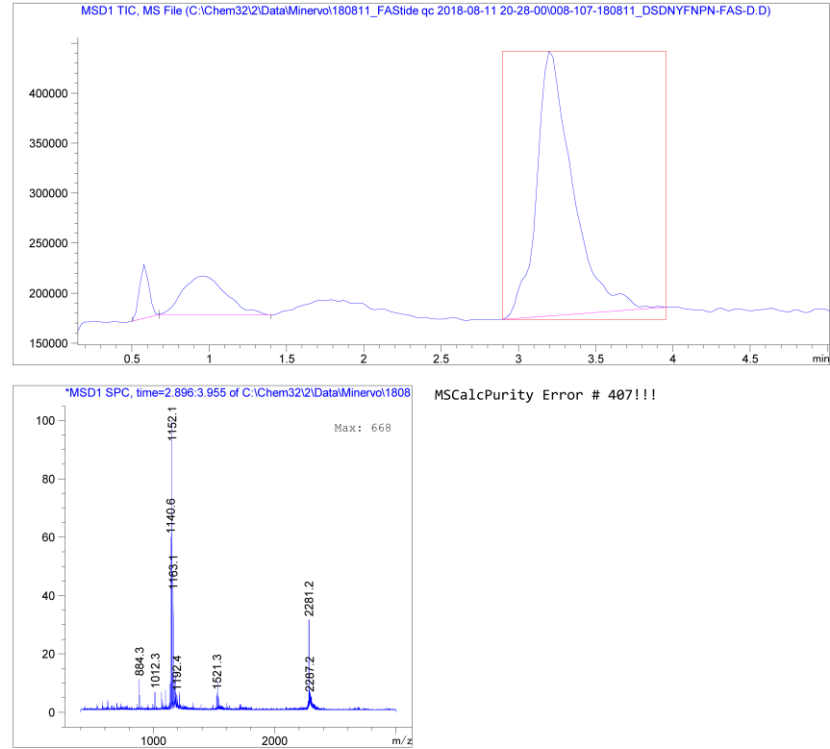
## MS Peak Purity Range Report

Data File C:\Chem32\...811\_FASTide qc 2018-08-11 20-28-00\008-107-180811\_DSDNYFNP-N-FAS-D.D  
Sample Name: 180811\_DSDNYFNP-N-FAS-D



## MS Peak Purity Range Report

Data File C:\Chem32\...811\_FASTide qc 2018-08-11 20-28-00\008-107-180811\_DSDNYFNP-N-FAS-D.D  
Sample Name: 180811\_DSDNYFNP-N-FAS-D



=====  
Fraction Information  
=====  
Fraction collection off  
=====  
No Fractions found.  
=====

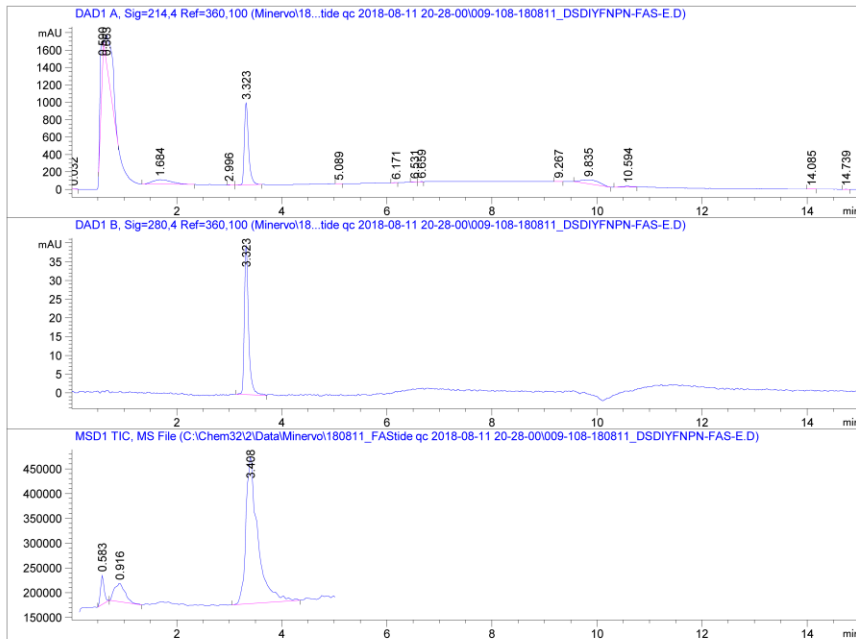
Peak #3 at 3.206 min ( 2.896 to 3.956 min)  
-> No purity results available. <-

\*\*\* End of Report \*\*\*

# FASTide-E

## MS Peak Purity Range Report

Data File C:\Chem32\...811\_FASTide qc 2018-08-11 20-28-00\009-108-180811\_DSDIYFNP-N-FAS-E.D  
Sample Name: 180811\_DSDIYFNP-N-FAS-E



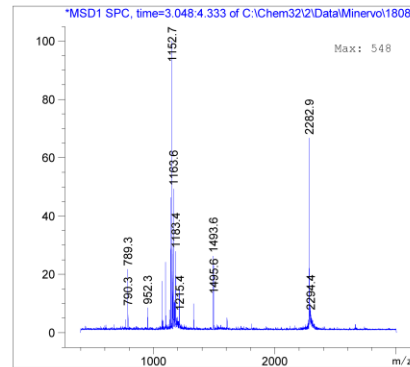
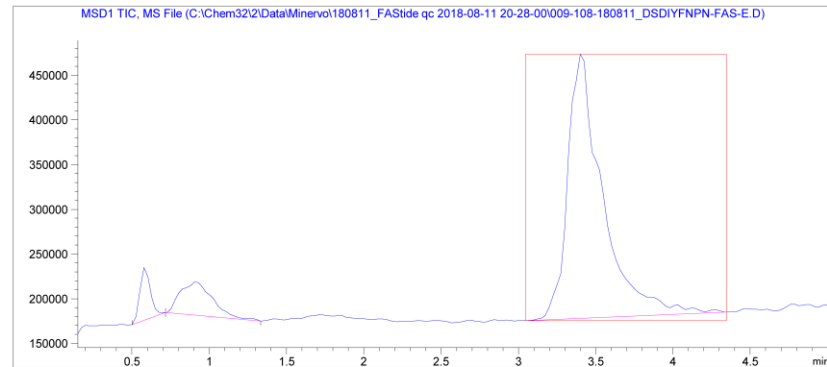
### Fraction Information

Fraction collection off

No Fractions found.

## MS Peak Purity Range Report

Data File C:\Chem32\...811\_FASTide qc 2018-08-11 20-28-00\009-108-180811\_DSDIYFNP-N-FAS-E.D  
Sample Name: 180811\_DSDIYFNP-N-FAS-E



MSCalcPurity Error # 407!!!

Peak #3 at 3.408 min ( 3.048 to 4.345 min)

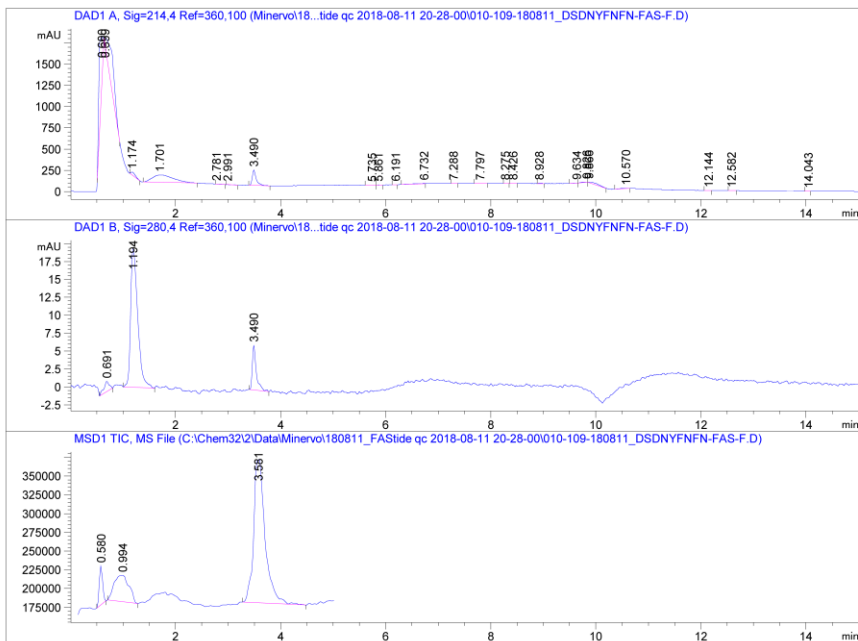
-> No purity results available. <-

\*\*\* End of Report \*\*\*

# FASTide-F

## MS Peak Purity Range Report

Data File C:\Chem32\...811\_FASTide qc 2018-08-11 20-28-00\010-109-180811\_DSDNYFNFN-FAS-F.D  
Sample Name: 180811\_DSDNYFNFN-FAS-F



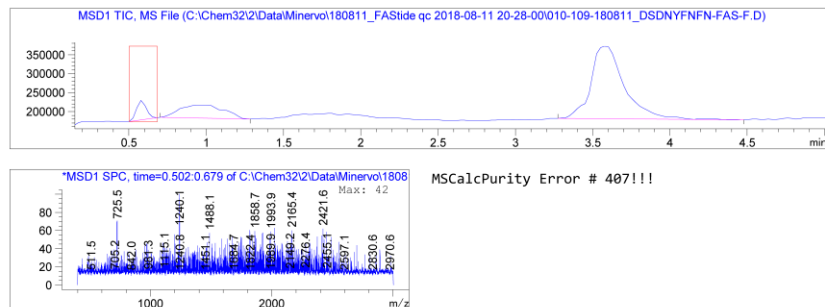
### Fraction Information

Fraction collection off

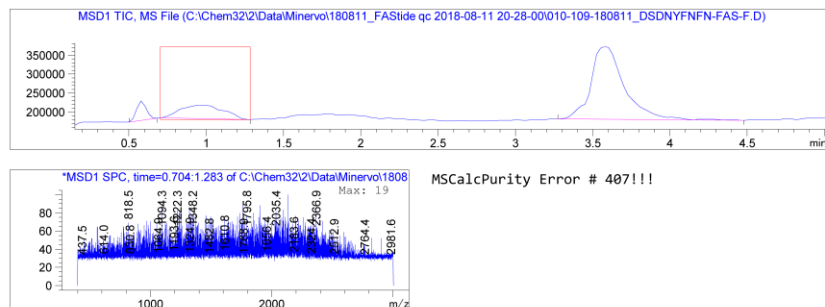
No Fractions found.

## MS Peak Purity Range Report

Data File C:\Chem32\...811\_FASTide qc 2018-08-11 20-28-00\010-109-180811\_DSDNYFNFN-FAS-F.D  
Sample Name: 180811\_DSDNYFNFN-FAS-F

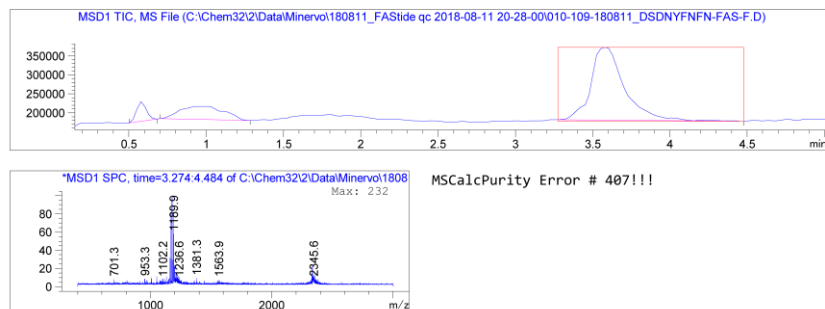


Peak #1 at 0.580 min ( 0.502 to 0.682 min)  
-> No purity results available. <-



## MS Peak Purity Range Report

Data File C:\Chem32\...811\_FASTide qc 2018-08-11 20-28-00\010-109-180811\_DSDNYFNFN-FAS-F.D  
Sample Name: 180811\_DSDNYFNFN-FAS-F



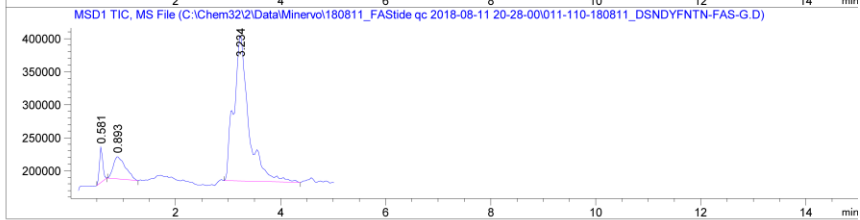
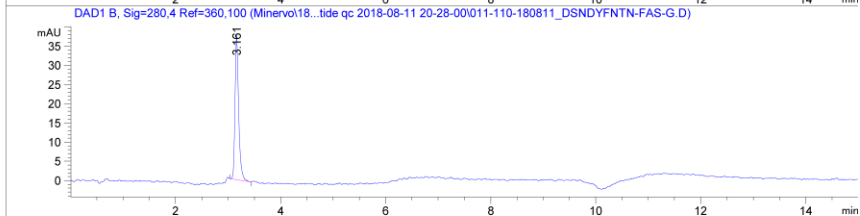
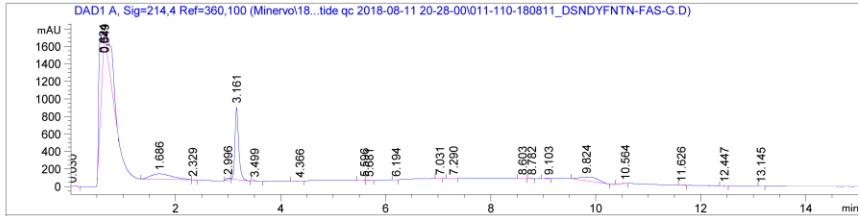


# FASTide-G

## MS Peak Purity Range Report

Data File C:\Chem32\...811\_FASTide qc 2018-08-11 20-28-00\011-110-180811\_DSNDYFNTN-FAS-G.D

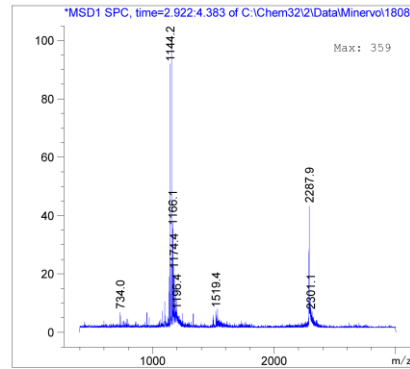
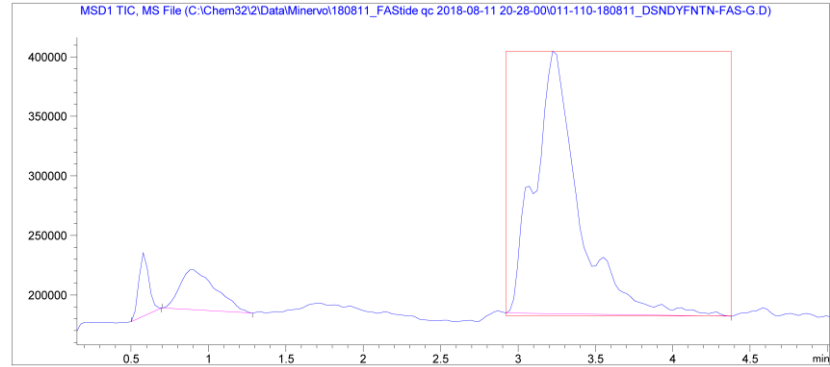
Sample Name: 180811\_DSNDYFNTN-FAS-G



## MS Peak Purity Range Report

Data File C:\Chem32\...811\_FASTide qc 2018-08-11 20-28-00\011-110-180811\_DSNDYFNTN-FAS-G.D

Sample Name: 180811\_DSNDYFNTN-FAS-G



MSCalcPurity Error # 407!!!

Peak #3 at 3.234 min ( 2.922 to 4.377 min)  
-> No purity results available. <-

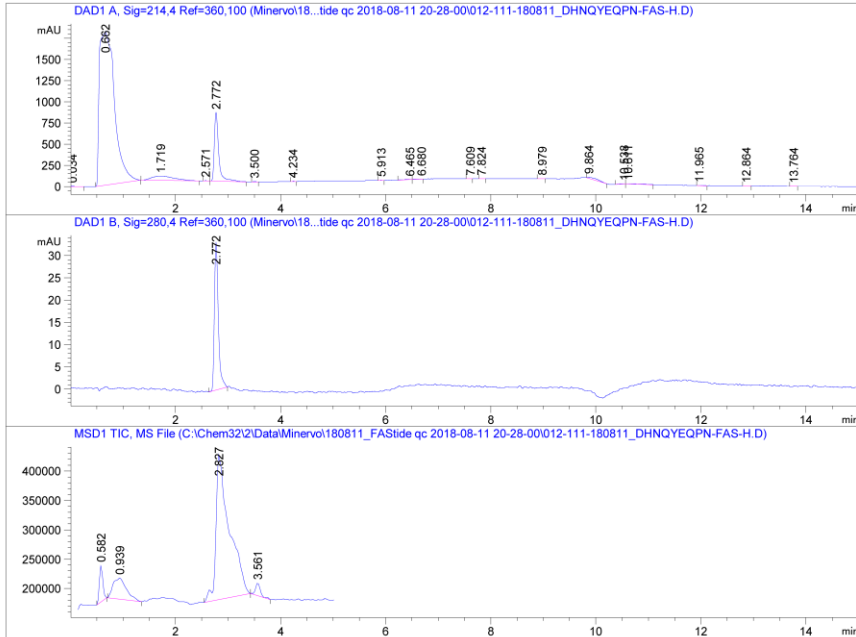
\*\*\* End of Report \*\*\*

=====  
Fraction Information  
=====  
Fraction collection off  
=====  
No Fractions found.  
=====

# FASTide-H

## MS Peak Purity Range Report

Data File C:\Chem32\...811\_FASTide qc 2018-08-11 20-28-00\012-111-180811\_DHNQYEQPN-FAS-H.D  
 Sample Name: 180811\_DHNQYEQPN-FAS-H

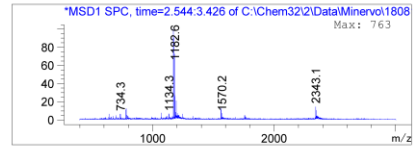
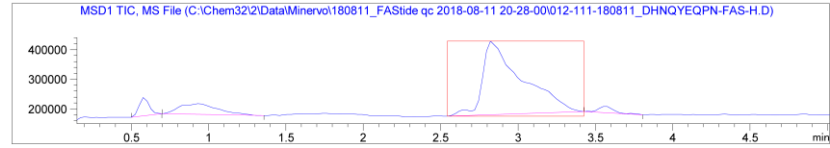


```

=====
                          Fraction Information
=====
Fraction collection off
=====
No Fractions found.
=====
  
```

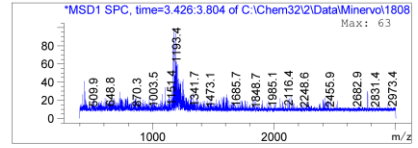
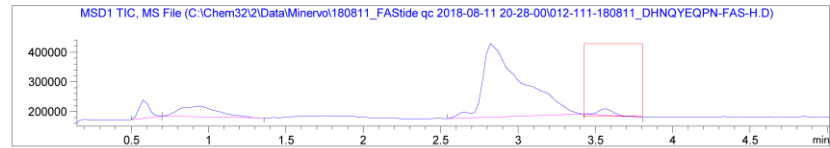
## MS Peak Purity Range Report

Data File C:\Chem32\...811\_FASTide qc 2018-08-11 20-28-00\012-111-180811\_DHNQYEQPN-FAS-H.D  
 Sample Name: 180811\_DHNQYEQPN-FAS-H



MSCalcPurity Error # 407!!!

Peak #3 at 2.827 min ( 2.544 to 3.426 min)  
 -> No purity results available. <-



MSCalcPurity Error # 407!!!

Peak #4 at 3.561 min ( 3.429 to 3.804 min)  
 -> No purity results available. <-

\*\*\* End of Report \*\*\*

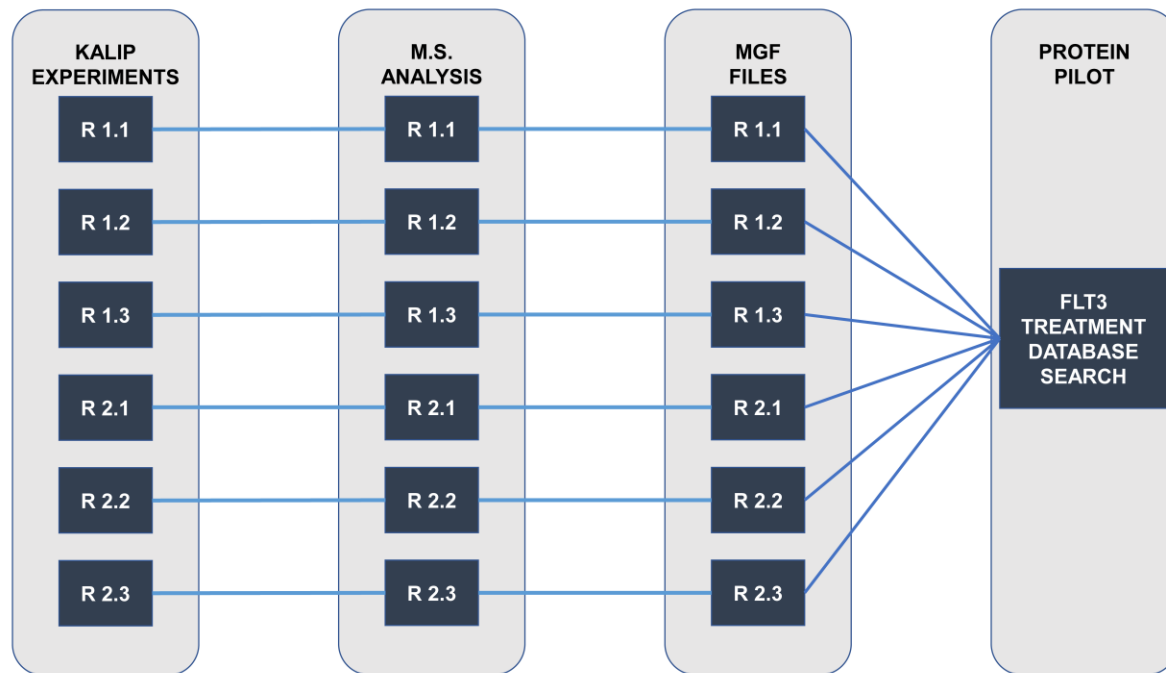


Figure S1 Schematic representation of raw mass spectrometer file combination for ProteinPilot database searches. Each KALIP kinase treatment (WT, D835Y and/or ITD) was performed with three biological replicates (R1). The KALIP process was then repeated later to generate a second independent KALIP technical experiment (R2). Replicates were individually analyzed on the mass spectrometer and then converted to MGF files, ProteinPilot 5.0 database search consisted of six mass spectrometer files for each kinase treatment (no kinase or kinase treatment).

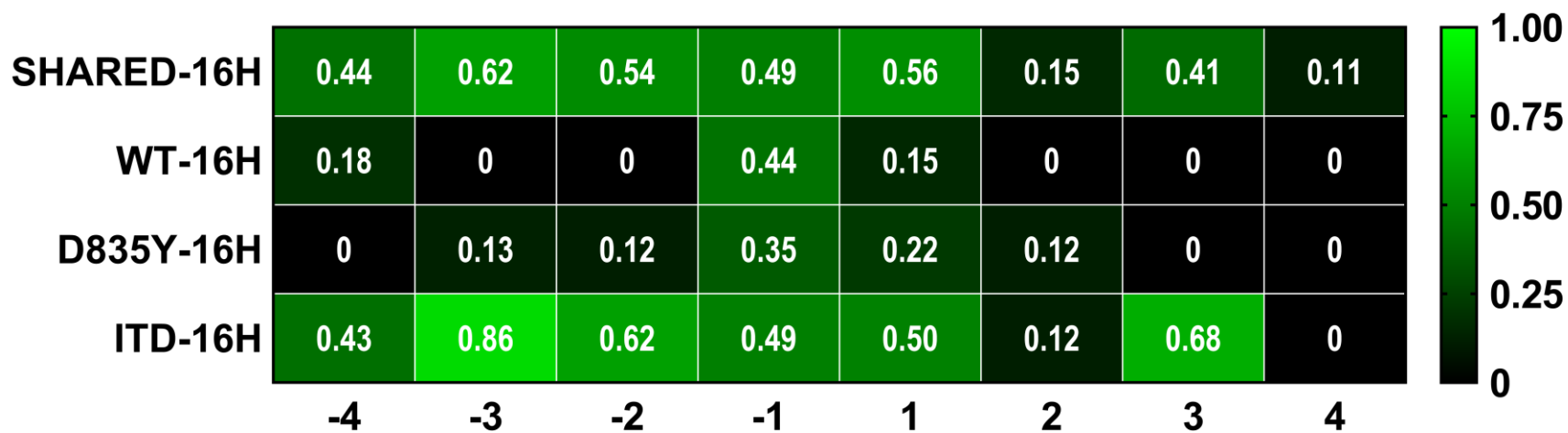
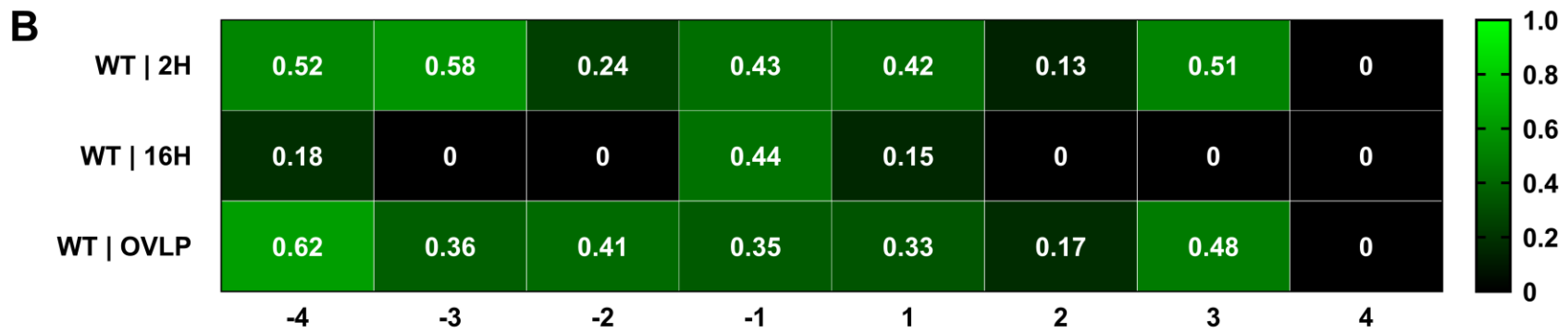
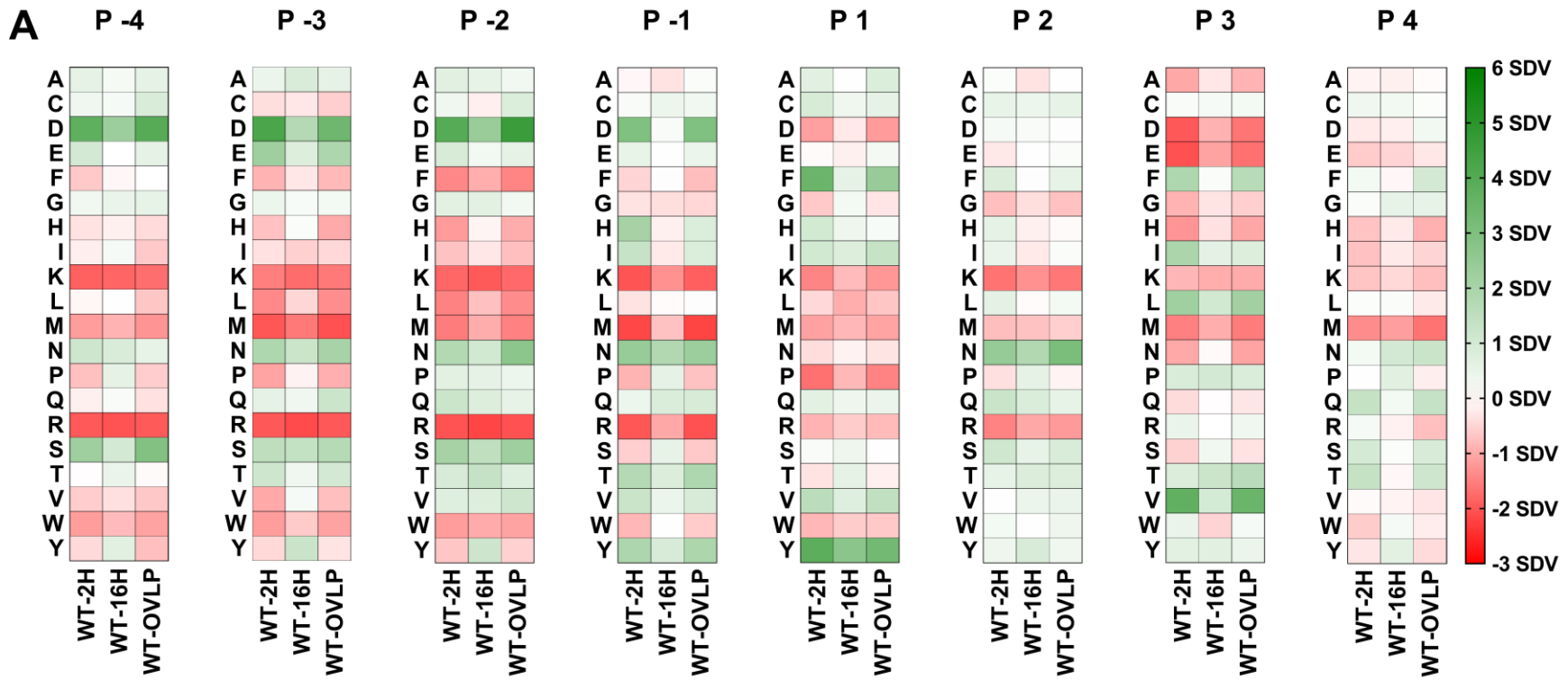


Figure S2 is a heat map representation of the Site Selectivity Matrix (SSM) values found in “Output file 2” and was generated as previously reported.<sup>31</sup> SSM values closer to 1 suggest that the kinase of interest would be more sensitive to changes in the particular residue at this position.



**Figure S3. Heat map representation of FLT3-WT time course KALIP experiment, Site Selectivity Matrix and artificial substrate library sequence scoring comparison.** (A) Observed representation of each amino acid at each position (-4 to +4 relative to phosphotyrosine) in the individual phosphoproteomics datasets for the kinase treatments at two hours (WT-2H) or sixteen hours (WT-16H), or for the sequences shared in the two datasets (WT-OVLP). Green = over-represented, white = neutral, red = under-represented. To summarize, differences were modest between the two treatment times. (B) We compared the three substrate lists' SSM values to identify positions with a value greater than 1, which is the previously reported threshold used to consider a position as "significant."<sup>1</sup> None of the KALIP dataset SSMs contained a position with a value greater than one, suggesting that all positions exhibited some flexibility for which particular amino acid was present.

## References

1. Lipchik, A. M. *et al.* KINATEST-ID: A pipeline to develop phosphorylation-dependent terbium sensitizing kinase assays. *J. Am. Chem. Soc.* **137**, 2484–2494 (2015).
2. Boughorbel, S., Jarray, F. & El-Anbari, M. Optimal classifier for imbalanced data using Matthews Correlation Coefficient metric. *PLoS One* **12**, e0177678 (2017).
3. Jiao, Y. & Du, P. Performance measures in evaluating machine learning based bioinformatics predictors for classifications. *Quant. Biol.* **4**, 320–330 (2016).
4. Narasimhan, H., Vaish, R. & Agarwal, S. On the Statistical Consistency of Plug-in Classifiers for Non-decomposable Performance Measures. *Proc. 27th Int. Conf. Neural Inf. Process. Syst.* 1493–1501 (2014). at <http://dl.acm.org/citation.cfm?id=2968826.2968993>
5. Breu, F., Guggenbichler, S. & Wollmann, J. *Machine Learning: ECML 2004. 15th European Conference on Machine Learning Pisa, Italy, September 2004 Proceedings* **3201**, (Springer Berlin Heidelberg, 2004).