



Leveraging MSI Data in NTA Workflows to Improve PFAS Discovery

David Schiessel
Babcock Labs, Inc.
dschiessel@babcocklabs.com

Outline



The Lop-Sided Nature of NTA Data Features

Case of No MS/MS data

Recent Developments in PFAS Prioritization

Kendrick, Kaufmann, and Zweigle

RT Models and Molecular Formula Decomposition

Conclusion



Non-Targeted Analyses generate results with varying confidence levels.

What is NTA?

Non-Targeted Analysis

BP4NTA - The characterization of the chemical composition of any given sample without the use of a priori knowledge regarding the sample's chemical content.

The framework by which a defined chemical space is investigated within a sample without a priori knowledge for the primary purpose of chemical discovery.

NTA Feature Distribution

Body of Knowledge Limited

Features in Electrospray may be chemicals, in-source fragments, or clusters/adducts

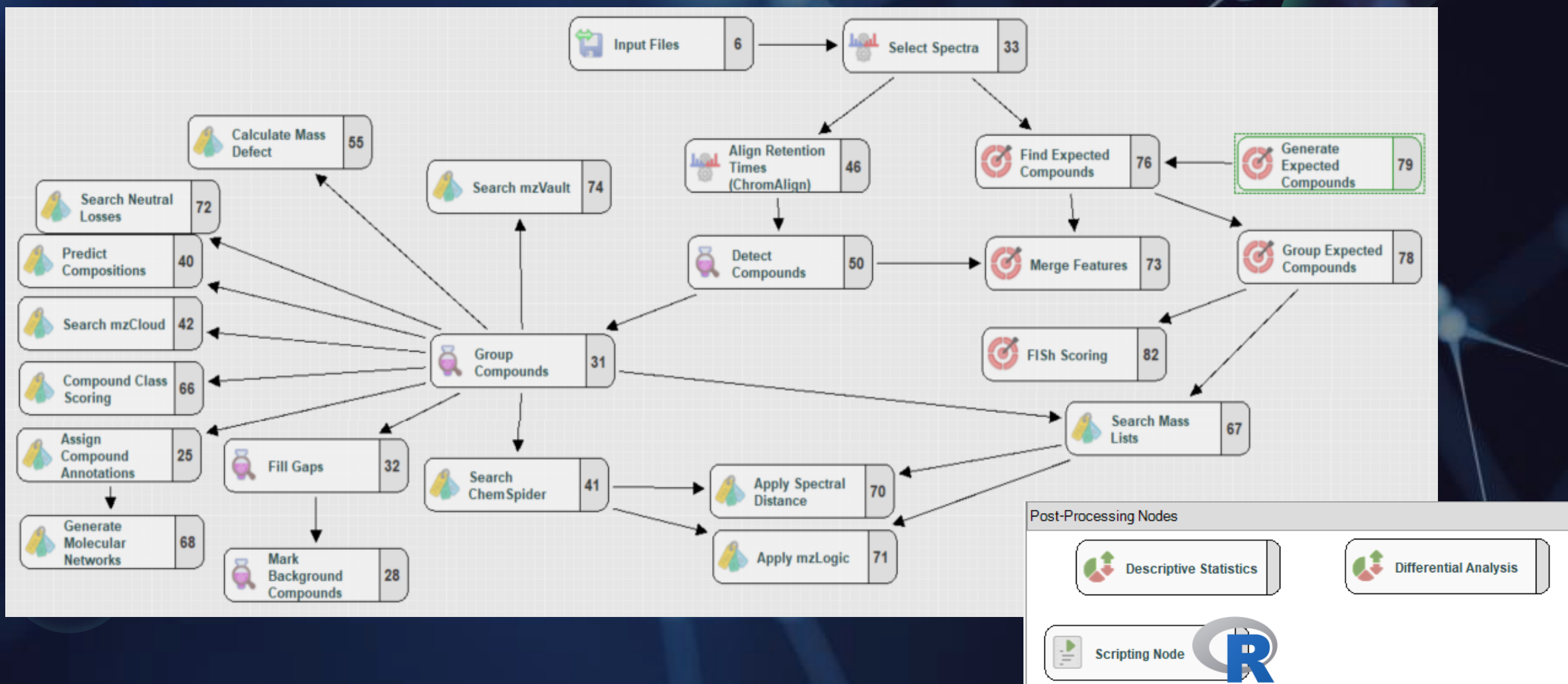
Chemical space of full scale NTA study is vast
(+9000 for drinking water)

More difficult to deconvolute In-source fragments than clusters
(once a bond breaks, the neutral is lost)

Confidence Annotation of features naturally heavily weighted on lower confidence scores
(Few level I and many Level 4/5)

Complete NTA Workflow

Example – Thermo Compound Discoverer 3.3



No MS/MS data

Conundrum of ddMS2 Workflows

Typical Full Scan NTA acquisition has a data-dependent MS2 component to it.

Triggered only when chemical feature is above threshold

Trade-off of ddMS2 over All ions fragmentation (AIF) is ease of deconvolution

Ultimately, many features do not have MS/MS data



Can we increase identification confidence for the many features that have no MS2 data?

Kendrick Mass / Mass Defect

Plot of Kendrick Mass Defect (KMD) vs. MW

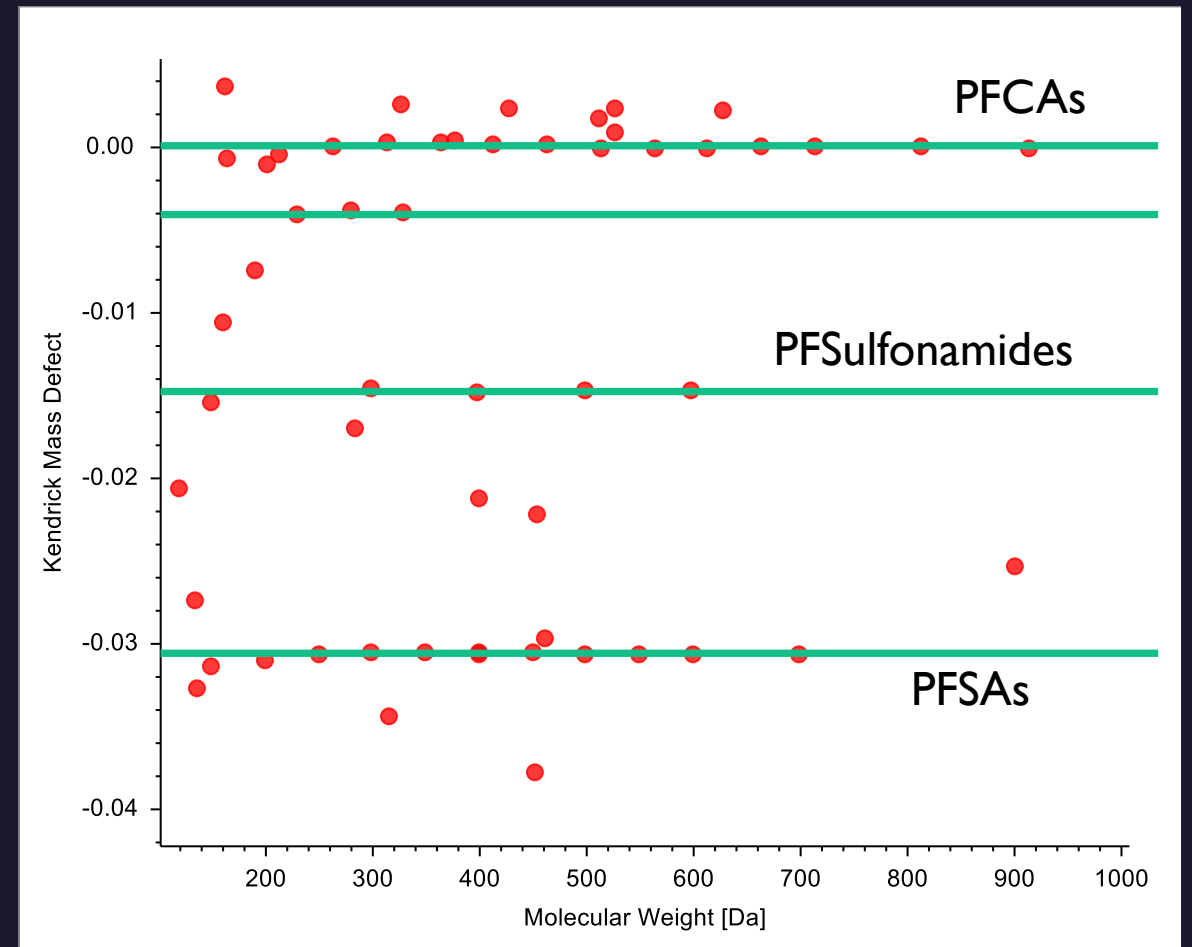
$$\text{KMD}(\text{CF}_2) = M * \text{round}(\text{KFM}) / \text{KFM} - \text{round}(M)$$

Where $\text{KFM} = 49.9968$ and $M = \text{mass}(\text{feature})$

Produces horizontal lines where CF_2 homologs occur

$$\text{Mass Defect} = \text{MW} - \text{floor}(\text{MW})$$

Disadvantage: Scales poorly to non CF_2 homologs



What's in a M/Z

Example: [M-H]⁻ ion for PFOA

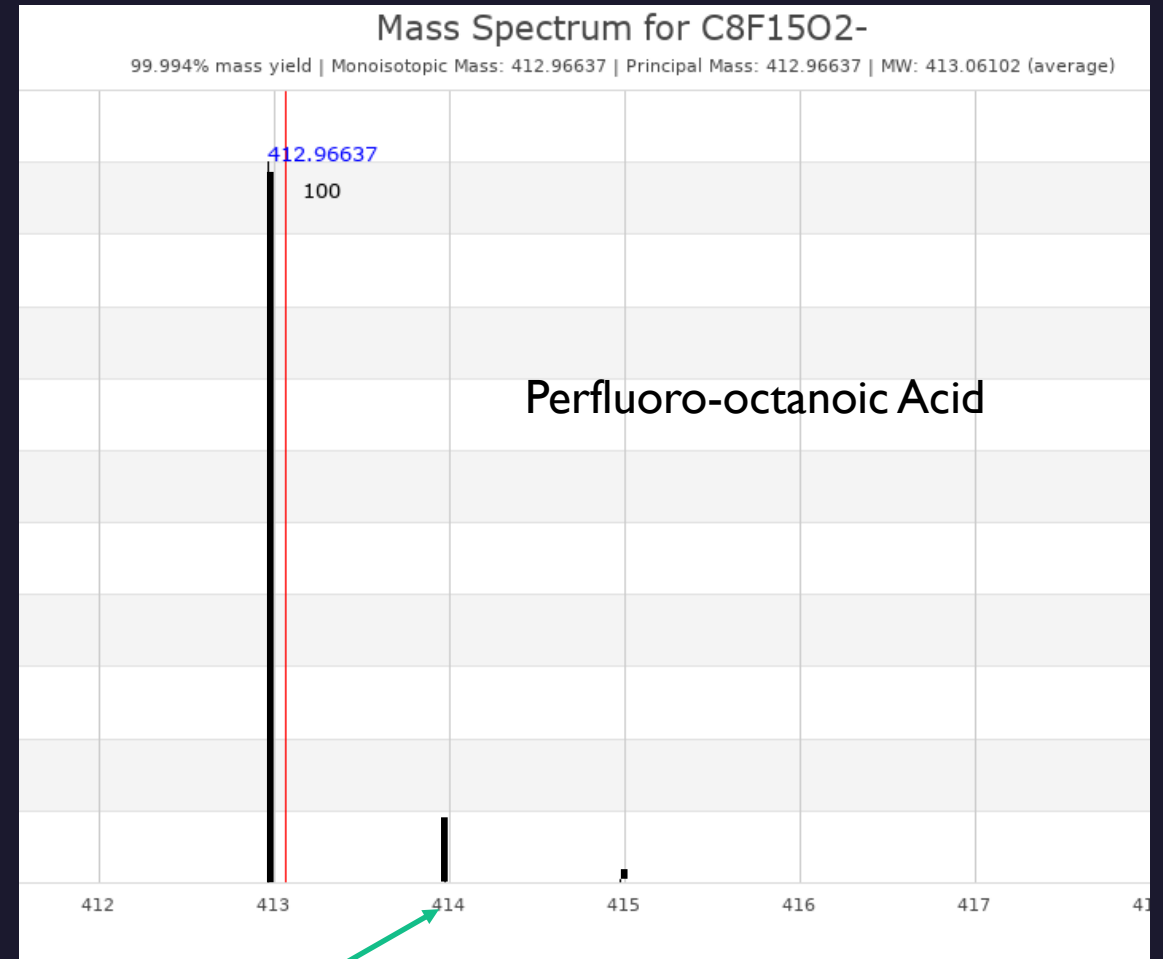
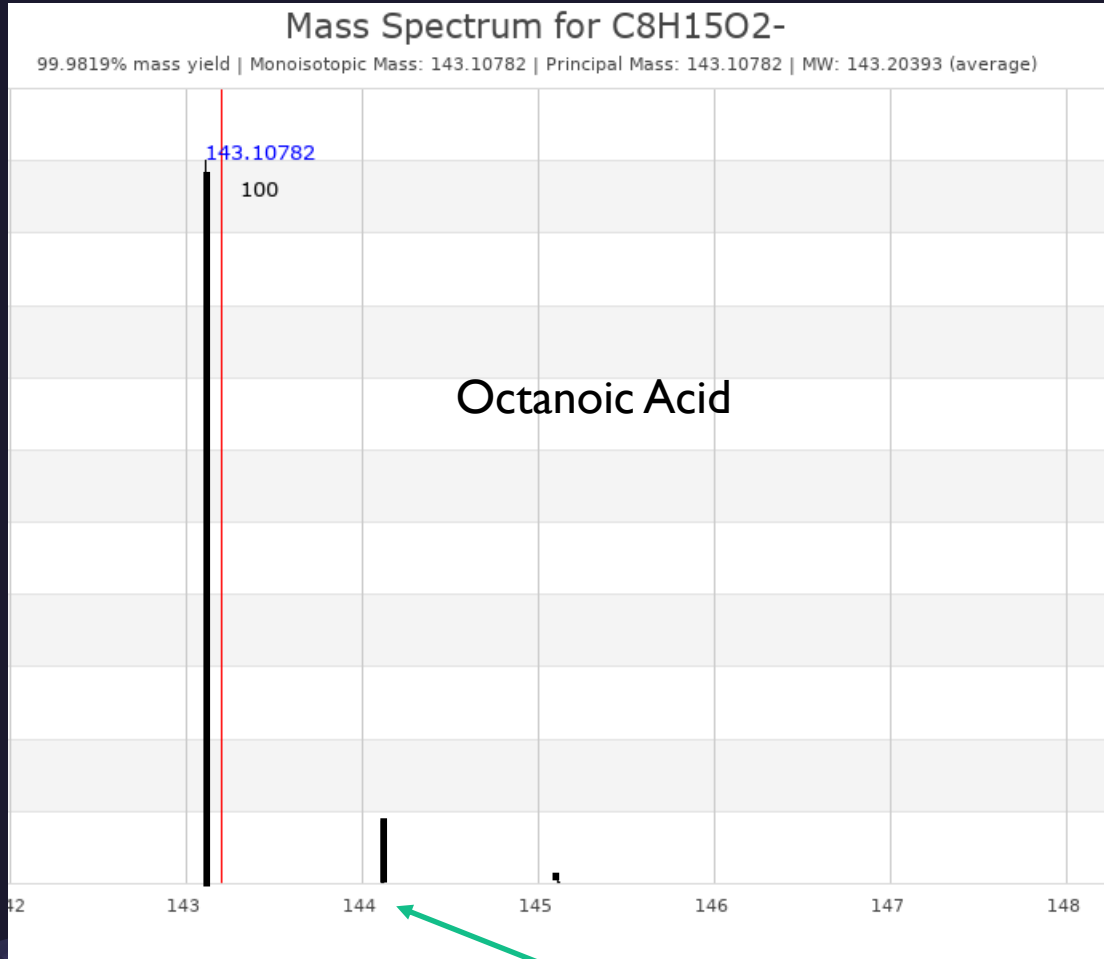
For PFAS:

- m/z is larger because F atoms have replaced H atoms
- Mass defect becomes larger as more F atoms added

412.96637

Mass Defect – Δ from integer

What's in a Mass Spectrum (MS1 Full Scan)



If Carbon 13 abundance is 8.73%, this feature has 8 carbons (independent of F atoms)

What's in a Mass Spectrum (MS1 Full Scan)

OCTANOIC ACID

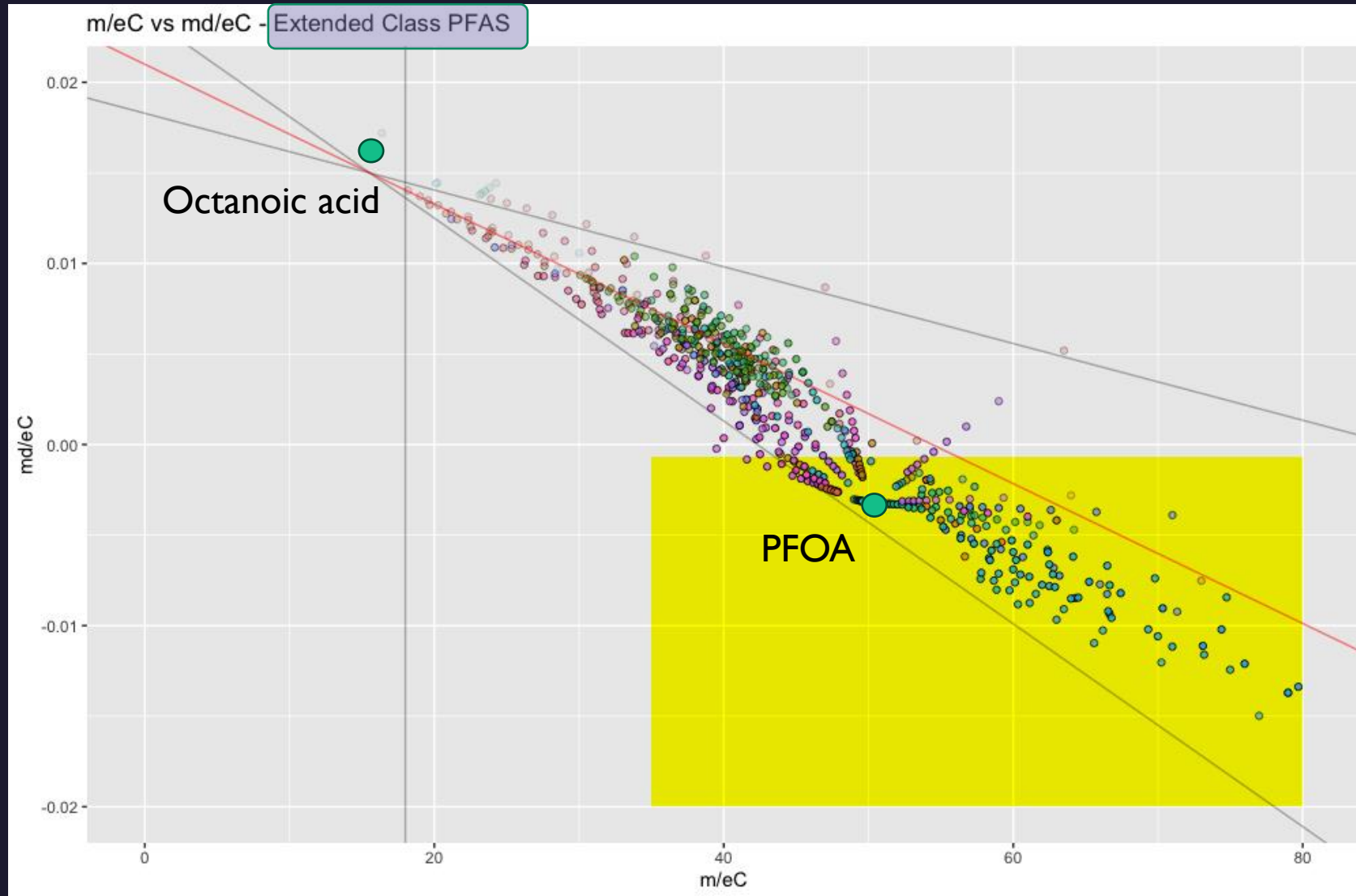
- MF = C₈ H₁₆ O₂
- MW ~ 143
- MW/Carbon ~ 18
- Mass Defect = 0.10782
- MD/Carbon = 0.0135

PFOA

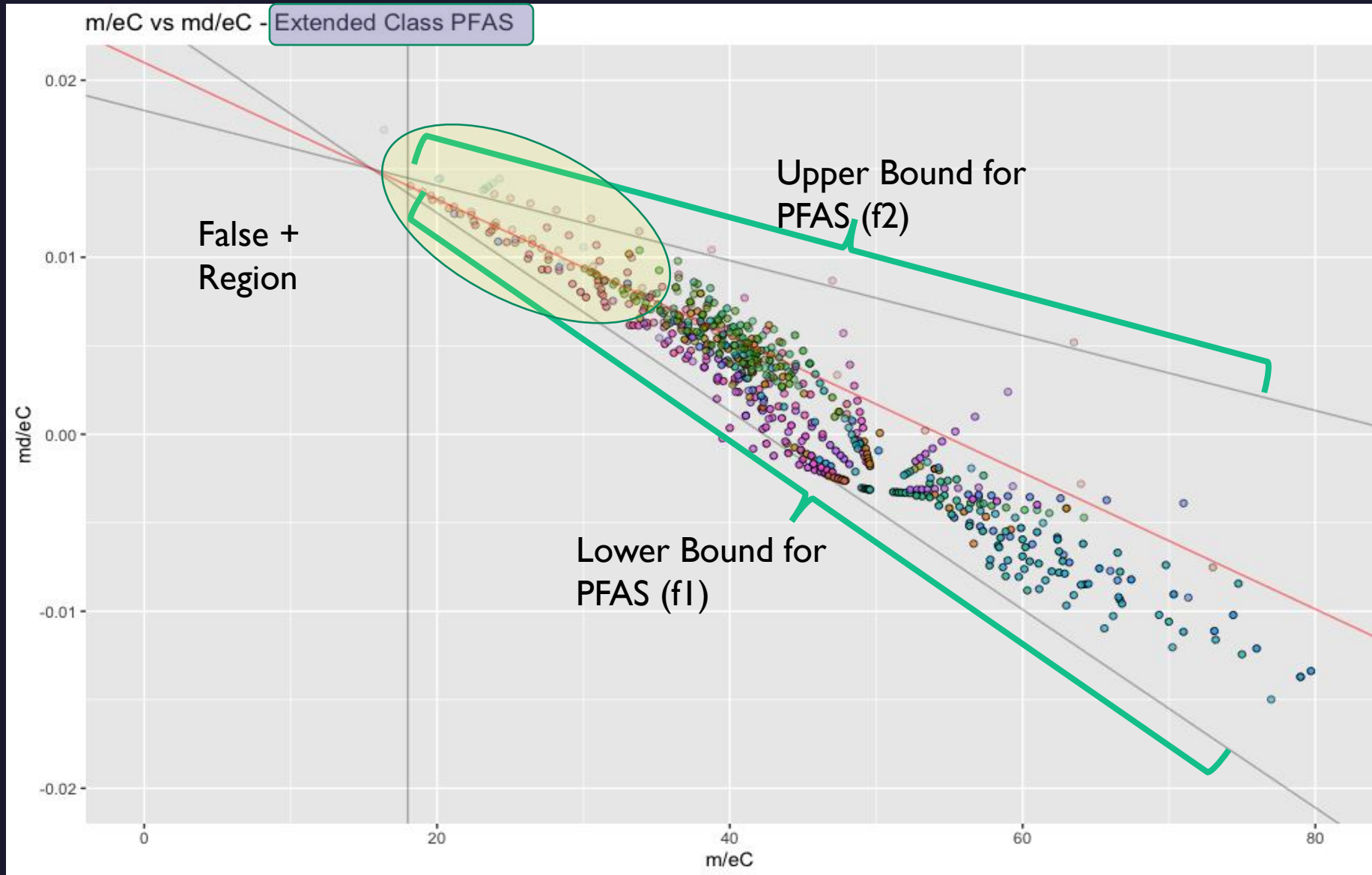
- MF = C₈ H F₁₅ O₂
- MW ~ 413
- MW/Carbon ~ 52
- Mass Defect = -0.03363
- MD/Carbon = -0.0042



Kaufmann Plot of PFAS (md/C vs. m/C)



Kaufmann Plot of PFAS (md/C vs. m/C)





Machine Learning Approach to RT Modeling

RT Modeling – How and Why

Using R packages: rcdk (QSAR), neuralnet to generate a 7:4 MLP

Pretty good prediction +/- 1.0 min, some outliers

Not universal and not directly transferable

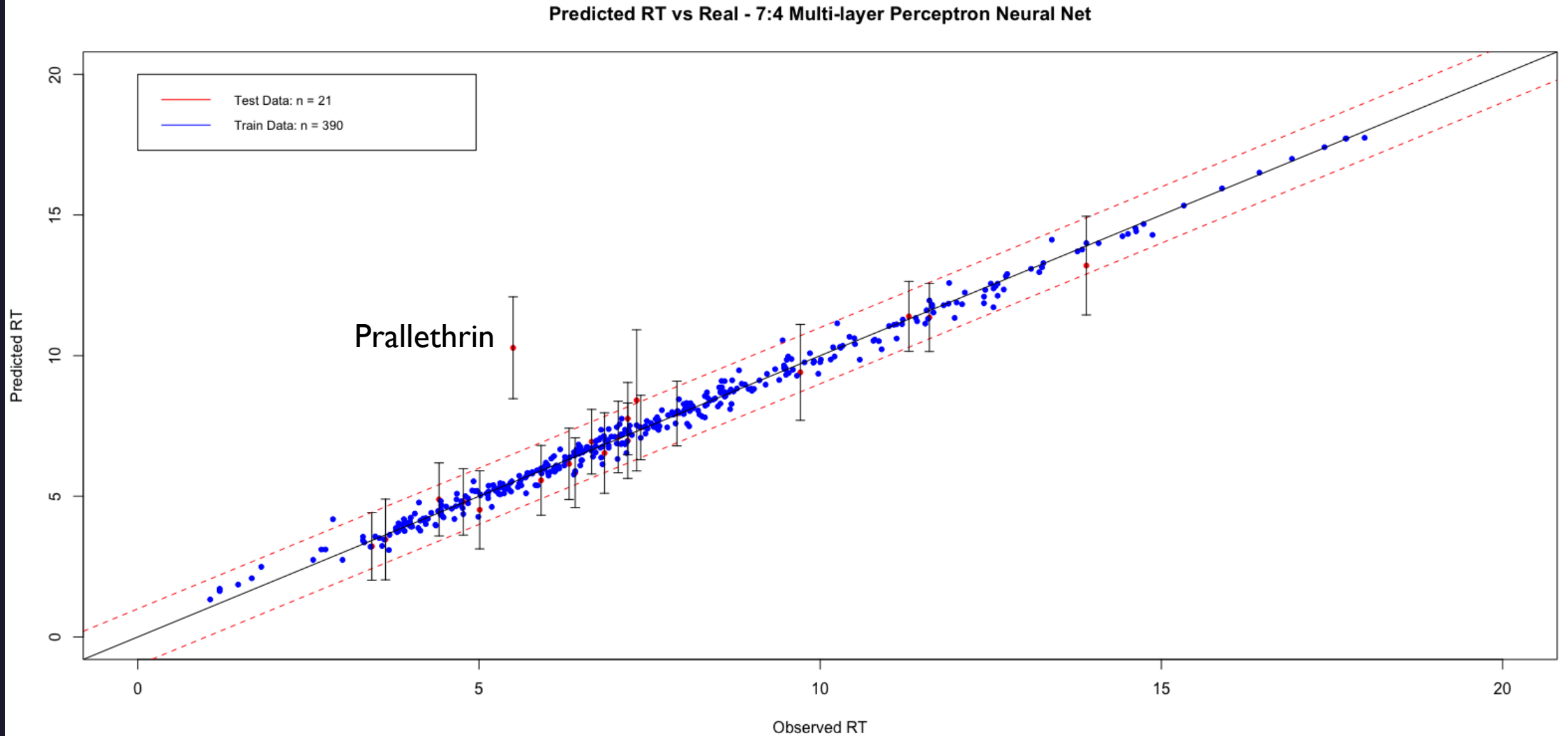
Overly trained for PFAS (80+ PFAS) but includes pesticides, CECs

Did initial Cross-Validation from 10% to 90% training, run with 5 replicates of randomized data. RMSE minimization confirmed.

Dominant QSARs that predict RT are usual suspects (eg: AlogP)

Another tool – to help increase ID confidence

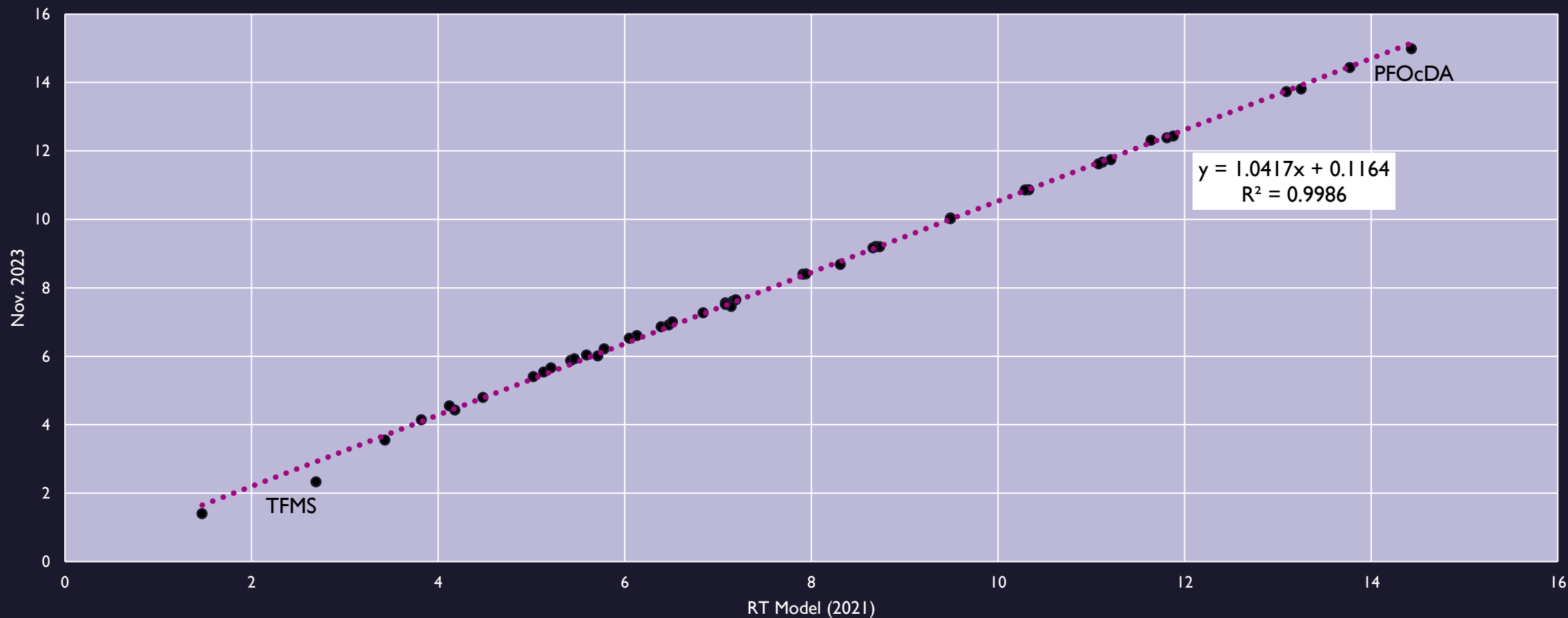
RT Model – Observed/Predicted



RT Model (2021) vs. 2023

47 compounds – mostly PFAS

Old Model vs 2023



What's Required for an RT Model

- Structures of all molecules (SMILES)
- Specific Quantum Structure Activity Relationship (QSAR) calculations. Examples: ALogP, nHB Donors, nHB Acceptors, Elemental counts, polarizability)
- Train the model
- Validate the Model using various test/train ratios. Determine if RMSE goes down as train ratio increases
- Store the model
- Use model on Feature candidate lists (one feature, many possible chemicals)



Molecular Networks To Explore Feature Relationships

Molecular Networks (Thermo CD)

massspecinfo/ x PubChem x Search ChemS... x A statistical ov... x PDR: Database x Mass Spectral x Table Explorer x ERD.png (4322 x Environmental x Determining th... x 20240706_202... x

File D:/CompDisc/20240627/20240706_20240627/index.html

Clock NetSuite - Custome... Water Boards FTP PFAS Broad Spectru... Qualtrax Thermo Fisher Dow... Ideagen Quality Ma... Adobe Acrobat DIMSpec for PFAS ... Mass Spectral Matc... Inventory Supplies...

SEARCH

Compound

Transformation

GRAPH INFO

FILTERS

Show Unknown Orphans

Show Identified Orphans

Require Transformation

Require MSn

THRESHOLDS

Score

Coverage

Matched Fragments

CLUSTERS

Node Links

Cluster Size

Isolation Depth

NODE STYLE

Show Confidence

Show Pie Chart

Show Structure

Show Name

Show Formula

Show Mass

Show RT

Triazines+TPs

PFCAs

PFSulfonamides

PFSulfonates

PFPeS

Formula: C5 H F11 O3 S

RT: 6.156 min

MW: 349.94713 Da

ΔMass: 0.00005 Da

Max. Area: 167,156,796

Fragments: 3

CCCC(F)(F)S(=O)(=O)F

SAMPLES

Sample	Count	Percentage
Blank2, Sample	18,474	0.01 %
BS, Sample	18,616,185	7.61 %
BSD, Sample	20,033,194	8.19 %
C4E2360-03, Sample	25,141	0.01 %
C4E2865-05, Sample	408,419	0.17 %
C4E2865-06, Sample	410,872	0.17 %
C4E3034-03, Sample	22,248	0.01 %
C4E3499-03, Sample	100,806	0.04 %
C4E3499-06, Sample	37,159	0.02 %
C4E3501-03, Sample

© Copyright 2018-2024 Thermo Fisher Scientific Inc. | Generated by Compound Discoverer | Powered by D3.js



Informed DeNovo Molecular Formula Generation

Molecular Formula Generation

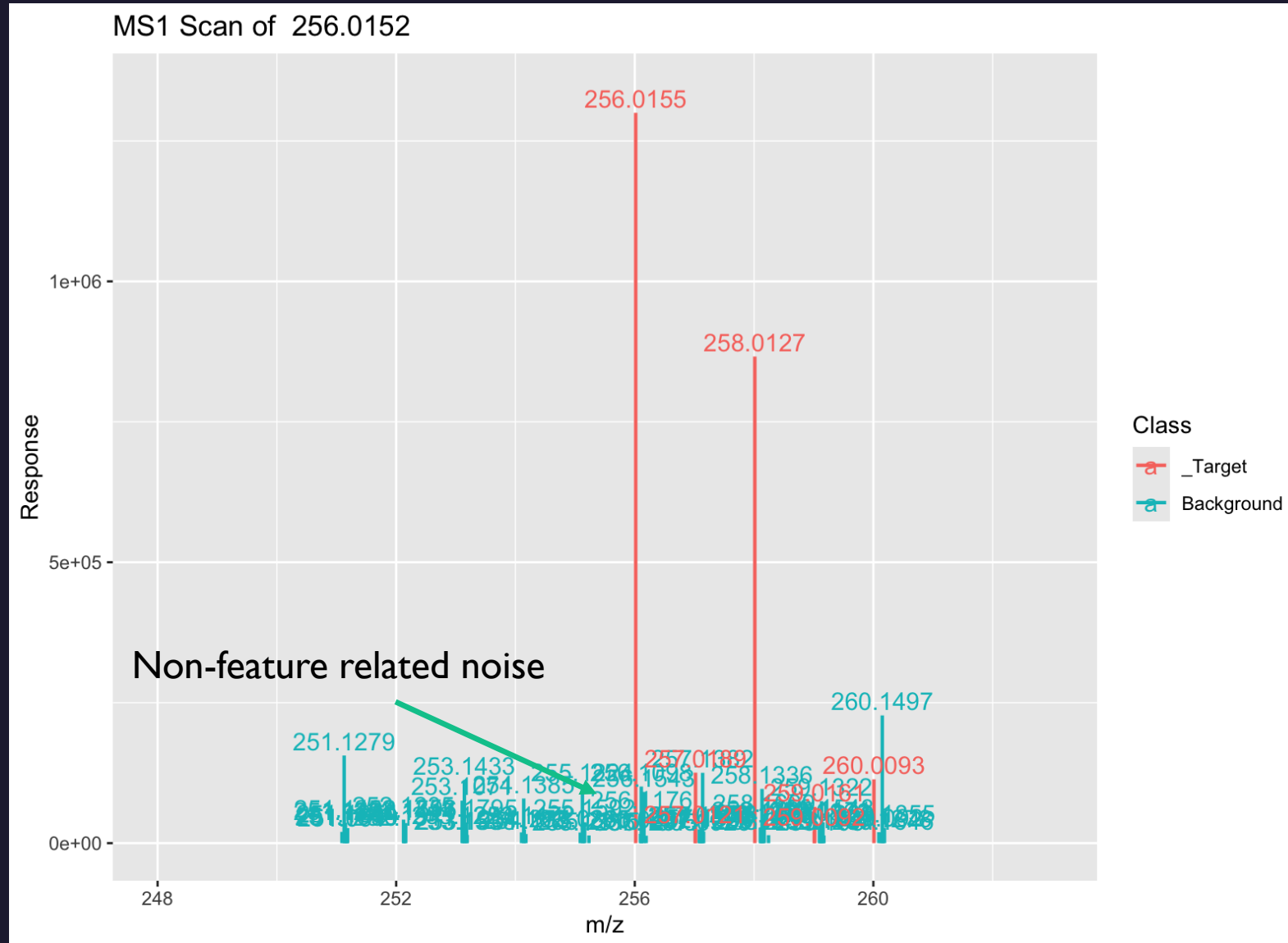
- DeNovo – “from the beginning”
- Decomposition computationally expensive / time consuming (last resort)
- Often produces more junk than useful information

	mz	MF	charge	RdisopScore	unsat	parity	error	nrule	ppm	SENIOR3	HtoC	NtoC
1	371.1009	C6H31O7N2S4	1	9.183108e-02	-7.5 e		-2.288724e-05	Valid	-0.0616739	3	5.1666667	0.3333333
5	371.1006	C13H23O10S	1	5.182191e-02	2.5 e		-2.393812e-04	Valid	-0.6450570	9	1.7692308	0.0000000
7	371.1006	C11H11N14S	1	5.010691e-02	13.5 e		-2.500732e-04	Valid	-0.6738686	45	1.0000000	1.2727272
8	371.1006	H15O010N2S	1	4.081205e-02	-73.0 e		-3.119322e-04	Valid	-0.8405591	-140	Inf	Inf
10	371.1012	C8H146O8	1	4.587981e-02	-64.0 e		3.399638e-04	Valid	0.9160953	-128	18.2500000	0.0000000
12	371.1005	C5H19O13N6	1	2.704001e-02	-0.5 e		-4.223192e-04	Valid	-1.1380173	5	3.8000000	1.2000000
14	371.1005	C3H7O3N2O	1	2.590080e-02	10.5 e		-4.330112e-04	Valid	-1.1668289	41	2.3333333	6.6666666
17	371.1013	C6H23O8N6S2	1	3.147397e-02	-1.5 e		4.460708e-04	Valid	1.2020202	11	3.8333333	1.0000000
21	371.1014	CH154O5N2S3	1	2.005891e-02	-74.0 e		5.564578e-04	Valid	1.4994784	-134	154.0000000	2.0000000
23	371.1015	C14H27O5S3	1	1.442732e-02	1.5 e		6.290088e-04	Valid	1.6949805	15	1.9285714	0.0000000
26	371.1017	C7H35O2N2S6	1	4.583770e-03	-8.5 e		8.455028e-04	Valid	2.2783636	9	5.0000000	0.2857142
27	371.1000	C5H27O12N2S2	1	2.348783e-03	-6.5 e		-8.912772e-04	Valid	-2.4017114	-3	5.4000000	0.4000000
29	371.1000	C3H15O2N16S2	1	2.191293e-03	4.5 e		-9.019692e-04	Valid	-2.4305230	33	5.0000000	5.3333333
31	371.1018	C6H15O9N10	1	3.007650e-03	4.5 e		9.150288e-04	Valid	2.4657143	19	2.5000000	1.6666666
34	371.0999	C18H11O2N8	1	1.755212e-03	17.5 e		-9.35252e-04	Valid	-2.5209459	43	0.6111111	0.4444444
37	371.0999	C5H138O2N10	1	1.063501e-03	-58.0 e		-1.008076e-03	Valid	-2.7164480	-106	27.6000000	2.0000000
39	371.1019	CH146O6N6S	1	1.459270e-03	-68.0 e		1.025416e-03	Valid	2.7631725	-126	146.0000000	6.0000000

Showing 1 to 17 of 64 entries. 22 total columns

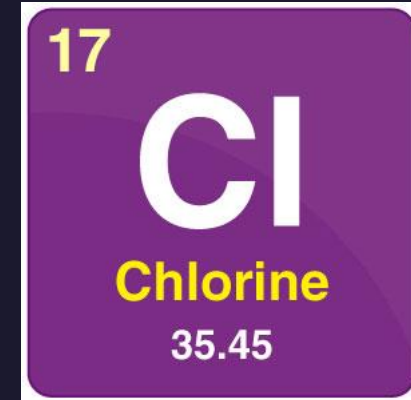
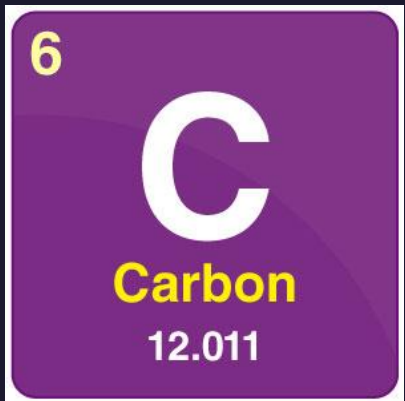
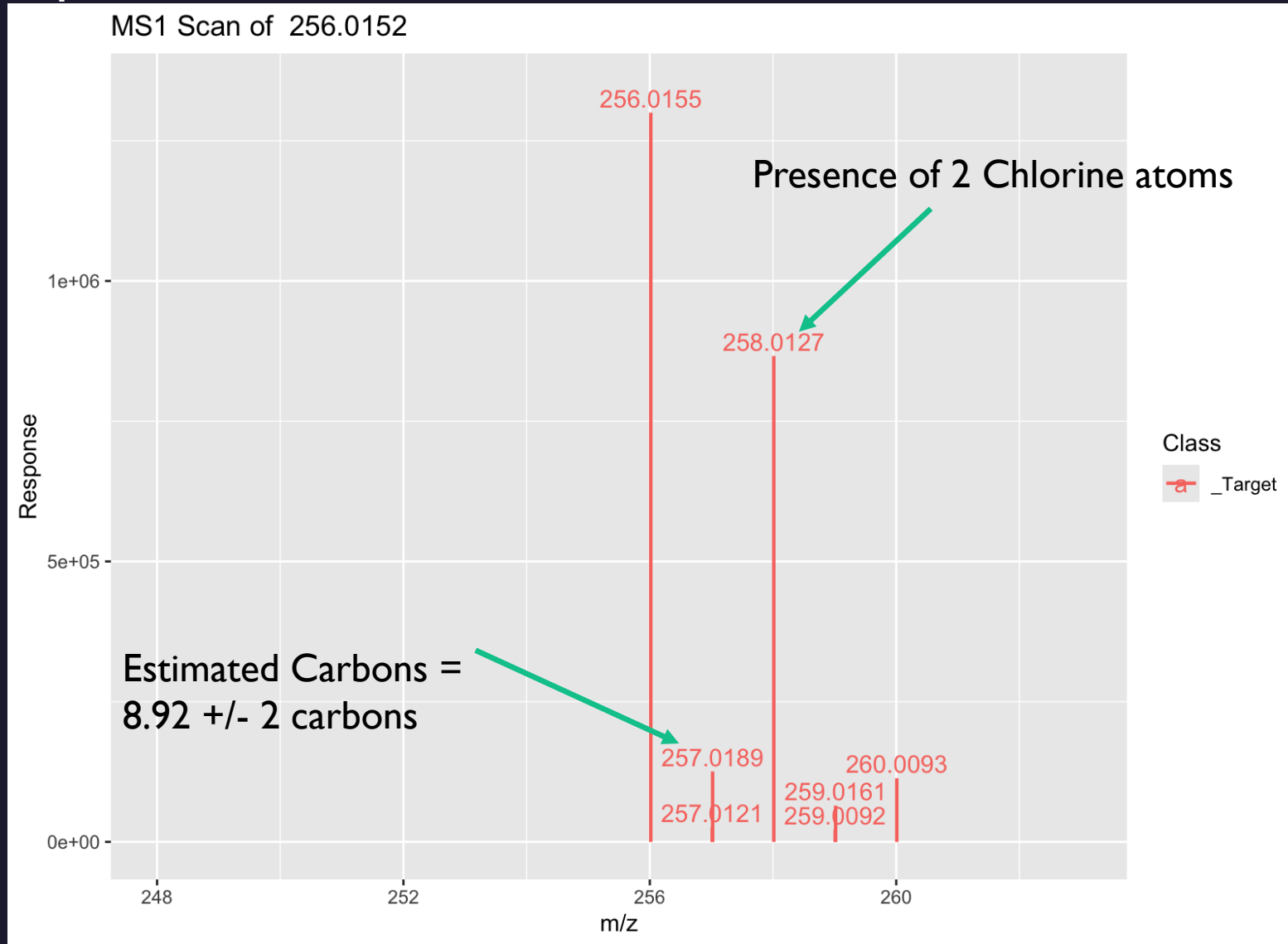
Molecular Formula Elemental Bounds

Raw MS1 Spectrum



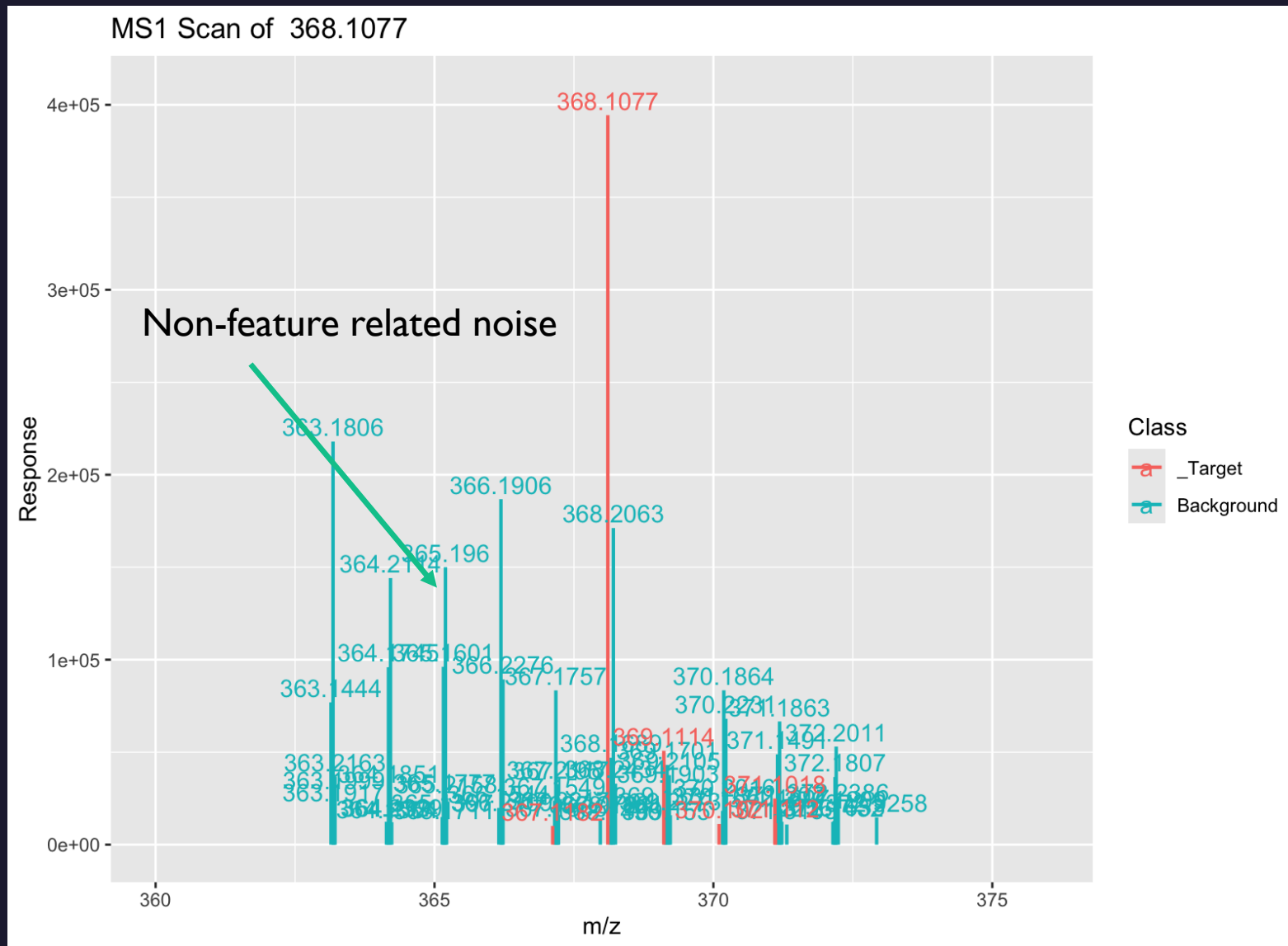
Molecular Formula Elemental Bounds

Filtered MS1 Spectrum



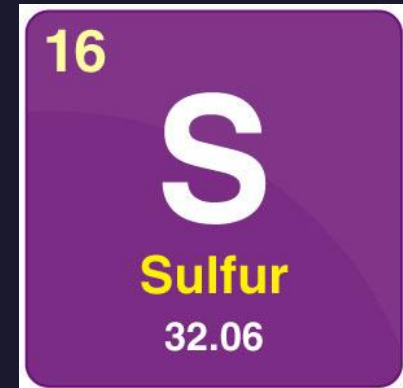
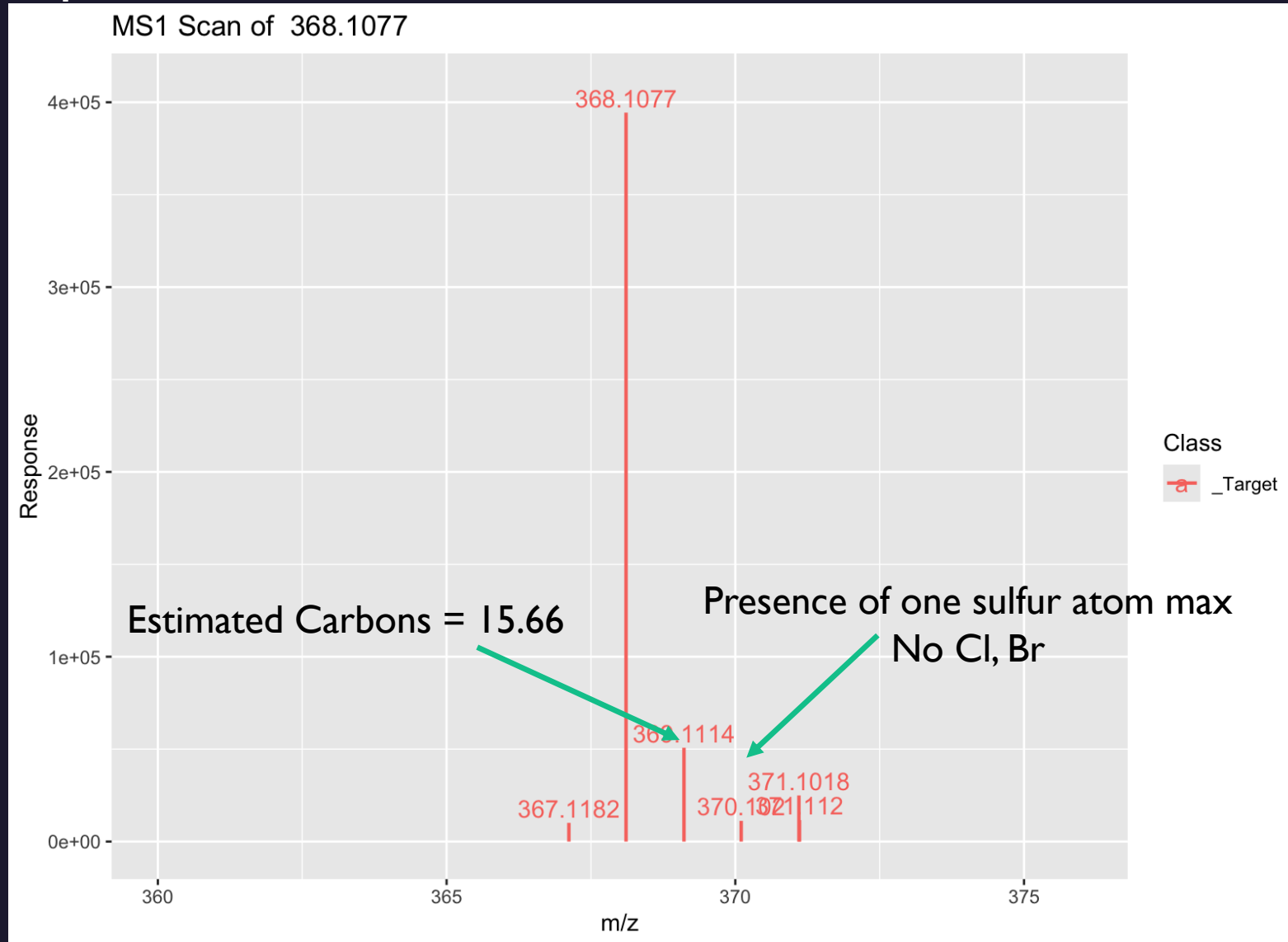
Molecular Formula Elemental Bounds

Raw MS1 Spectrum



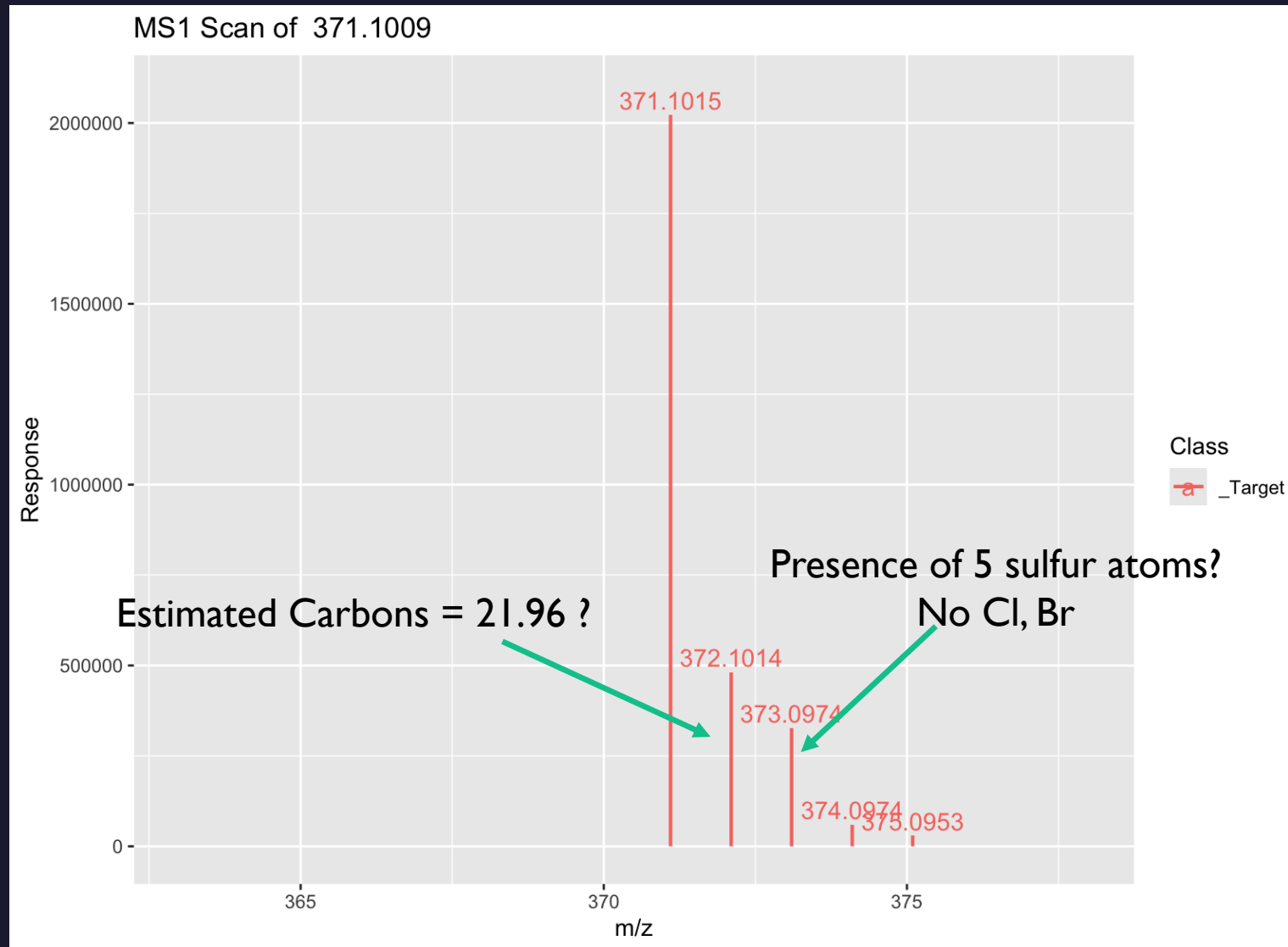
Molecular Formula Elemental Bounds

Filtered MS1 Spectrum



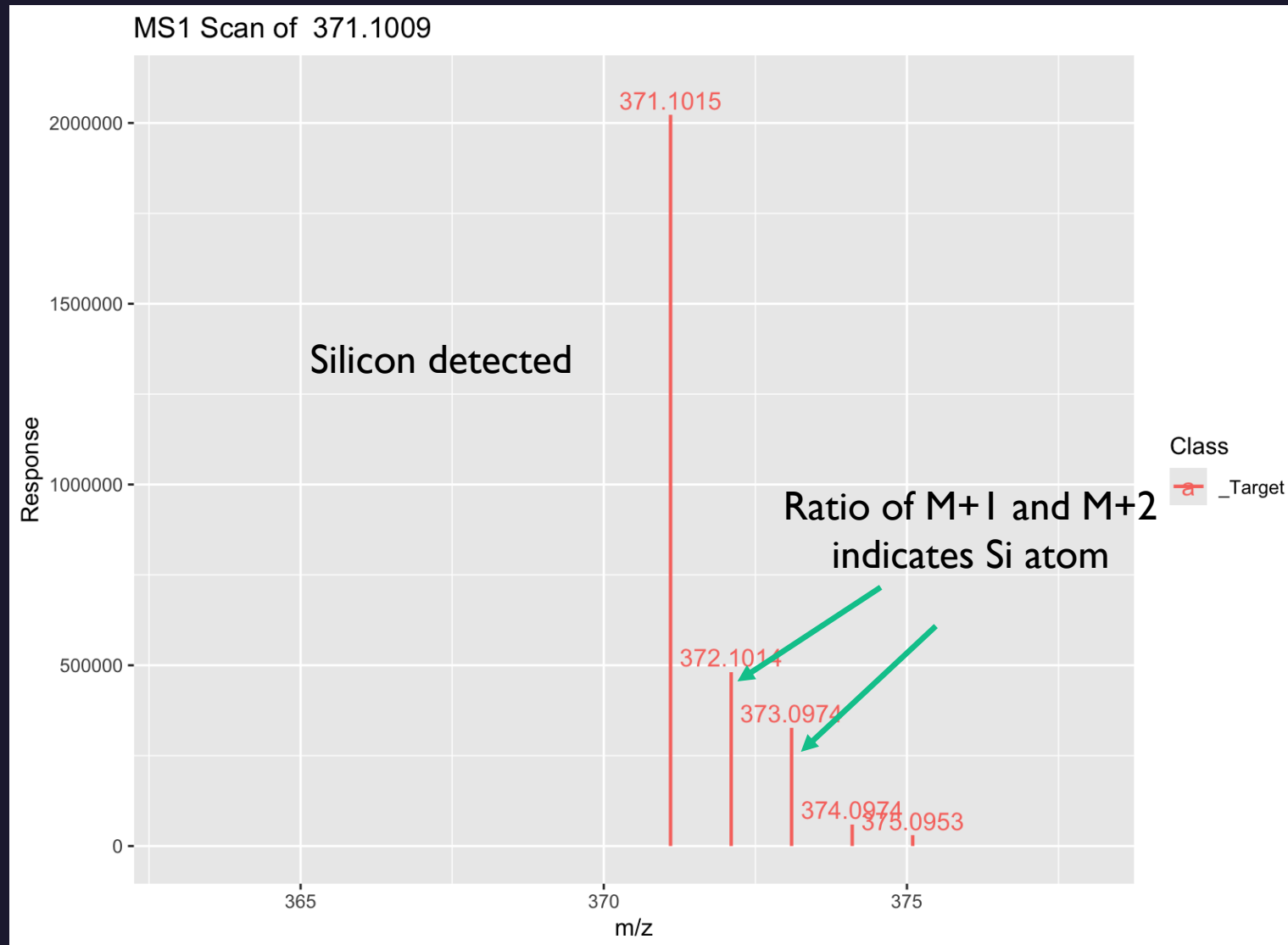
Molecular Formula Elemental Bounds

Filtered MS1 Spectrum



Molecular Formula Elemental Bounds

Filtered MS1 Spectrum



Presence of Silicon
throws off M/Carbon
ratio and easily
identified

Conclusions

- In the absence of MS/MS data, MSI spectra can be interrogated for more information
- Using High Res Mass Spec intrinsic values like m/z , mass defect, and $[^{13}\text{C}]$ ratios can be used to calculate “PFAS-ness” and can be prioritized.
- Retention time prediction models provide an orthogonal technique to confirm or reject potential features
- Plotting data with intrinsically determined values is very useful for PFAS prioritization (Kendrick or Kaufmann)
- Informed Molecular Formula predictions that properly decompose MSI spectra with elemental bounds more effective than agnostic ones

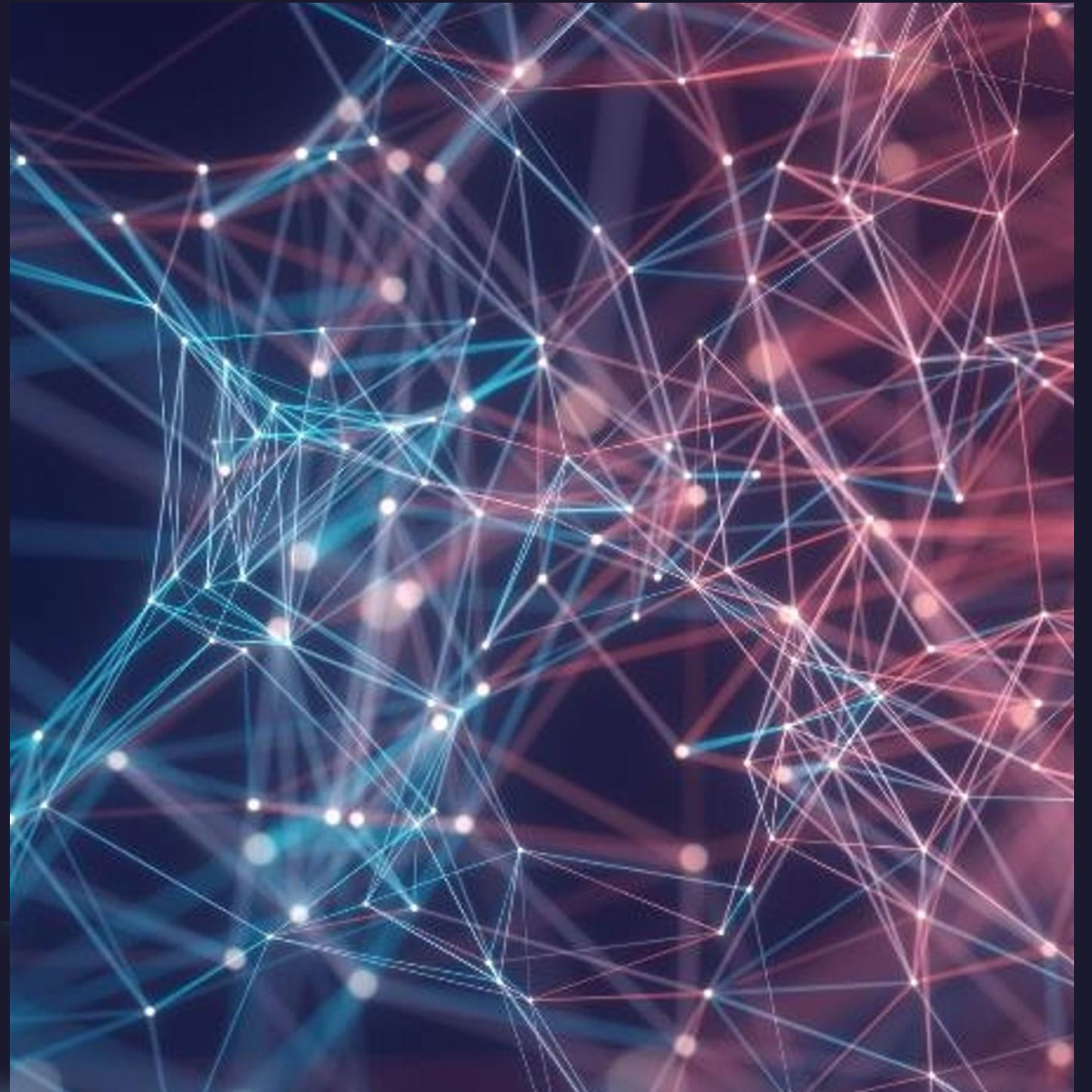
Thank you

David Schiessel

951 289 5278

dschiessel@babcocklabs.com

www.babcocklabs.com



References

- Bugsel, B.; Zwiener, C. LC-MS Screening of Poly- and Perfluoroalkyl Substances in Contaminated Soil by Kendrick Mass Analysis. *Anal Bioanal Chem* **2020**, *412* (20), 4797–4805. <https://doi.org/10.1007/s00216-019-02358-0>.
- Kaufmann, A.; Butcher, P.; Maden, K.; Walker, S.; Widmer, M. Simplifying Nontargeted Analysis of PFAS in Complex Food Matrixes. *Journal of AOAC INTERNATIONAL* **2022**, *105* (5), 1280–1287. <https://doi.org/10.1093/jaoacint/qsac071>.
- Zweigle, J.; Bugsel, B.; Zwiener, C. Efficient PFAS Prioritization in Non-Target HRMS Data: Systematic Evaluation of the Novel MD/C-m/C Approach. *Anal Bioanal Chem* **2023**, *415* (10), 1791–1801. <https://doi.org/10.1007/s00216-023-04601-1>.