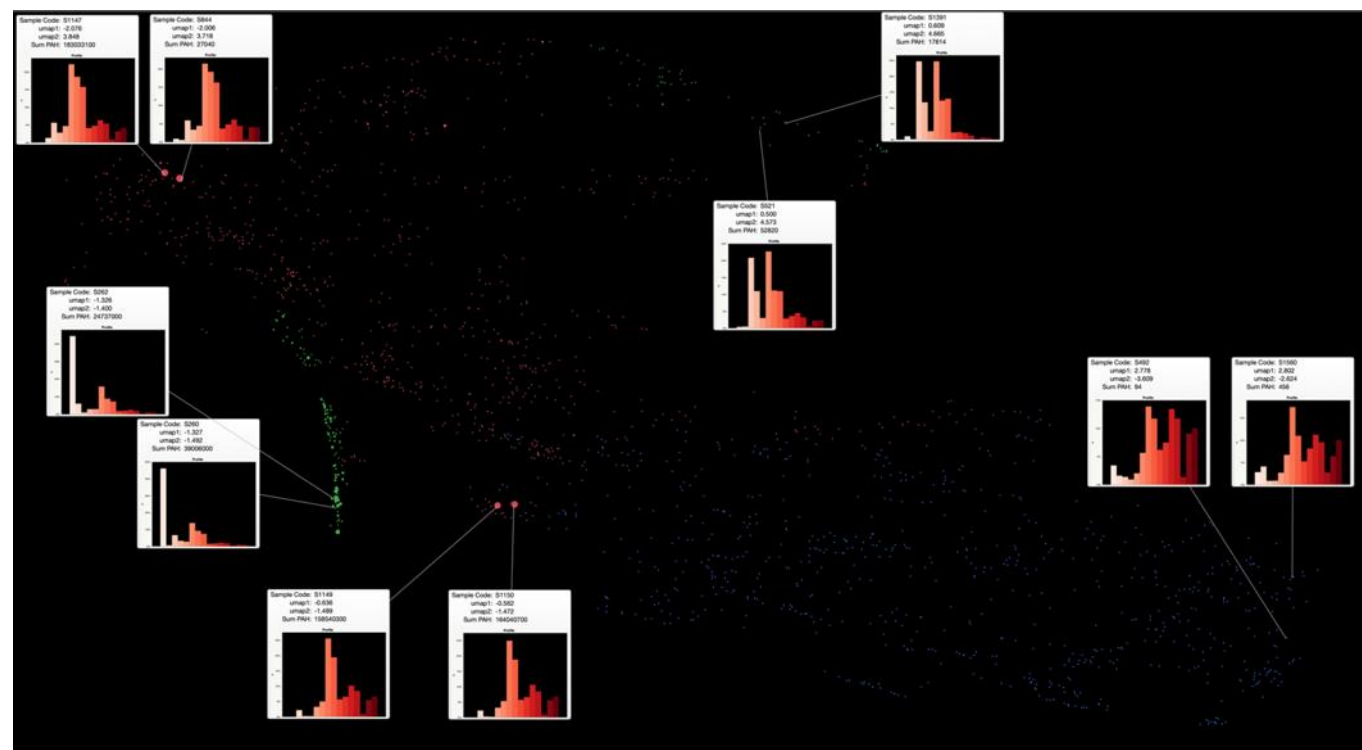




chemistry  
matters

*Making chemistry data meaningful*



# Managing data and data discovery for large scale Superfund sites and releases

*Court Sandau, PhD, PChem, FRSC and Paul Fuellbrandt, BSc, PAg, PMP*

October 12, 2023



# Companies



## Chemistry Matters Inc.

- **Niche chemistry consulting company specializing in environmental forensics, geoforensics, biomonitoring and arson investigations**
- Big environmental datasets have become our specialty to merge chemistry understanding with statistical interpretation and science communication
- Virtual office-based business

## Statvis Analytics Inc.

- Software development company specializing in data and statistical visualization software customized for environmental industry
  - **Environmental Data Intelligence System (EDIS)<sup>TM</sup>**
- Turning complex data and statistics into easy-to-understand graphics and pictures
- Helping people communicate science

**ALBERTA TEST SITE** Up to date: 40 samples from Sep 30, 2008 - Dec 05, 2018. Edit Custom Site Guidelines Edit Generic Site Guidelines Upload CSV

**Geospatial Data Exploration**

Search Samples:  **FULLSCREEN** **HIDE TABLE**

SAMPLE NAME	THALLIUM	BENZENE	TOLUENE	ETHYLBENZ.	XYLENES	F1	F2	F3	F4	NAPHTHA...	ACENAPH...	FLUORENE	ANTHRAC...	PHENANT...
BH41	-	0.009	0.05	0.01	0.05	10	10	147	504	0.01	0.05	0.05	0.004	0.05
BH41	0.31	0.005	0.24	0.05	0.05	10	20	190	425	0.01	0.05	0.05	0.004	0.05
BH40	-	0.005	0.02	0.005	0.03	10	50	50	100	-	-	-	-	-
BH40	0.15	0.005	0.02	0.005	0.03	10	50	50	100	-	-	-	-	-
BH41	-	0.005	0.02	0.005	0.03	10	50	107	462	-	-	-	-	-
BH41	0.005	0.02	0.005	0.03	10	50	50	100	100	-	-	-	-	-
BH42	0.19	0.005	0.02	0.01	0.05	10	11	40	52	0.01	0.05	0.05	0.004	0.05
BH42	-	0.01	0.02	0.01	0.04	100	100	100	140	-	-	-	-	-
BH42	-	0.005	0.02	0.005	0.03	10	50	122	500	-	-	-	-	-
BH42	-	0.005	0.02	0.005	0.03	10	50	50	100	-	-	-	-	-

**Your Site Soil Salinity**

Area Name	Location	Status	Depth	Sample Date	Latitude	Longitude	Cation	Chloride	Electrical Conductivity	Magnesium	pH	Potassium	Sulfate	Sodium	Sulfate + Sodium	Sulfur	Temperature of Sample
Area 1	Location 1	Status 1	Depth 1	Sample Date 1	Latitude 1	Longitude 1	Cation 1	Chloride 1	Electrical Conductivity 1	Magnesium 1	pH 1	Potassium 1	Sulfate 1	Sodium 1	Sulfate + Sodium 1	Sulfur 1	Temperature of Sample 1

**Automated Tabulation**

**ALBERTA TEST SITE** Up to date: 40 samples from Sep 30, 2008 - Dec 05, 2018. Edit Custom Site Guidelines Edit Generic Site Guidelines Upload CSV

**Annotation and Markup**

Search Samples:  **FULLSCREEN** **HIDE TABLE**

**ALBERTA TEST SITE** Up to date: 40 samples from Sep 30, 2008 - Dec 05, 2018. Edit Custom Site Guidelines Edit Generic Site Guidelines Upload CSV

**Trending**

Search Samples:  **FULLSCREEN** **HIDE TABLE**

**ALBERTA TEST SITE** Up to date: 40 samples from Sep 30, 2008 - Dec 05, 2018. Edit Custom Site Guidelines Edit Generic Site Guidelines Upload CSV

**Salinity Fingerprinting**

Search Samples:  **FULLSCREEN** **HIDE TABLE**

**ALBERTA TEST SITE** Up to date: 40 samples from Sep 30, 2008 - Dec 05, 2018. Edit Custom Site Guidelines Edit Generic Site Guidelines Upload CSV

**PAH Fingerprinting**

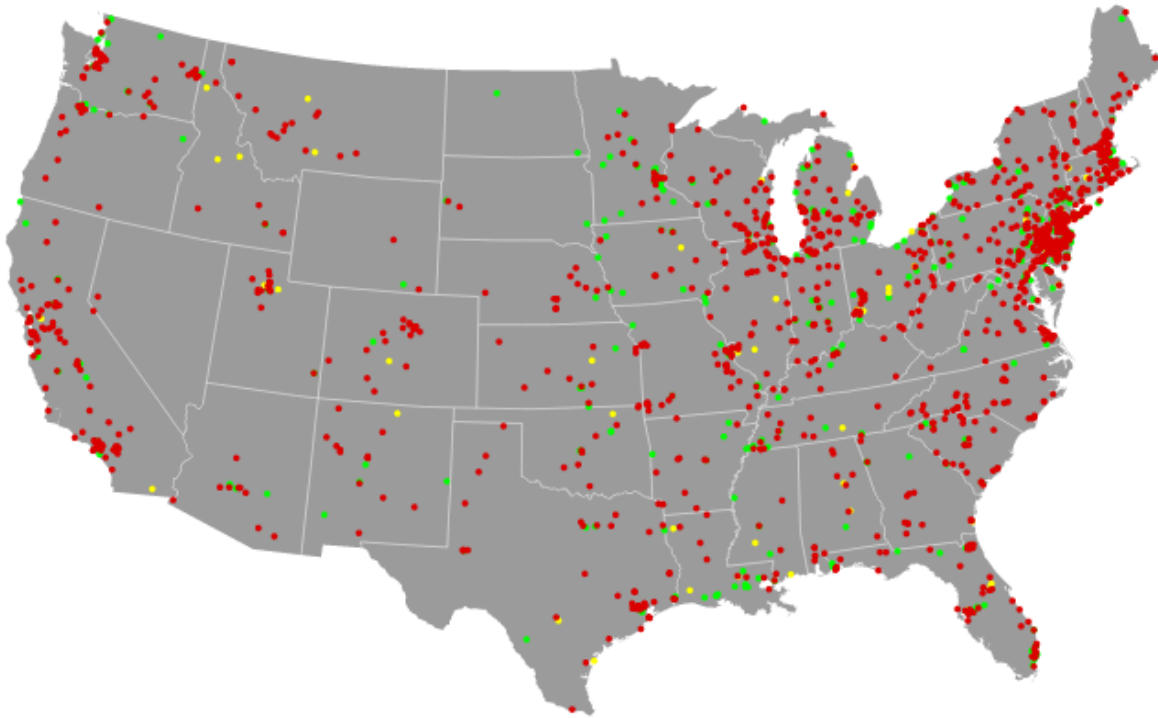
Search Samples:  **FULLSCREEN** **HIDE TABLE**

Get in touch at [info@Statvis.com](mailto:info@Statvis.com) to learn more about the only

**Environmental Data Intelligence System™**

**statvis** Trust Through Data

# United States Superfund Sites



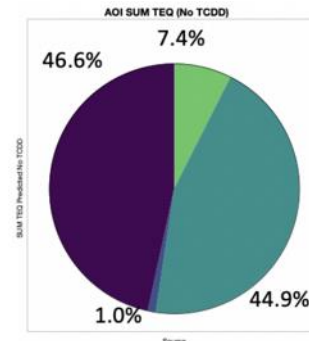
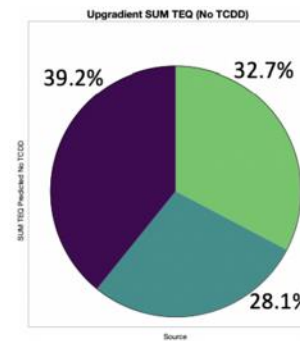
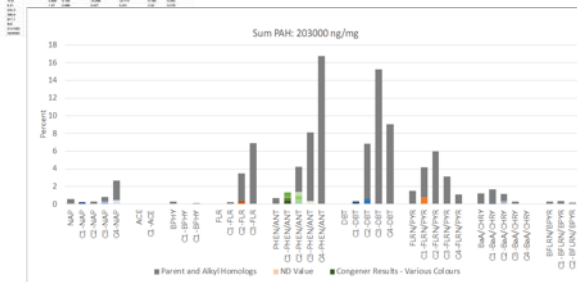
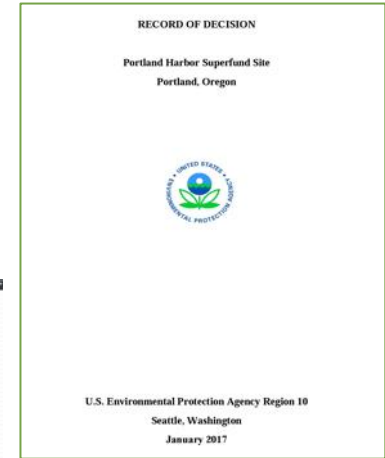
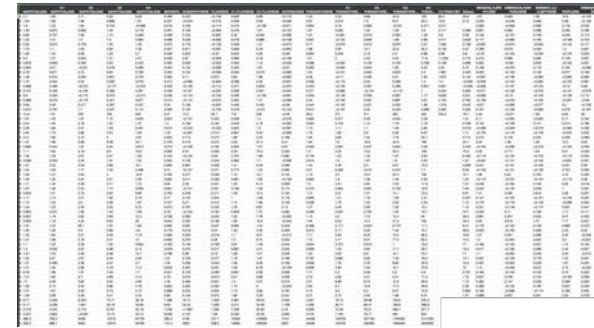
- Sites under EPA control due to hazardous materials that require remediation
- ~1400 sites across US
- All different types of contaminants
  - > 600 chemicals
  - Pd (43%), trichloroethylene (42%), Cr (35%), benzene (34%), PERC (28%), As (28%) and toluene (27%)
  - PCBs, PCDD/Fs, PAHs, PFAS/PFOA
- Some sites have 100+ years of industrial operations

# Superfund Sites – Remediation Estimates



# Typical Superfund Site Project Workflow

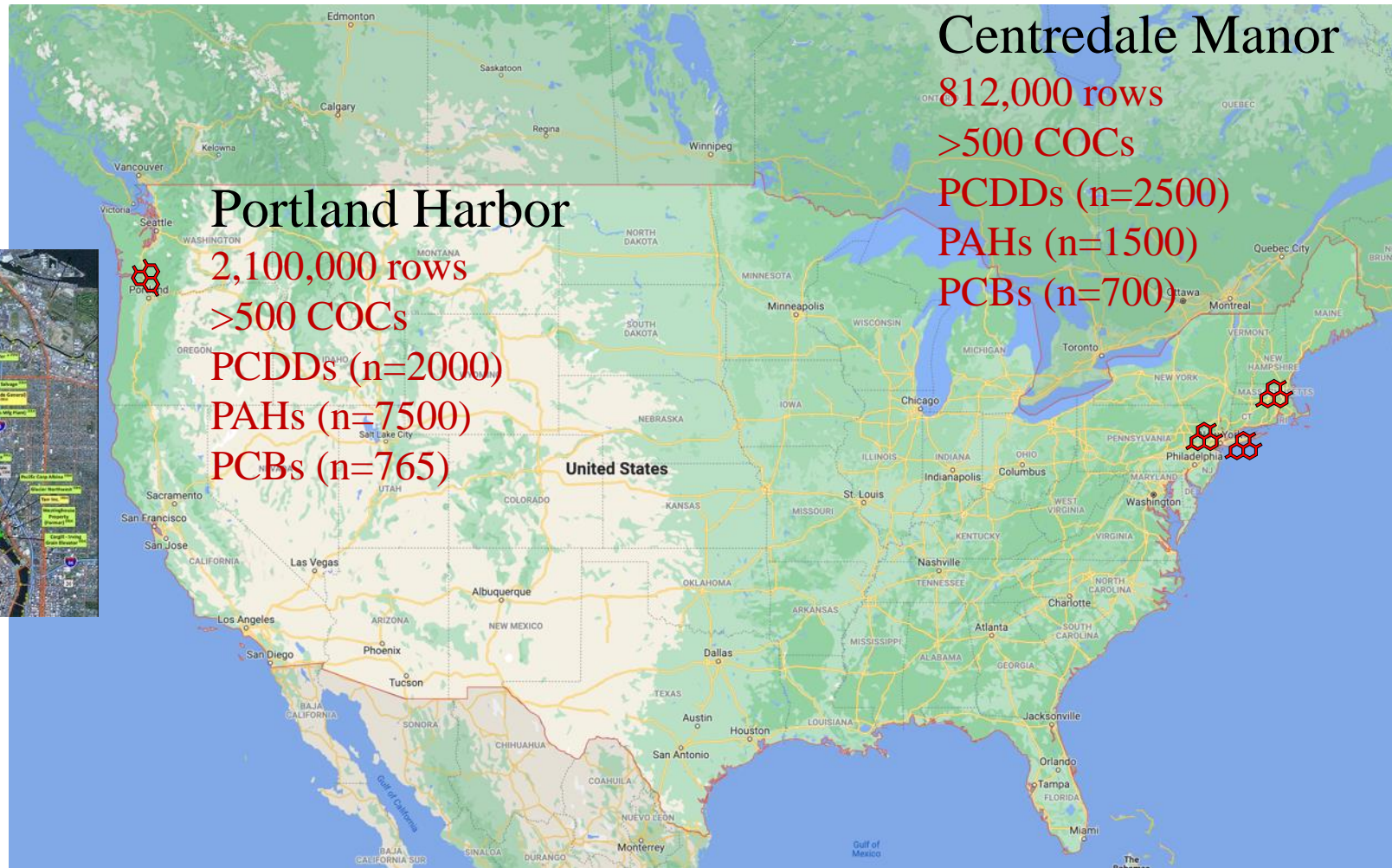
- Review historical reports and data
- Compile and wrangle historical data
- Process data through interactive data workflow
  - Receptor modelling/source apportionment
- Report and present findings
- Depositions, rebuttal reports and testifying



# Data Wrangling

- Data import and merge
- Dataset familiarization
- MAR – Missing at Random values
- Duplicates, replicates, repeat analytical results
- Data tidying (unit consolidation, co-elutions)
- Detection Limits –  $\frac{1}{2}$  DL, KM, Multiplicative Simple Replacement, Multiple Lognormal Replacement (Maximum Likelihood or Robust Estimates)
- Data screening
  - minimum detects per analyte
  - minimum detects per sample
- Fundamental data coding – up/downstream, labs, investigations, locations
- Initial data review

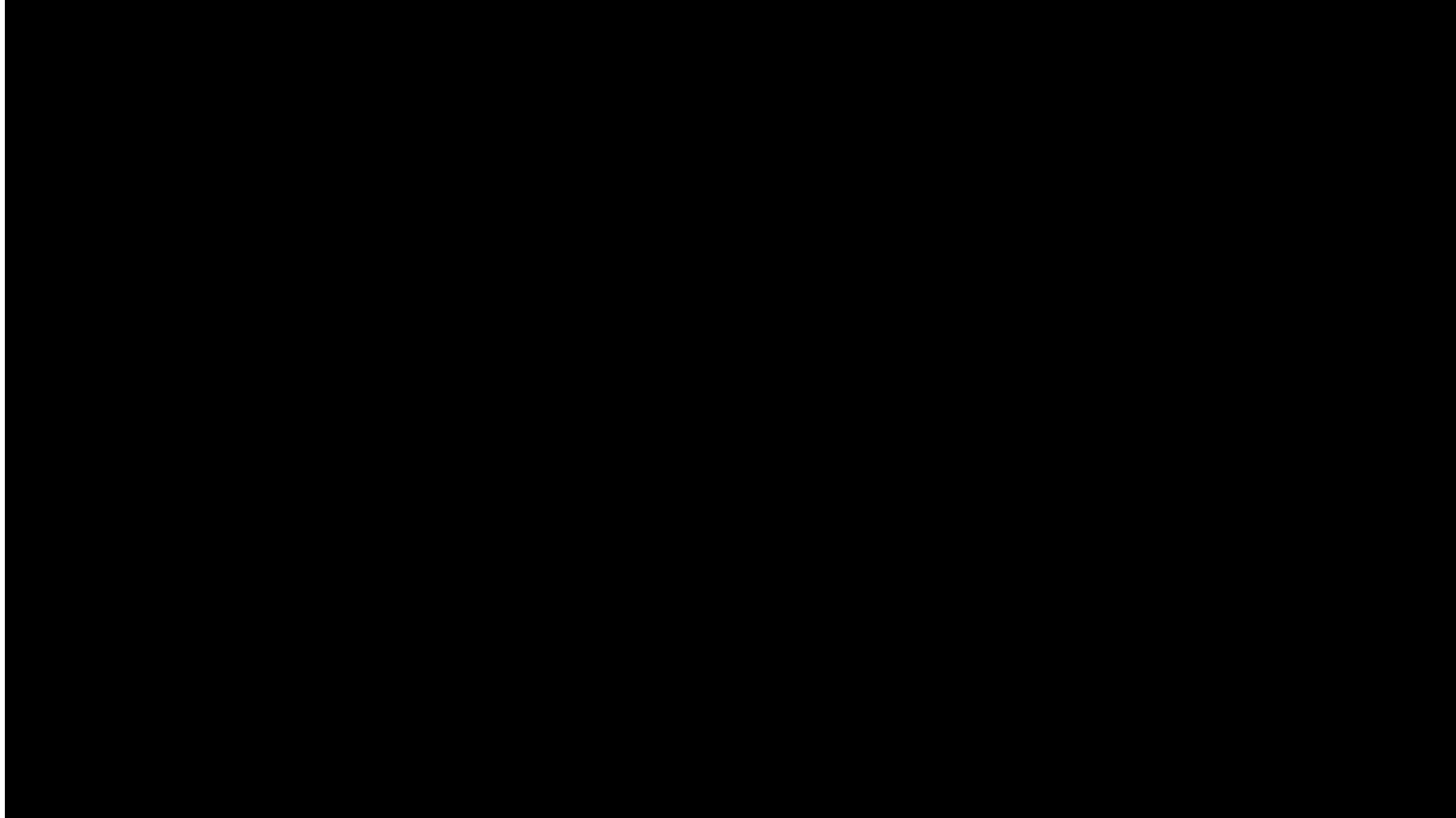
# Superfund Sites – Lots of Data to Work With







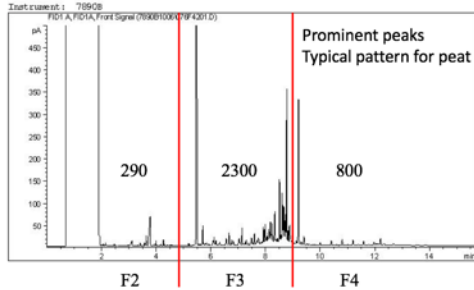
# Portland Harbour Superfund Site - PAH Data



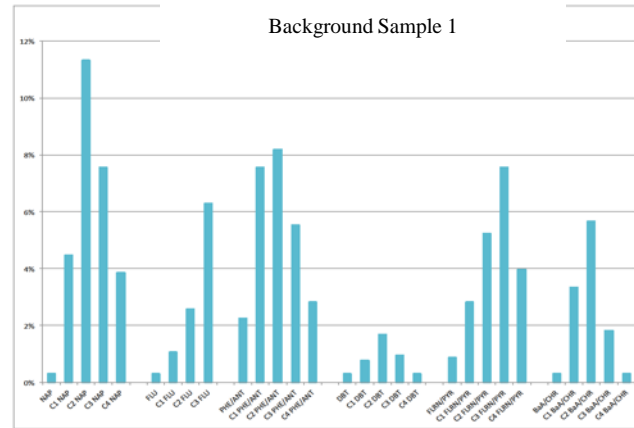
# Requirements for Data Discovery

- **Imperative to have data interactivity to discover what the data is telling you and to develop knowledge**
  - *The goal is 'data mastery'*
- This can only be accomplished with commercial statistical software, like SAS JMP
  - Can be expensive
  - Licensed per user
  - Hard to share with clients or pass along data knowledge
- We are currently developing these data discovery tools for everyday sites in Statvis

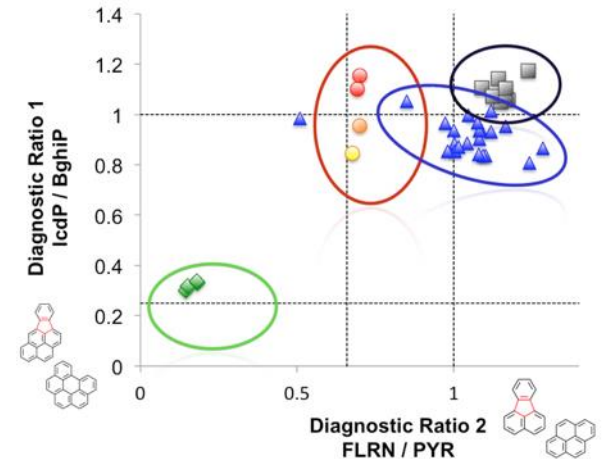
# Progression of Chemical Fingerprinting in Environmental Forensics



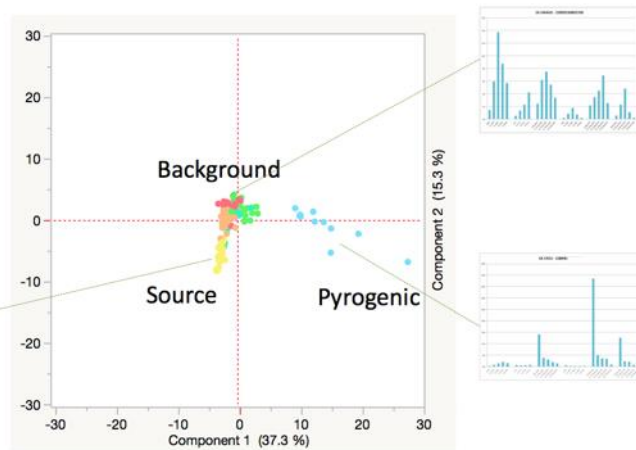
Chromatogram Analysis



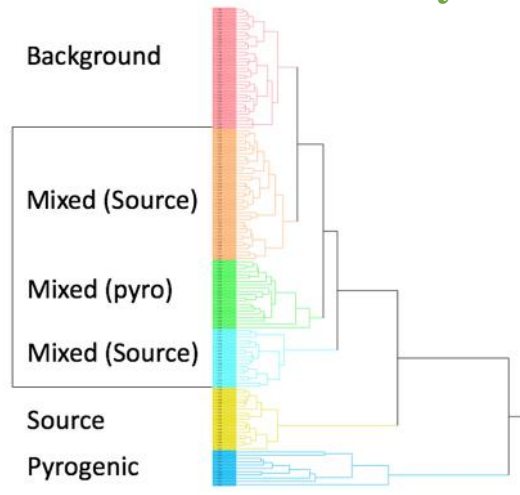
Individual Sample Pattern Analysis



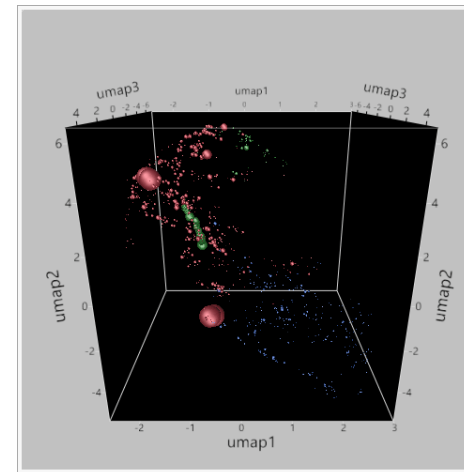
Double Ratio Plots / Diagnostic Ratios



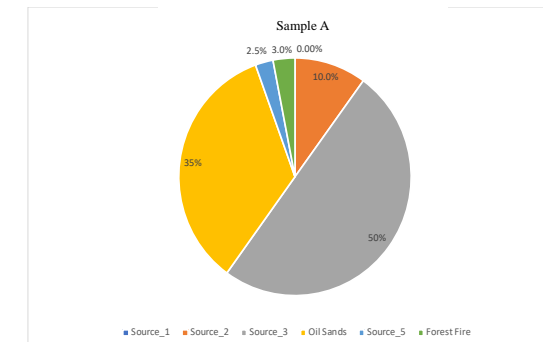
PCA



HCA



UMAP

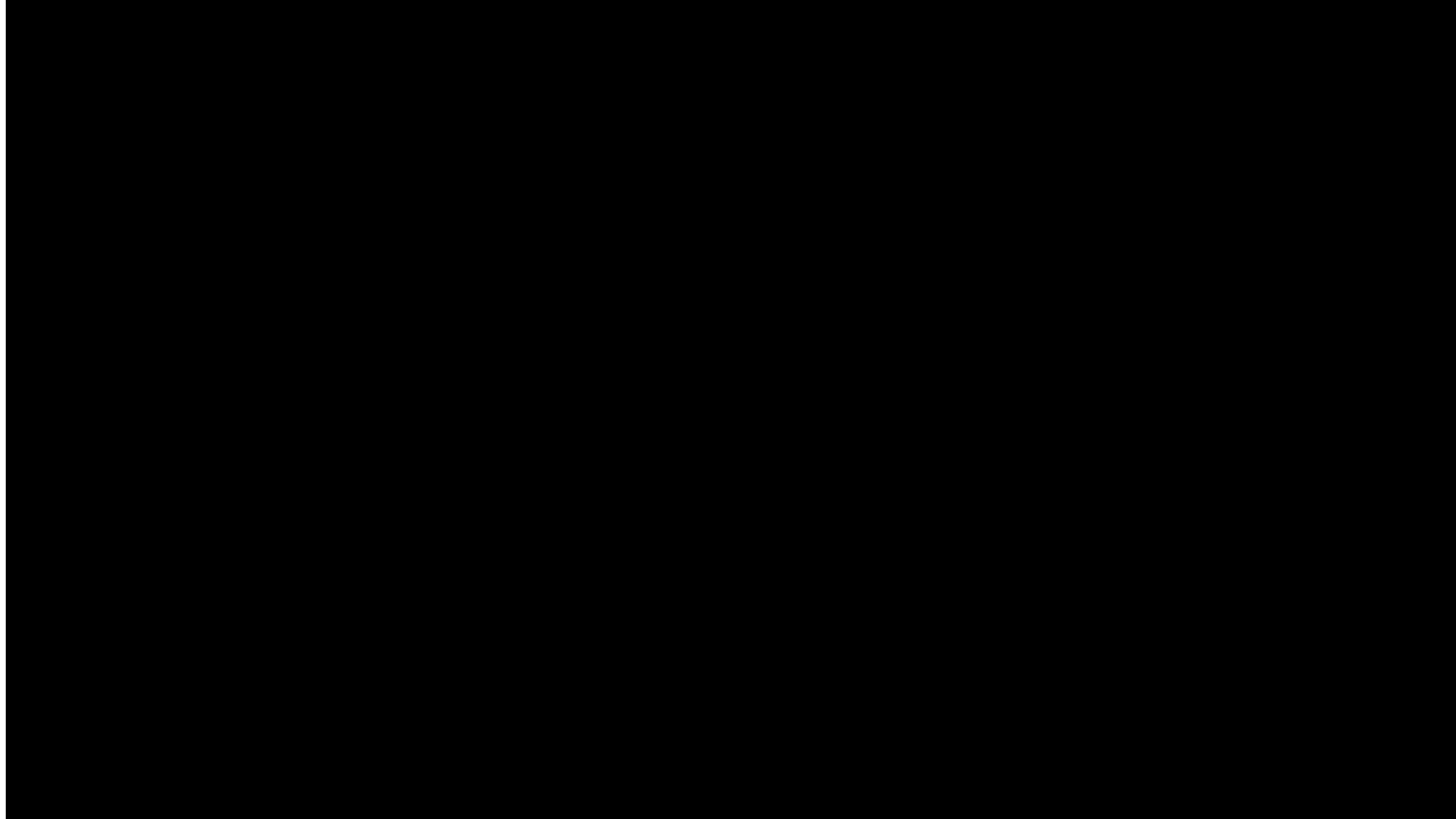


Source Apportionment



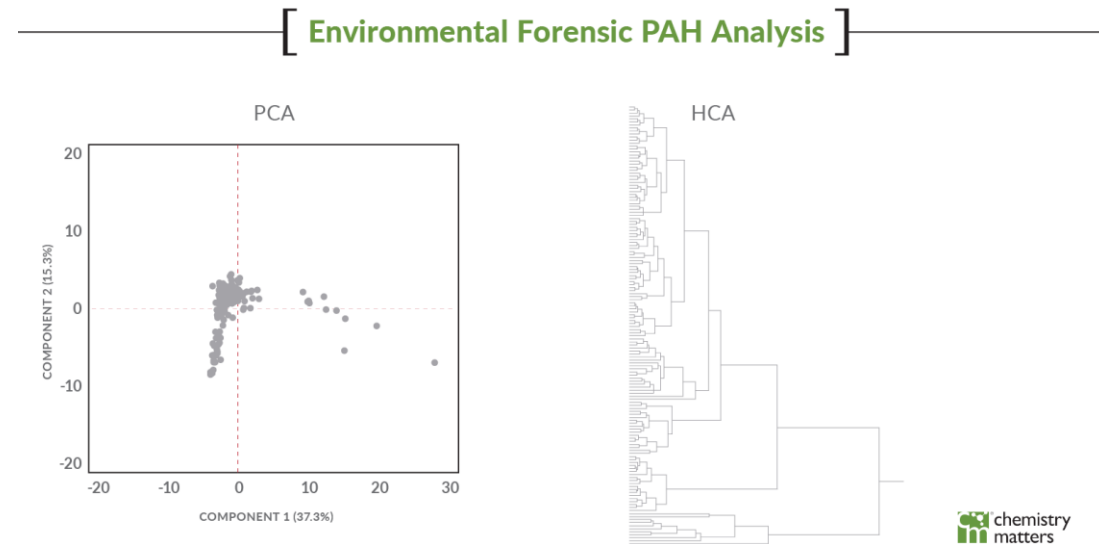


# Portland Harbour PAH Fingerprinting



# Integration and interactivity is the key to scientific communication

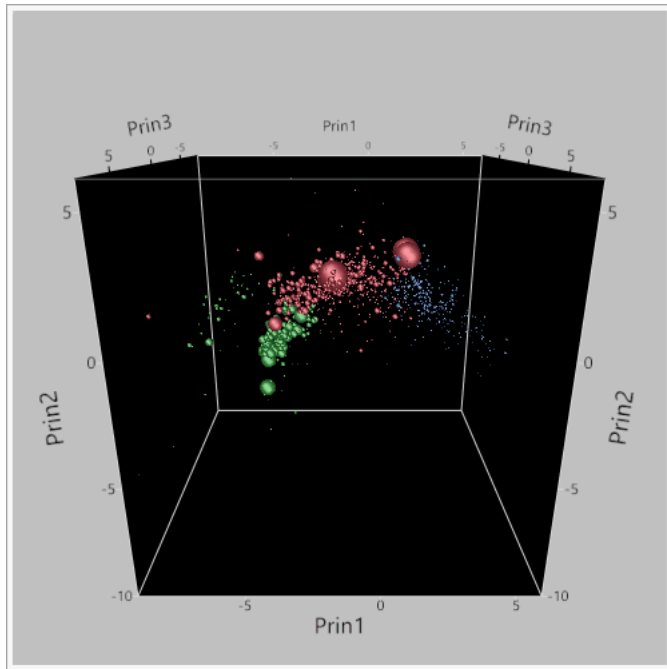
- To build trust and buy-in, data mastery is essential
  - Data at your fingertips
  - Able to ‘show’ the answer
    - Seeing is believing, even with complex data
- Integrated statistical analysis coupled with geospatial analysis with visualization allows learning while you are analyzing
- An interactive data portal that allows the user to interact and demonstrate results all while communicating and explaining its meaning



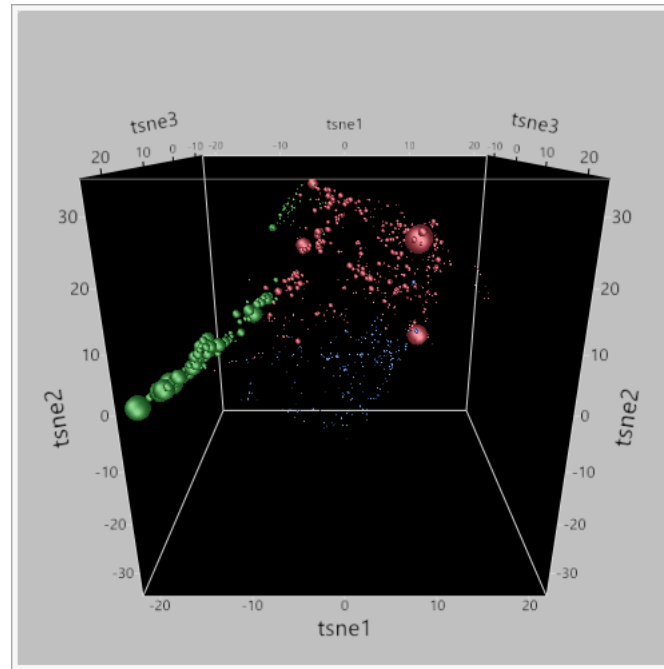
# Different Techniques for Dimensionality Reduction

## Data Visualizations

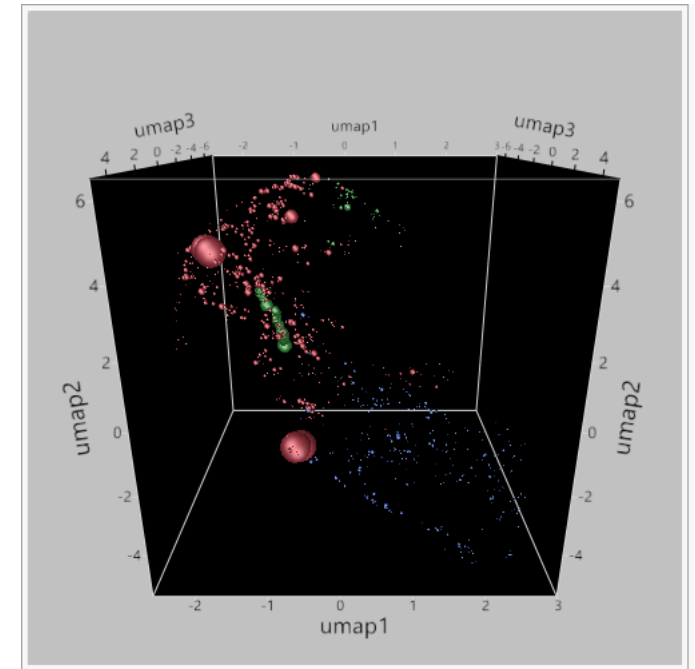
Principal Component Analysis



t-Distributed Stochastic Neighbor Embedding



Uniform Manifold Approximation and Projection

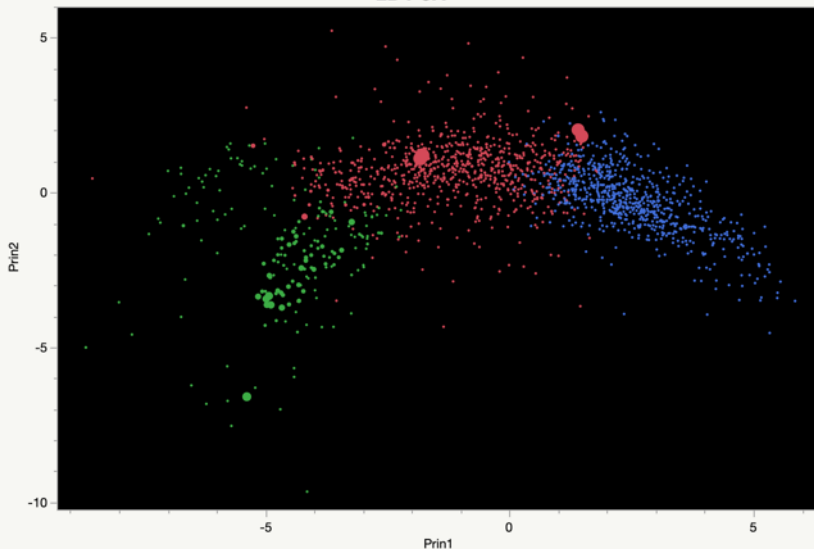


- Each is a different technique used to reduce number of dimensions but keep data variance and/or data structure
- Helps visualize data in 2 or 3D space

# Different Techniques for Data Visualizations

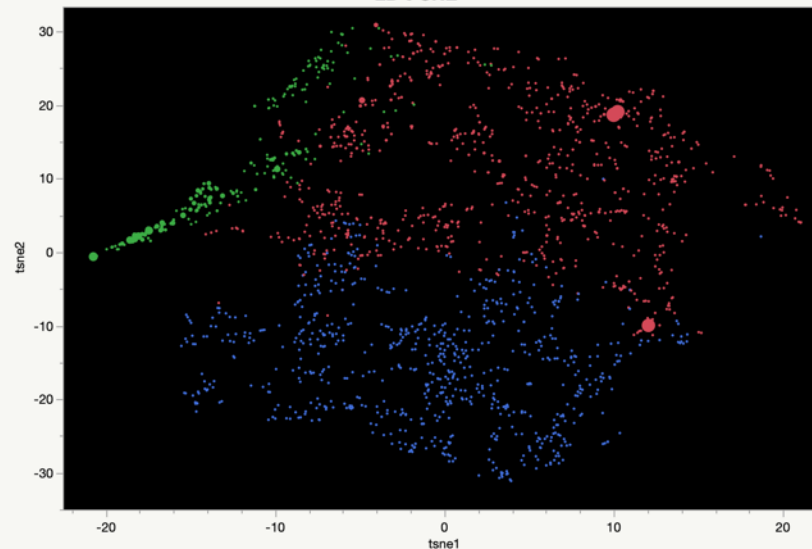
## PCA

2D PCA



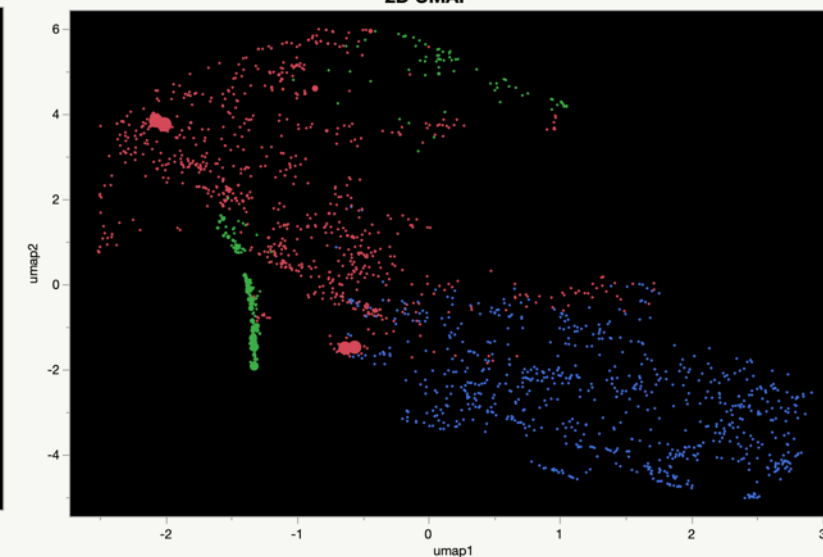
## t-SNE

2D t-SNE



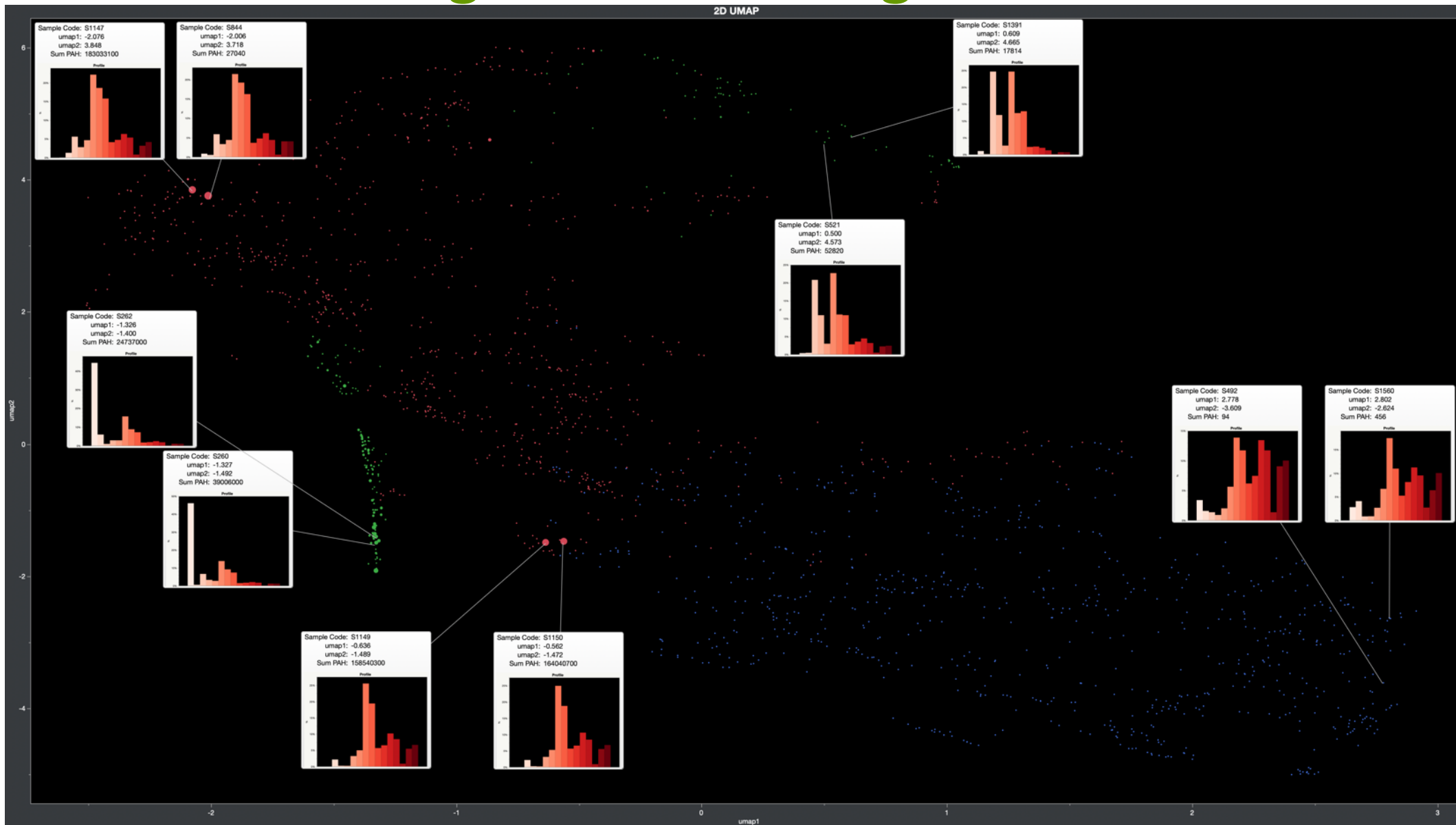
## UMAP

2D UMAP



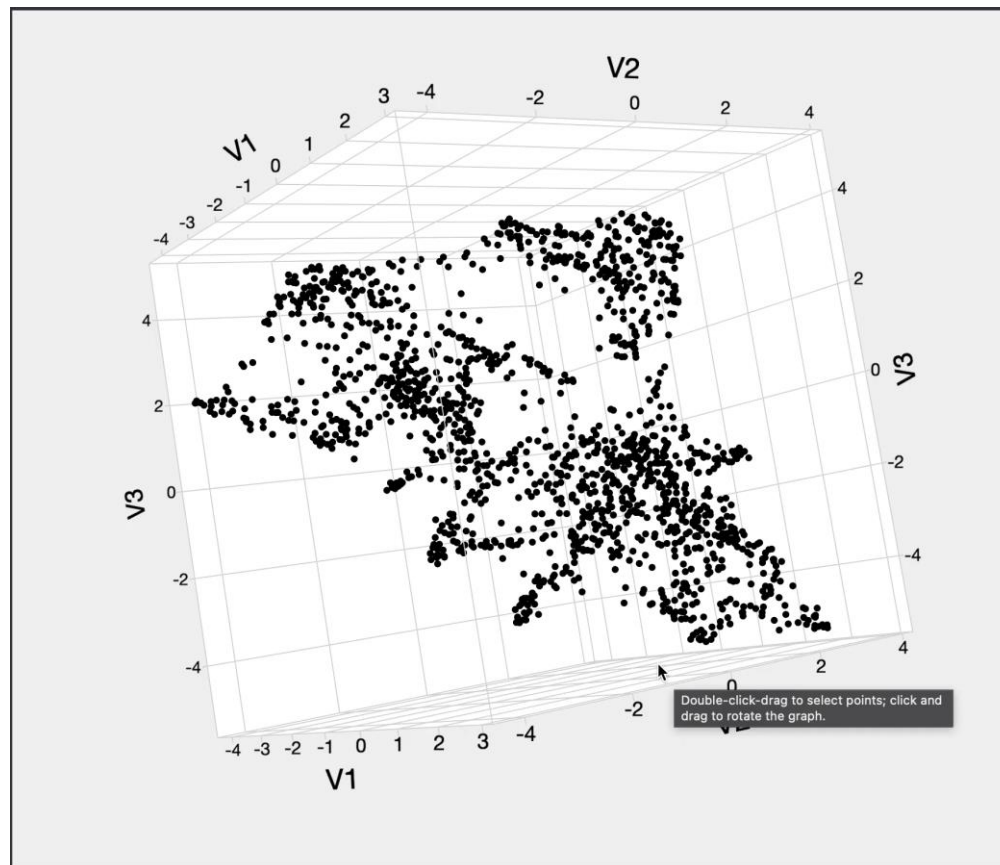
- Provides visualization of all compounds – patterns and concentrations
- Each has different math to derive visualization
- UMAP is newest kid on the block and my personal favourite

# Visualizing and Interacting with the Data



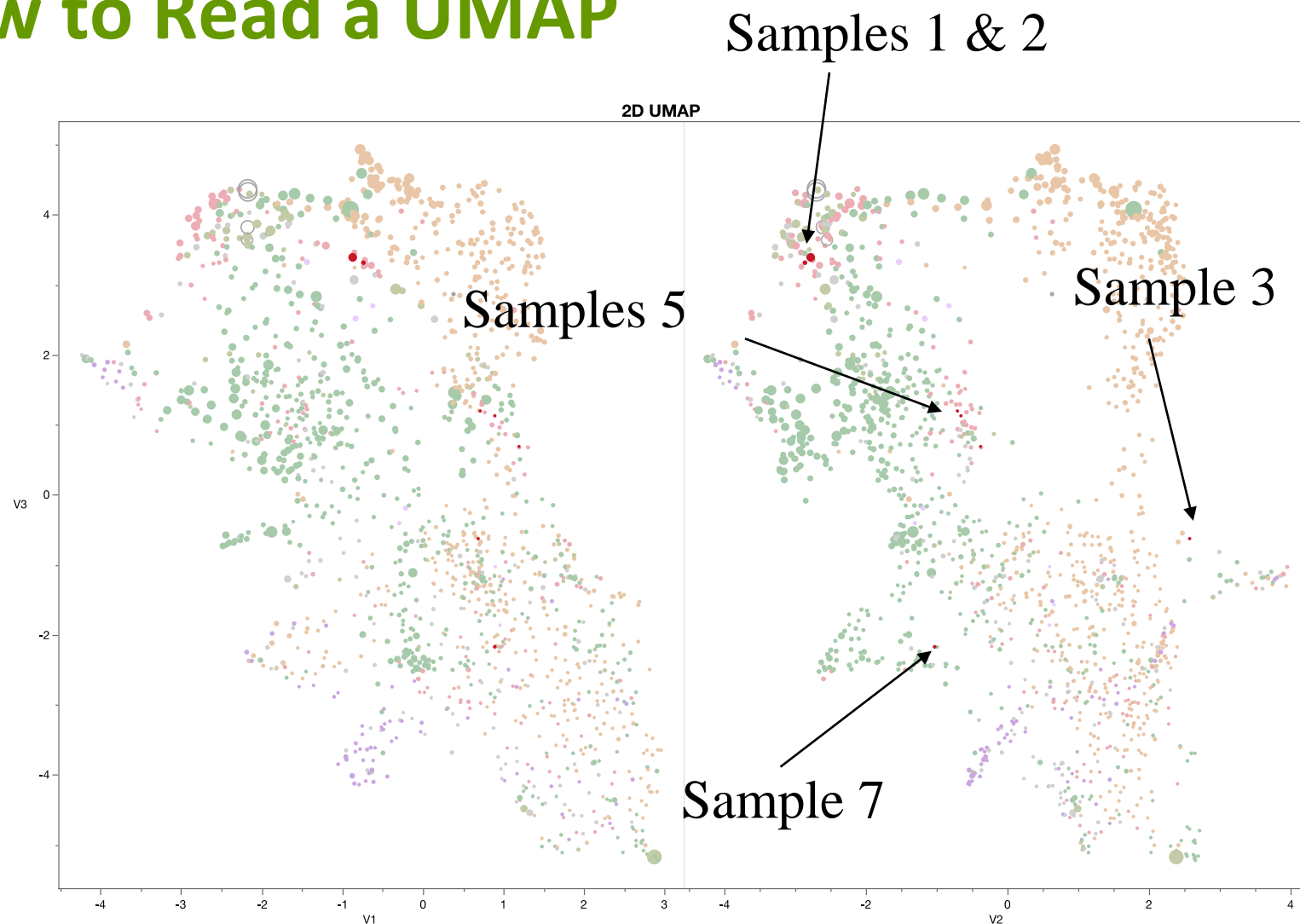


# How to Read a UMAP



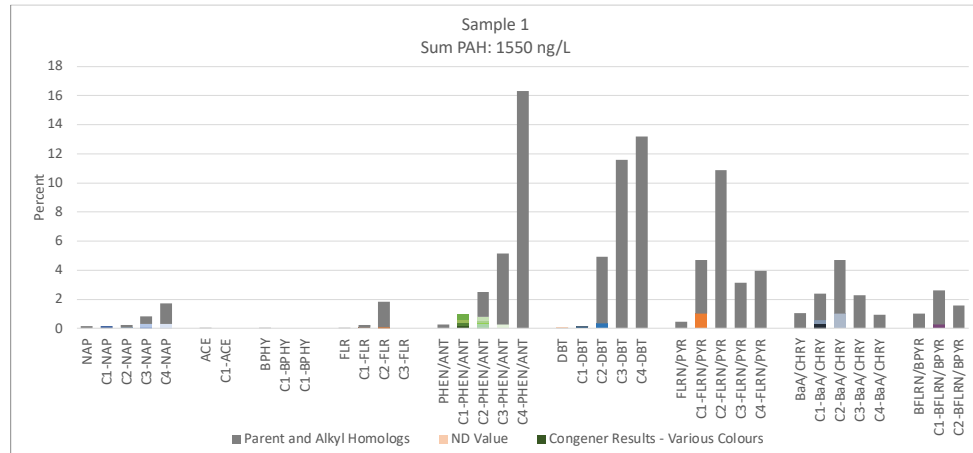
# How to Read a UMAP

- Samples close together are more similar
- Samples that are far away from each other are different
- Comparing all samples to Sample 1

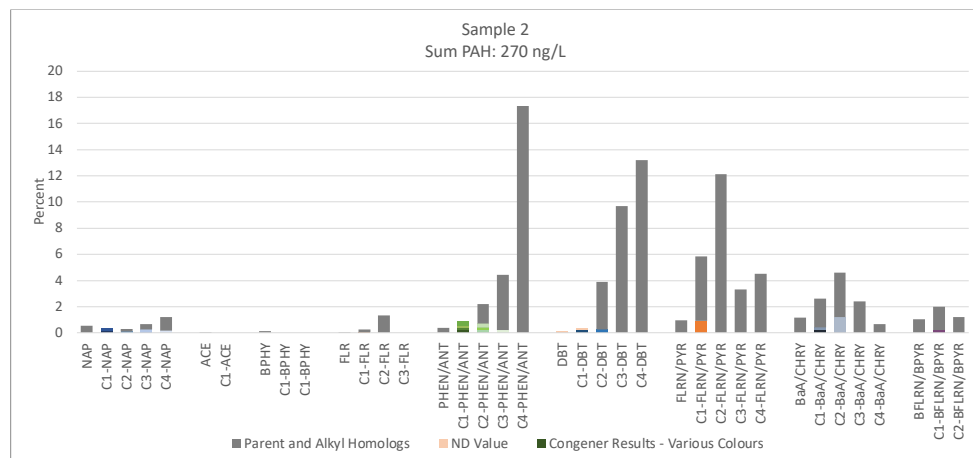


# How to Read a UMAP

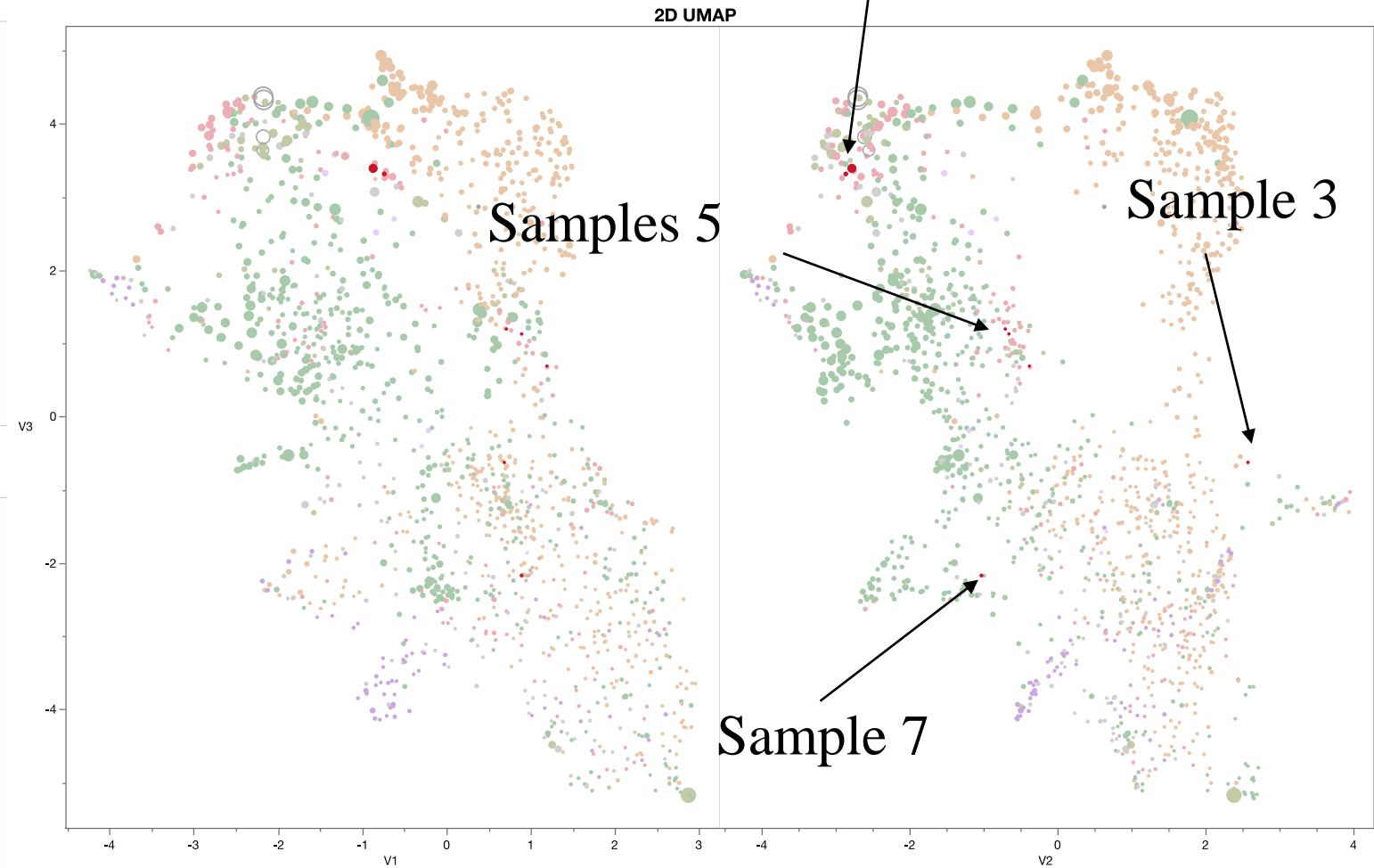
## Sample 1



## Sample 2

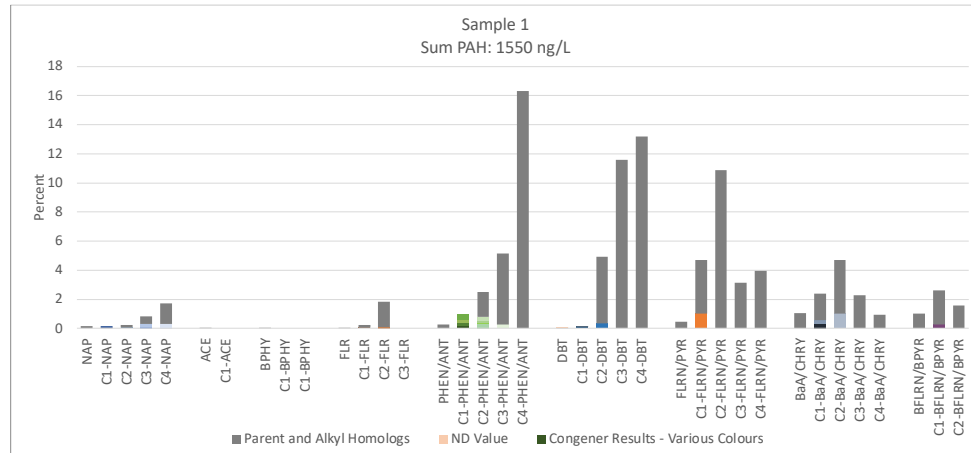


$\cos(\theta) : 0.995$

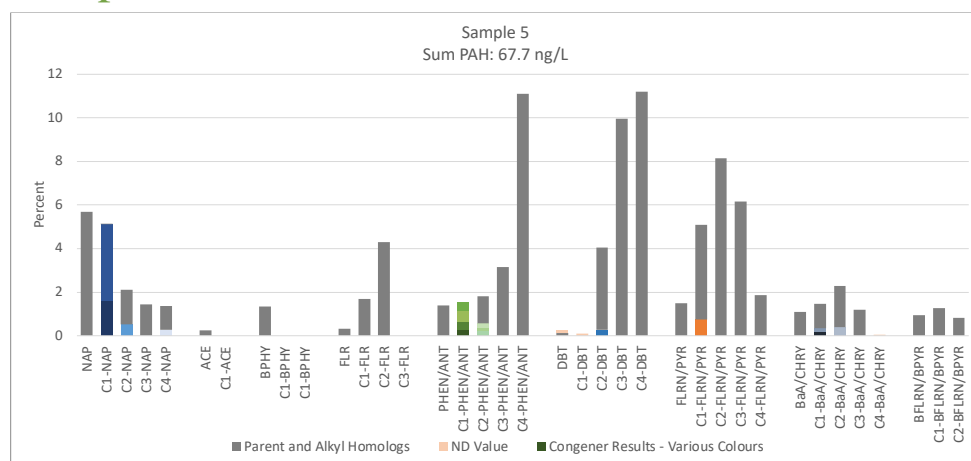


# How to Read a UMAP

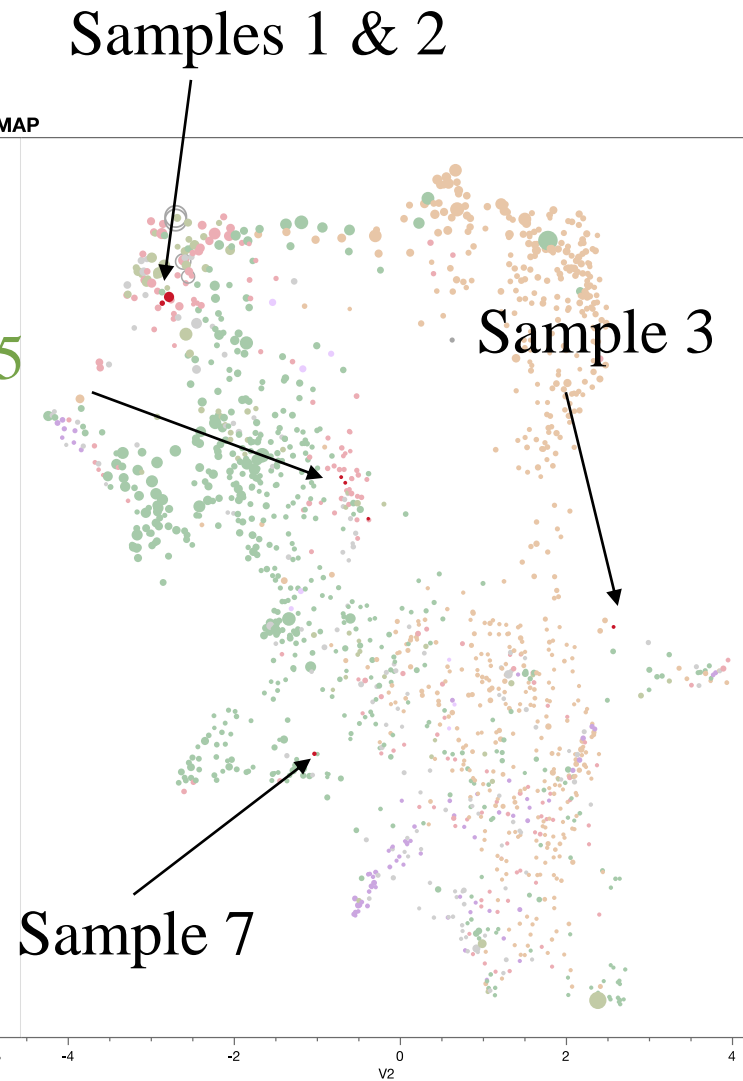
Sample 1



Sample 5

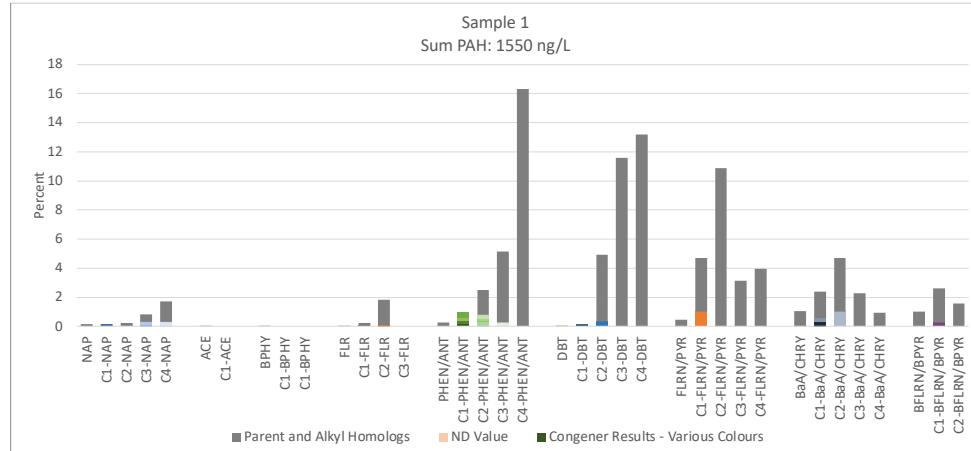


$\cos(\theta) : 0.904$

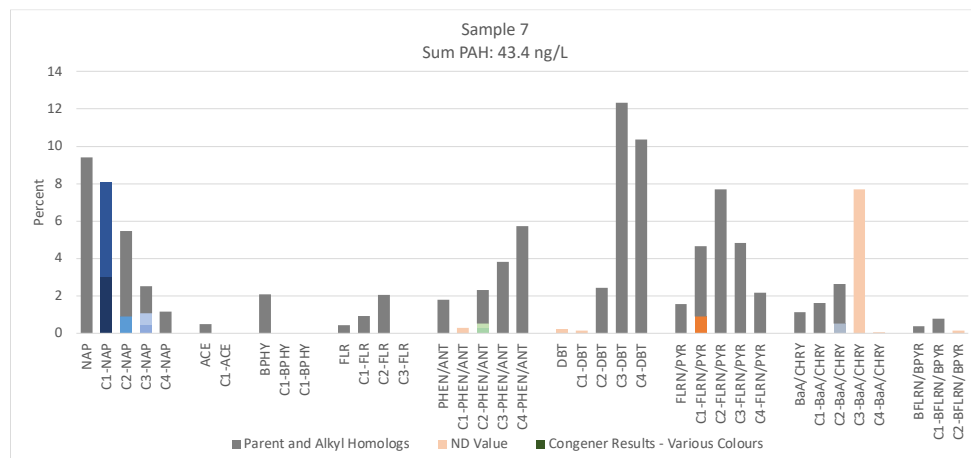


# How to Read a UMAP

## Sample 1



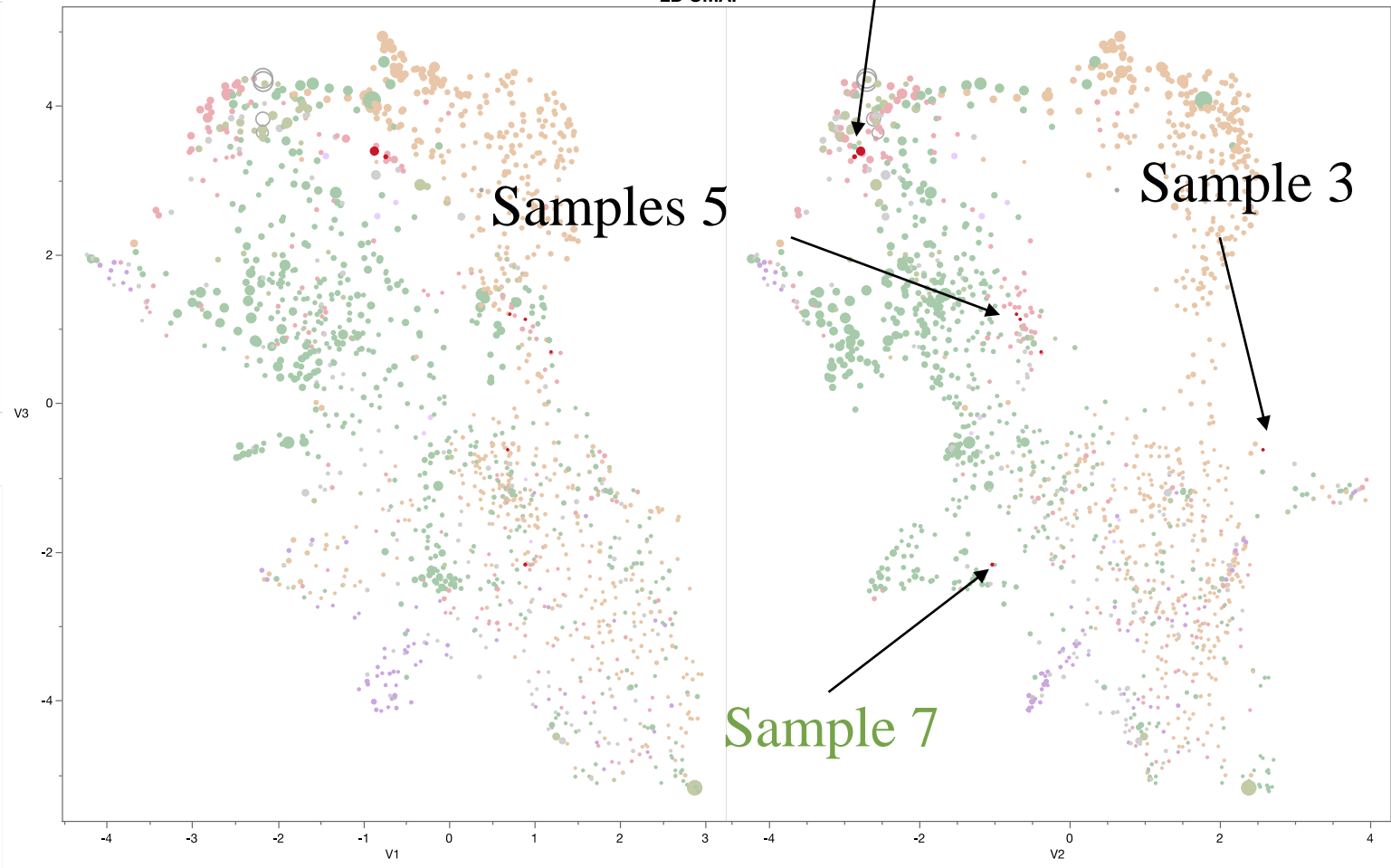
## Sample 7



$\cos(\theta) : 0.774$

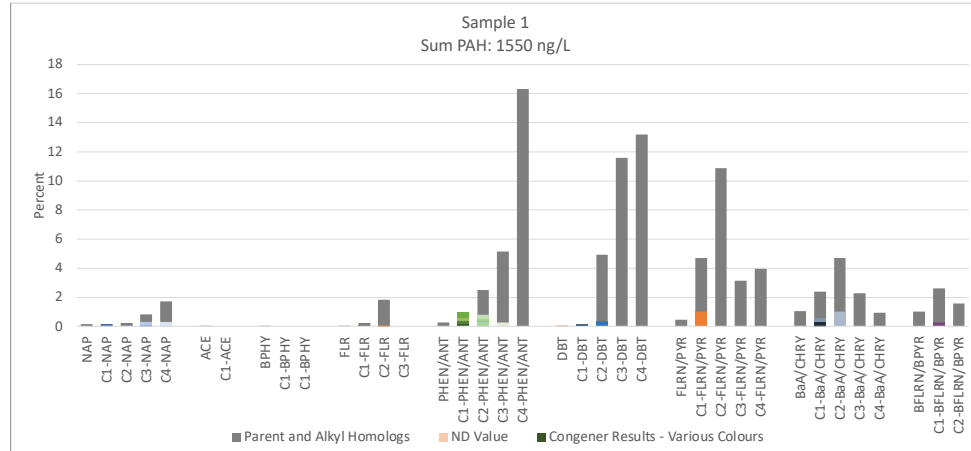
Samples 1 & 2

2D UMAP

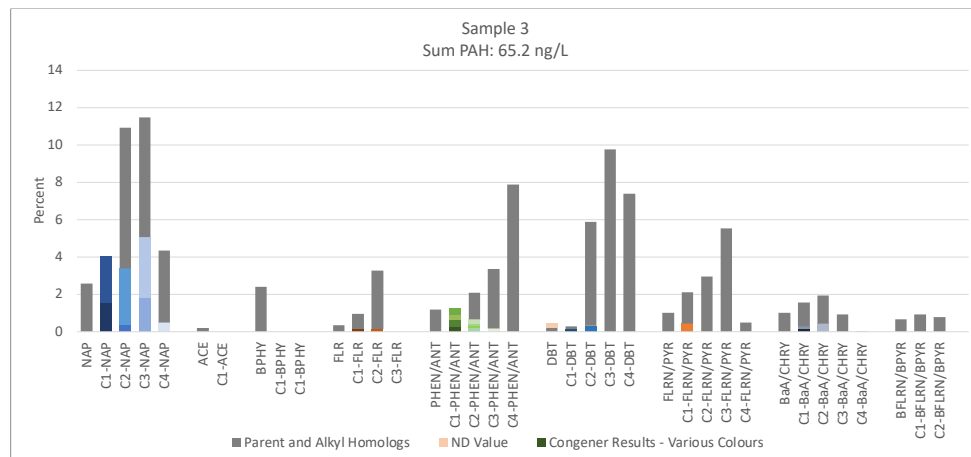


# How to Read a UMAP

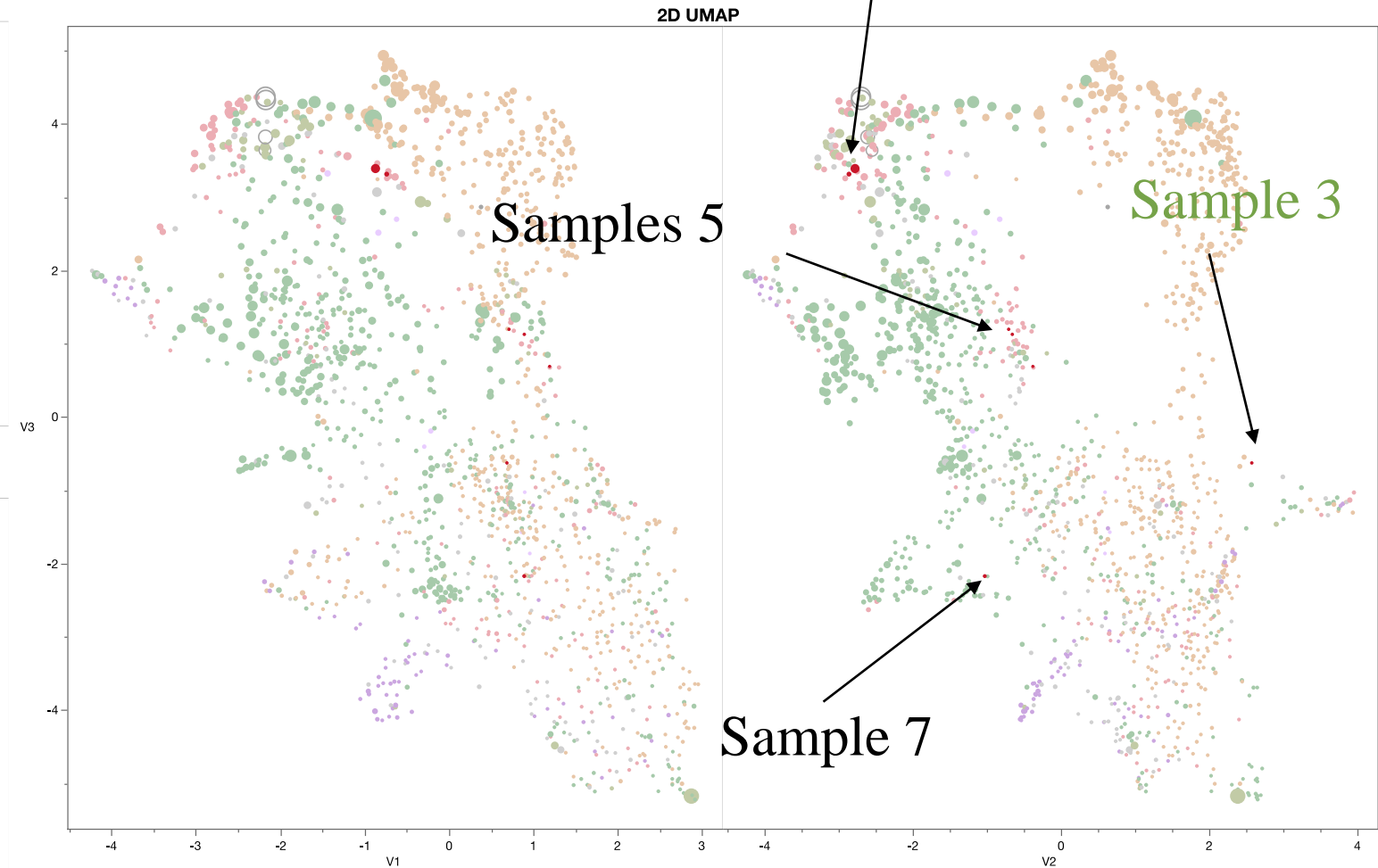
## Sample 1



## Sample 3

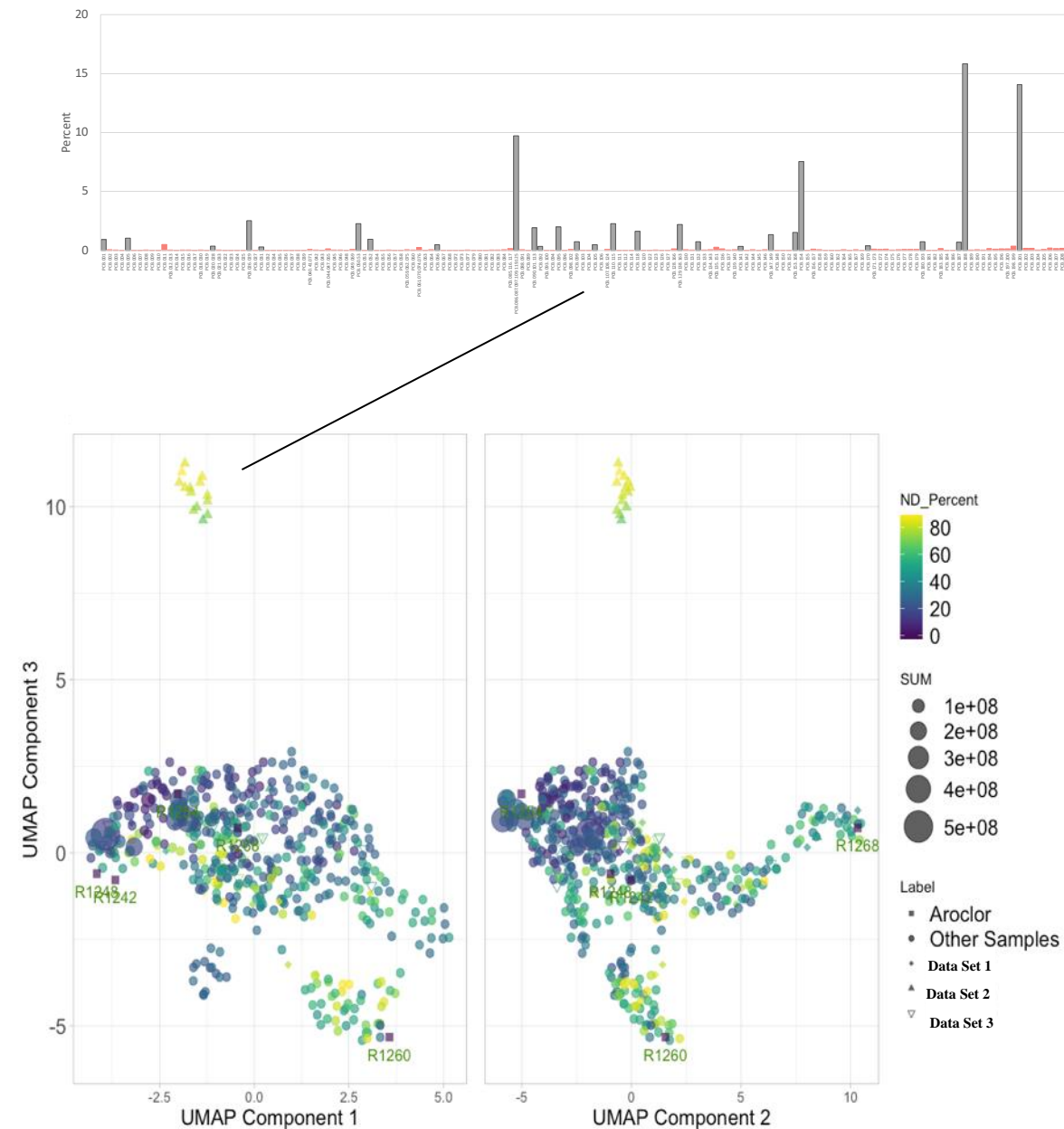


$\cos(\theta) : 0.675$



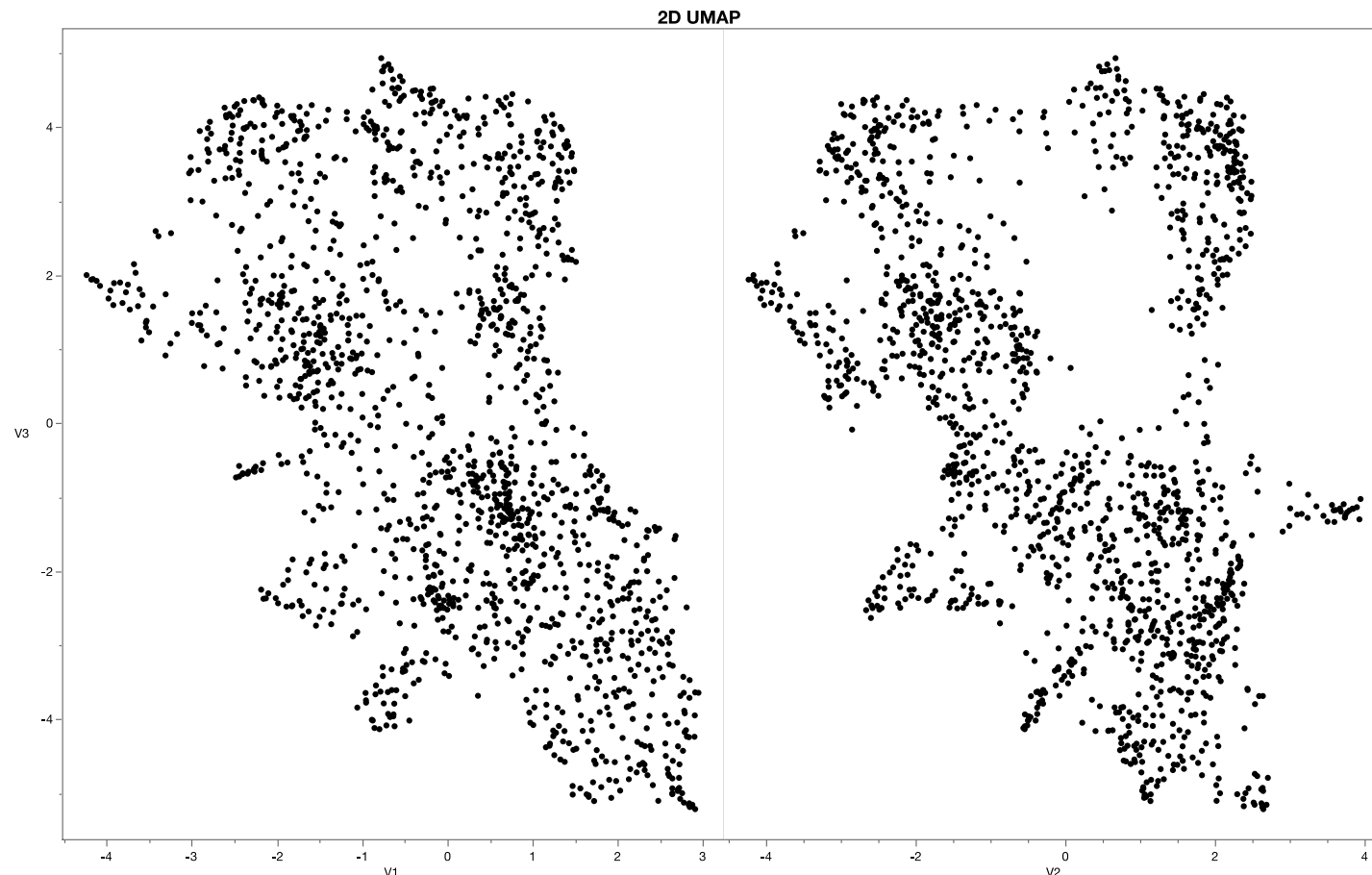
# UMAP – Finds Errors

- Preliminary look at data can help
- UMAP can find outliers and bad data
- We found a bad set of results for one site
- Odd fingerprint
- Data before and after this batch of samples looked normal (all same laboratory)
- Samples all collected from similar locations
- This dataset was erroneously reported by laboratory
- It was removed from further analysis



# Base UMAP

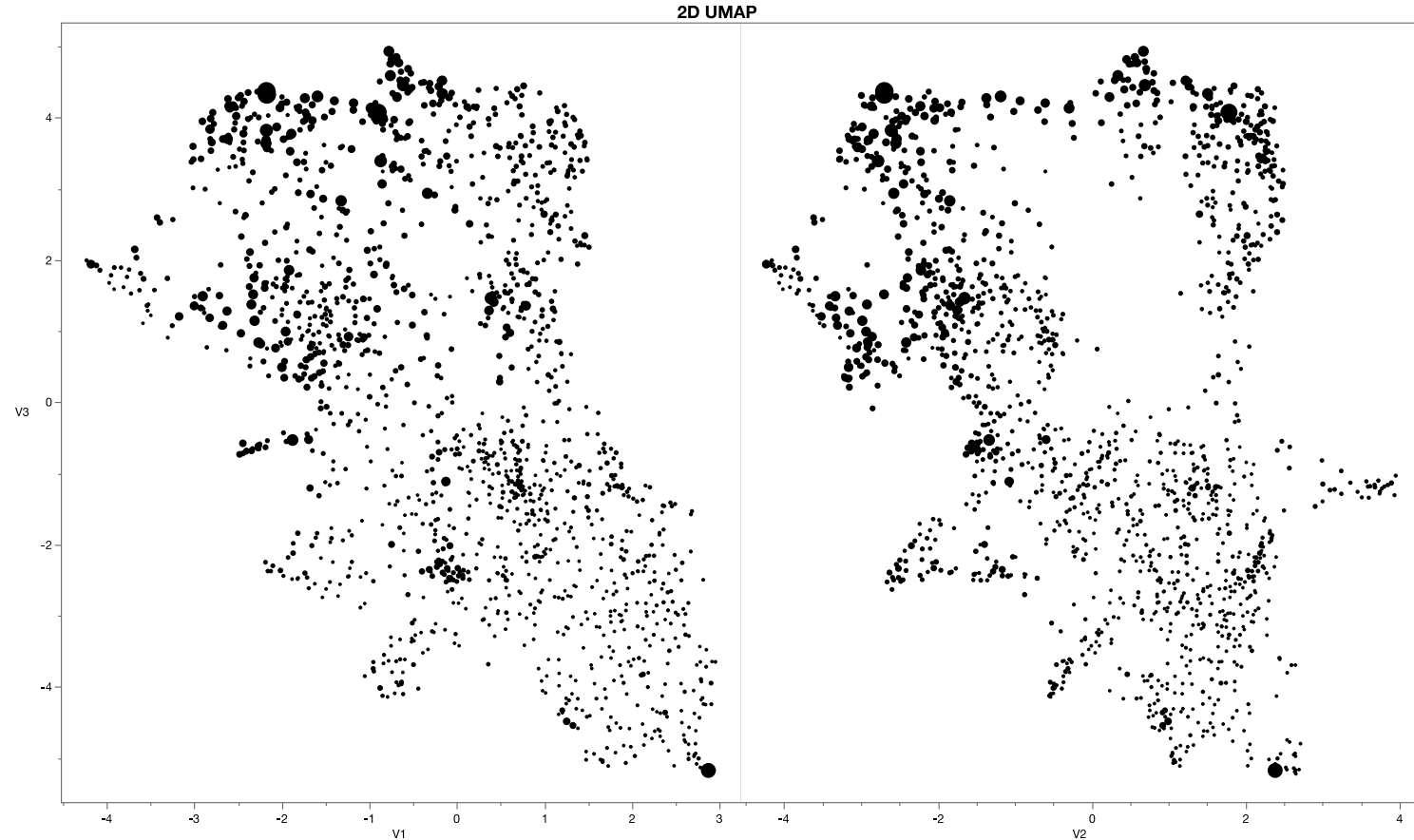
- Base UMAP shows unique data structure for all compounds (fingerprint) included in the database
- Data nodes (edges, extremes) provide clues as to the potential source-like samples in the dataset
- This dataset has multiple nodes





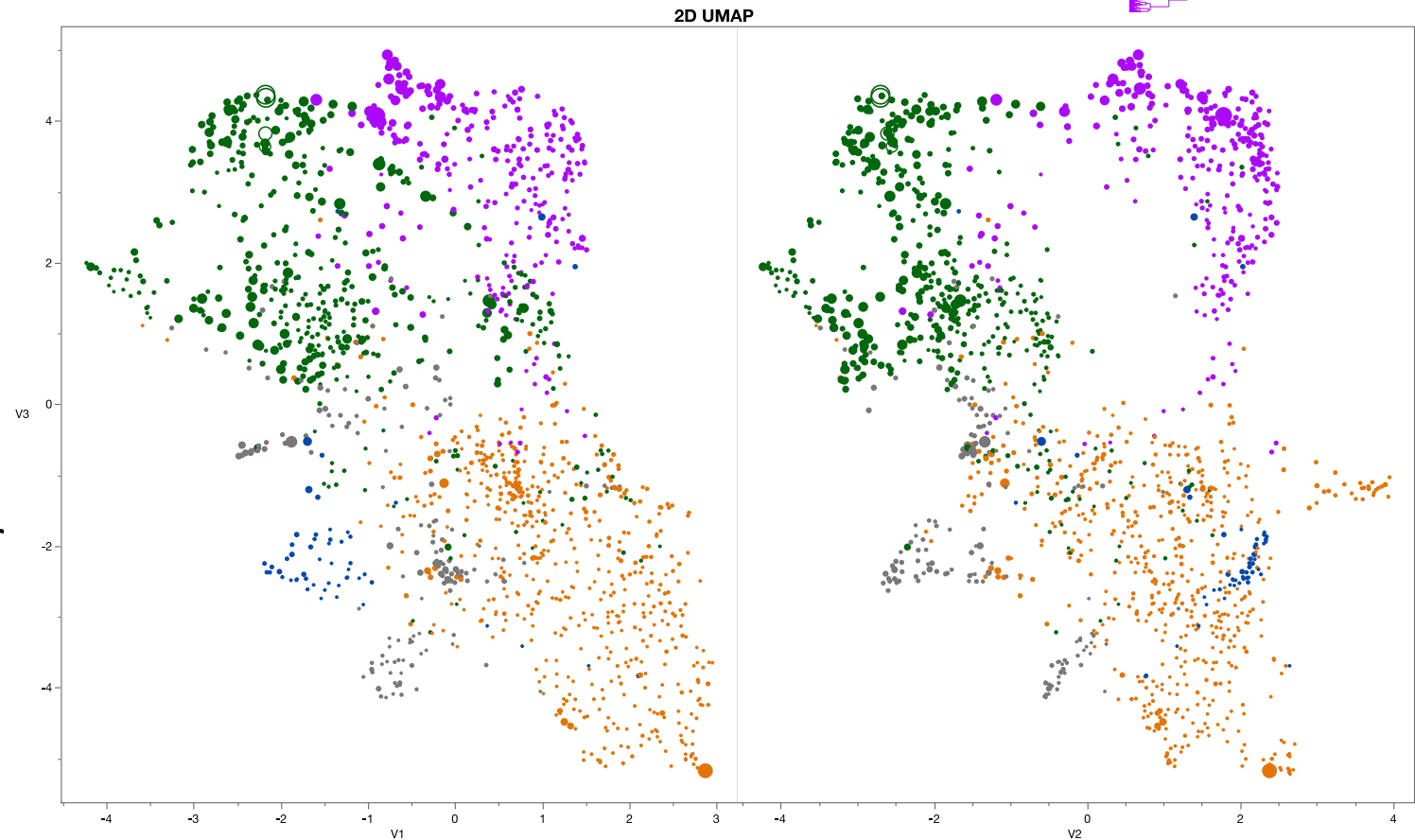
# Base UMAP – Size = Concentration

- Adding additional dimension to the dataset shows high concentration samples
- High concentrations generally associated with most source-like samples
  - Shows sources that are most important in this dataset
- Coupled with map, can show where highest concentration samples are located geospatially



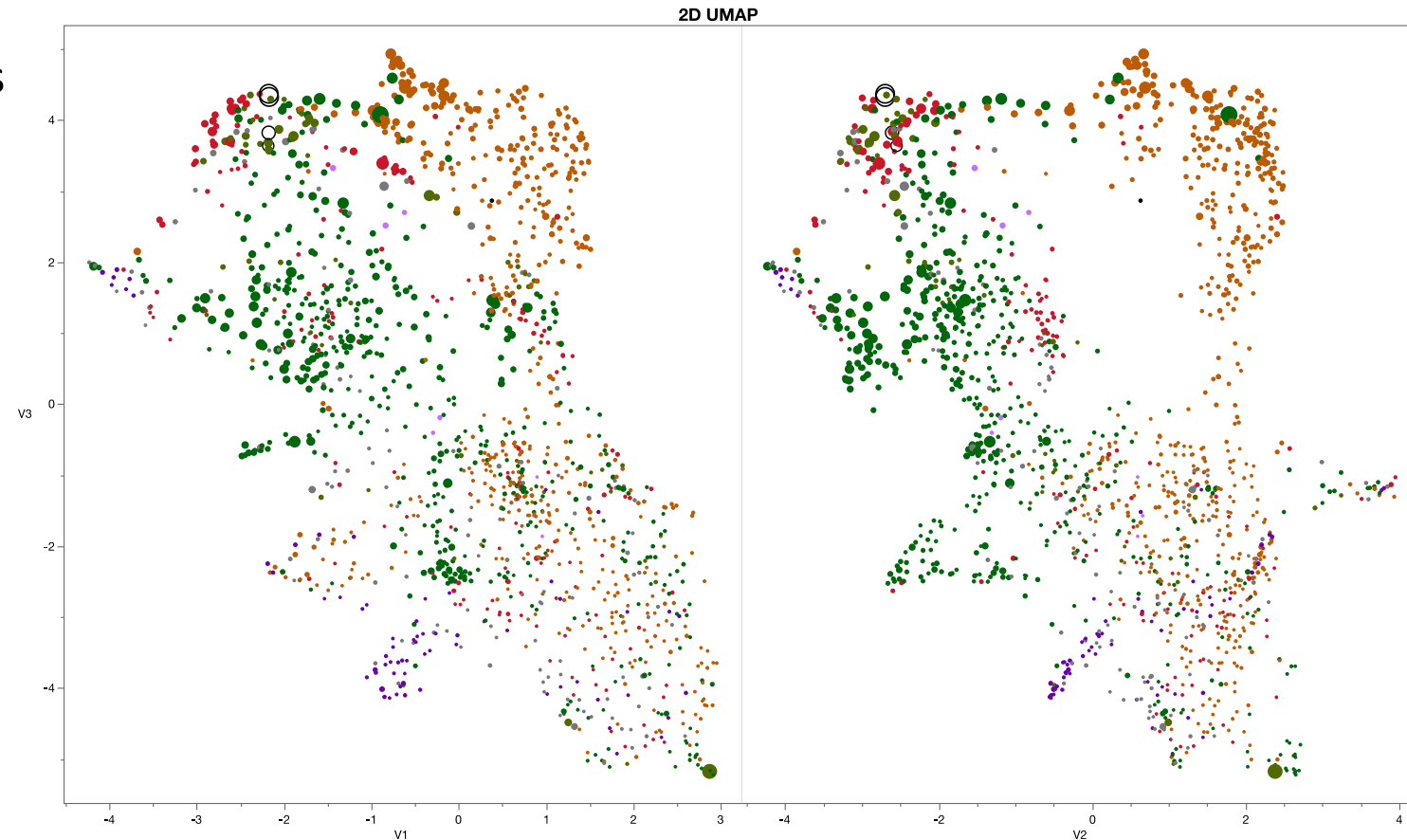
# UMAP – Add Sample Coding

- Coding data to match your dataset helps interpret what UMAP is telling you
  - preliminary hierarchical cluster analysis (HCA)
  - samples coded based on river system
- Other codings
  - sampling events
  - laboratories (to assess variability)
  - receptor modelling results
  - month/year of sampling
- Looking for groups, clusters, reasons for why the data is spread out the way it is



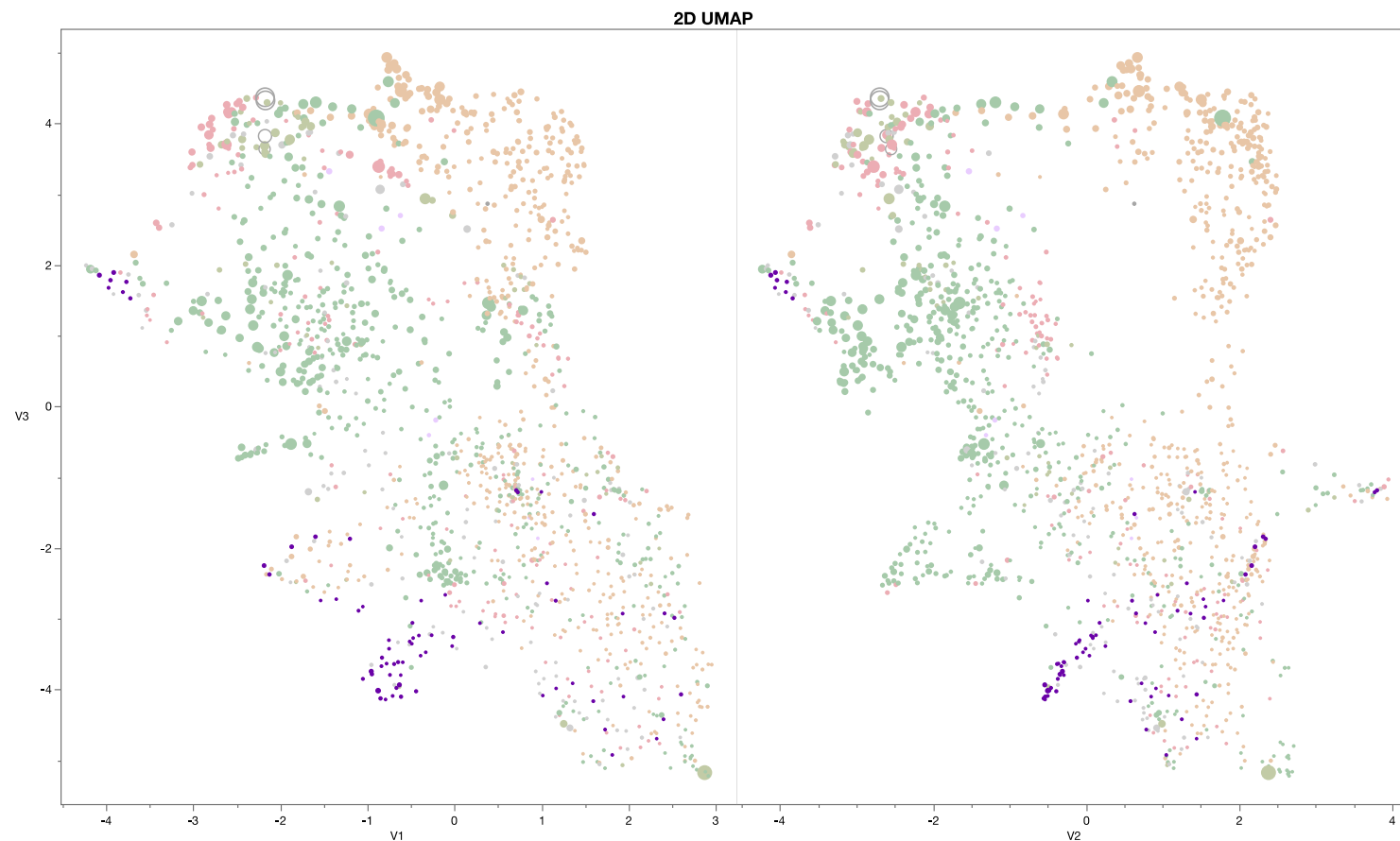
# UMAP – Add Sample Coding

- Coding data to match your dataset helps interpret what UMAP is telling you
  - preliminary hierarchical cluster analysis (HCA)
  - samples coded based on river system
- Other codings
  - sampling events
  - laboratories (to assess variability)
  - receptor modelling results
  - month/year of sampling
- Looking for groups, clusters, reasons for why the data is spread out the way it is



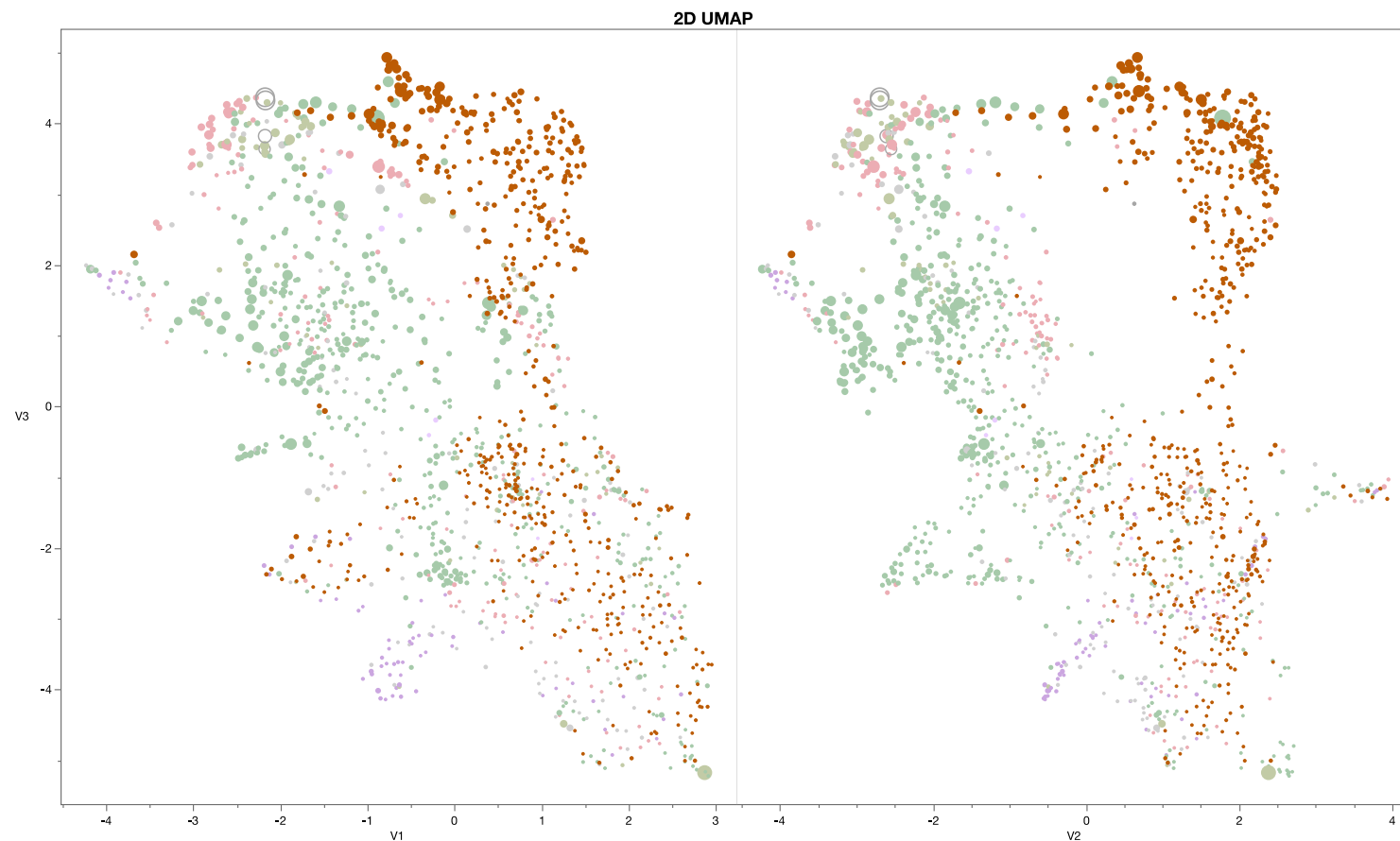
# UMAP – Scrolling through River Systems

- Ability to scroll through different river systems helps to identify what the chemical pattern of data might be telling you



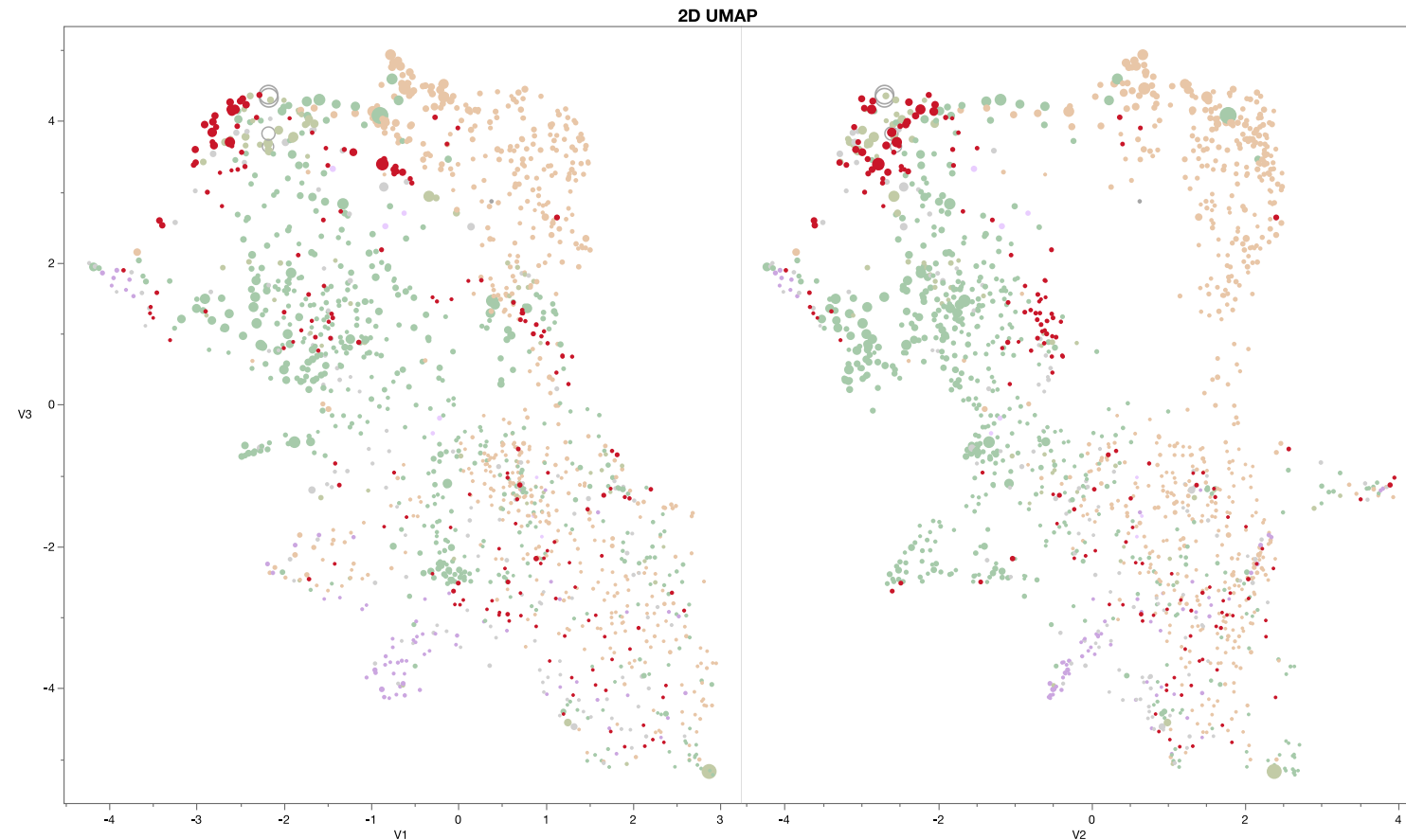
# UMAP – Scrolling through River Systems

- Ability to scroll through different river systems helps to identify what the chemical pattern of data might be telling you



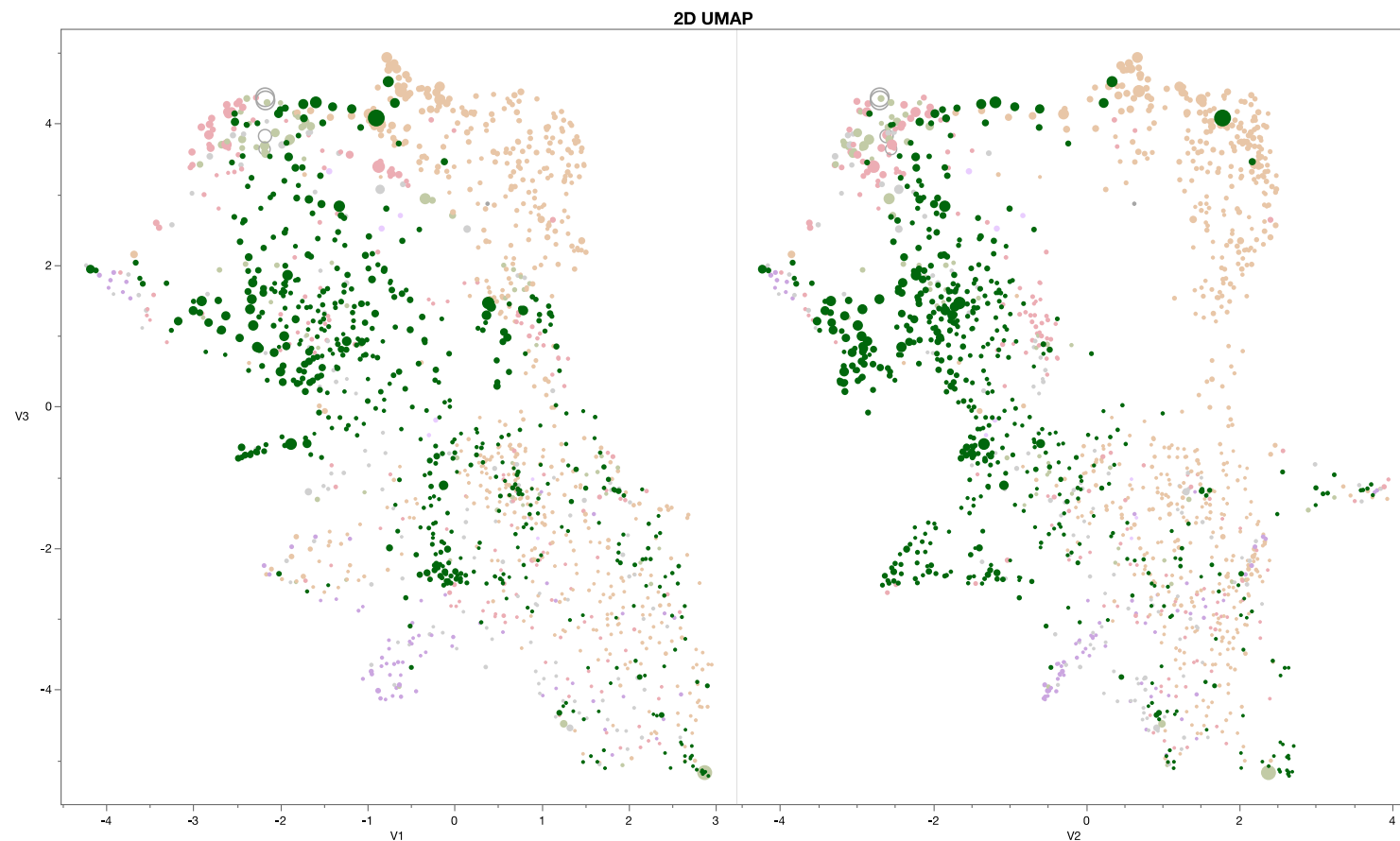
# UMAP – Scrolling through River Systems

- Ability to scroll through different river systems helps to identify what the chemical pattern of data might be telling you

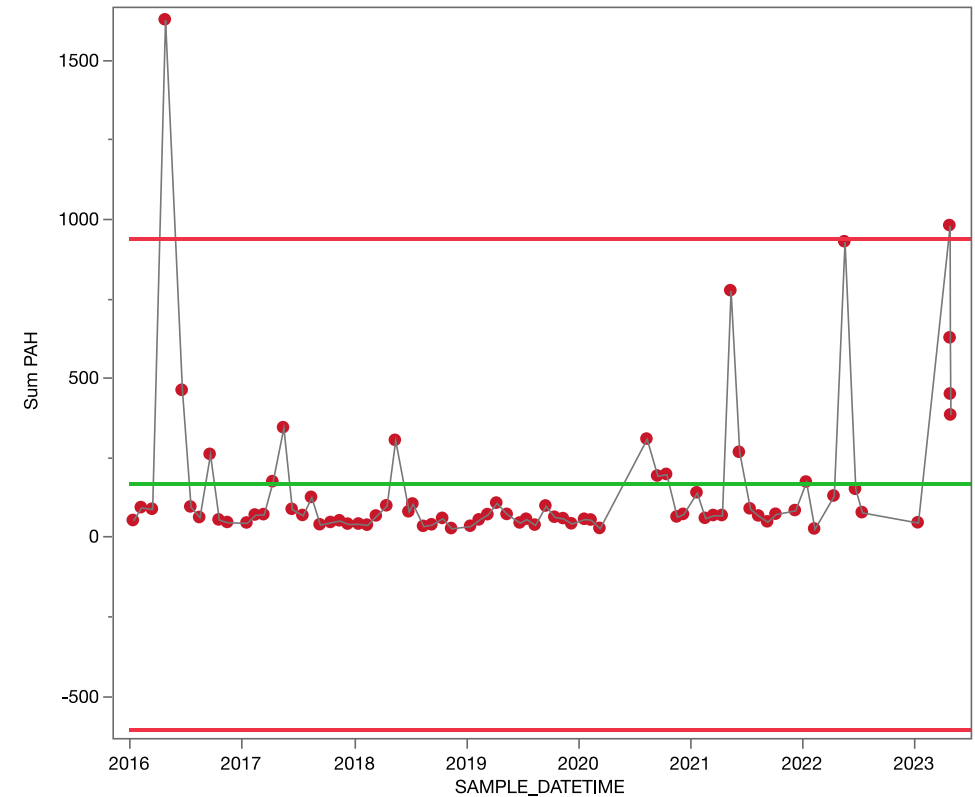
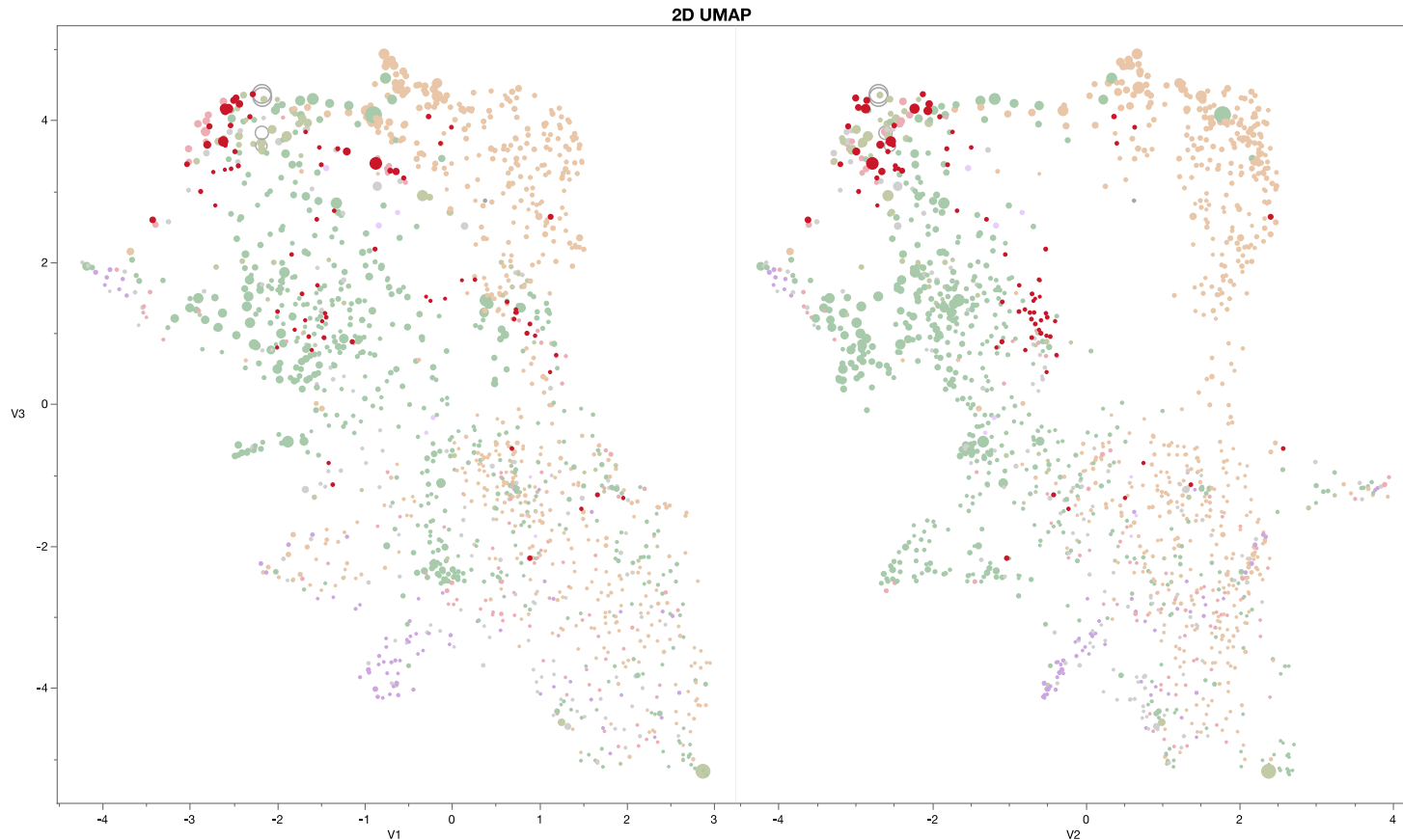


# UMAP – Scrolling through River Systems

- Ability to scroll through different river systems helps to identify what the chemical pattern of data might be telling you



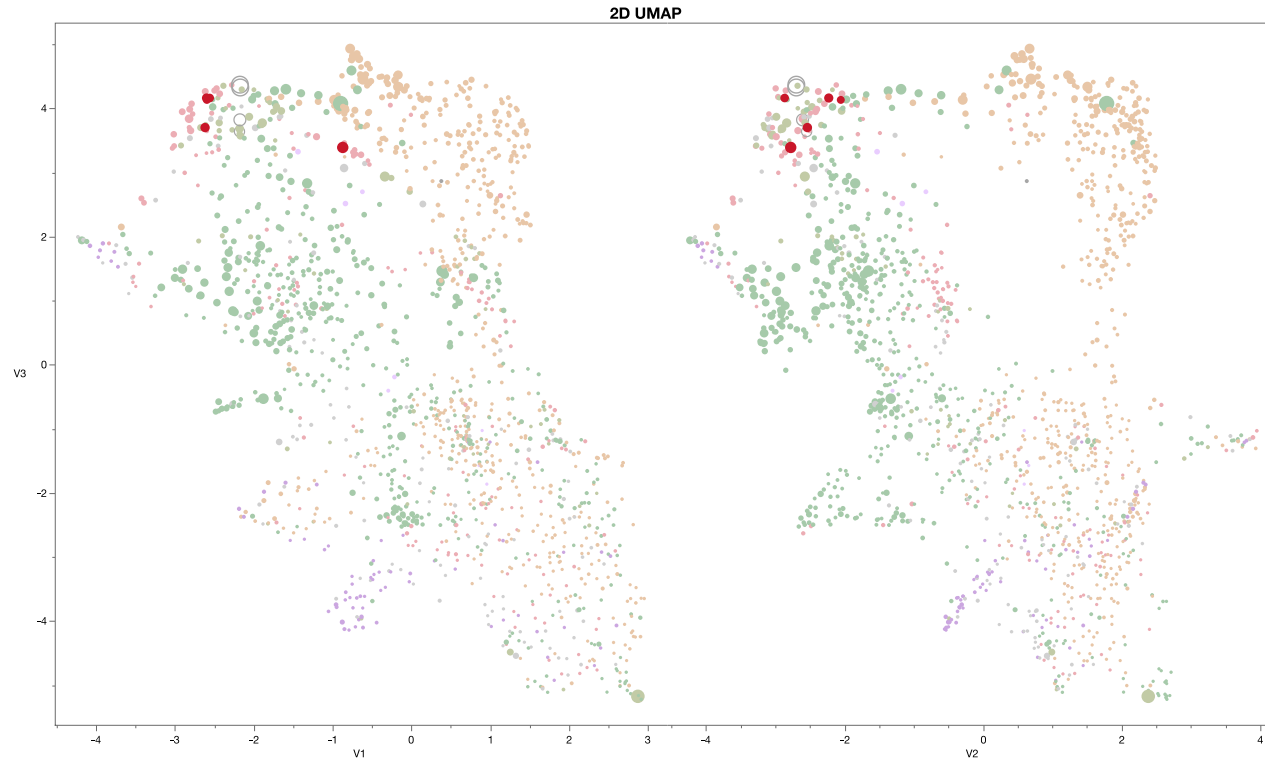
# Combining Monitoring with UMAP



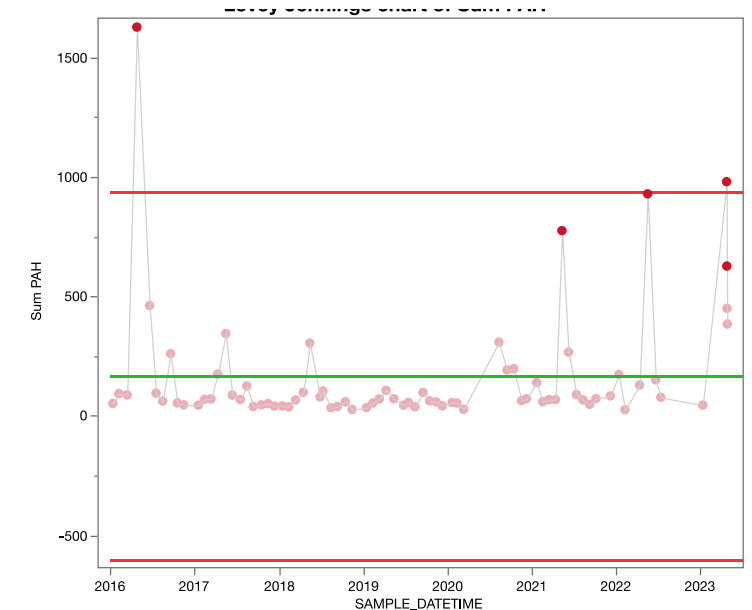
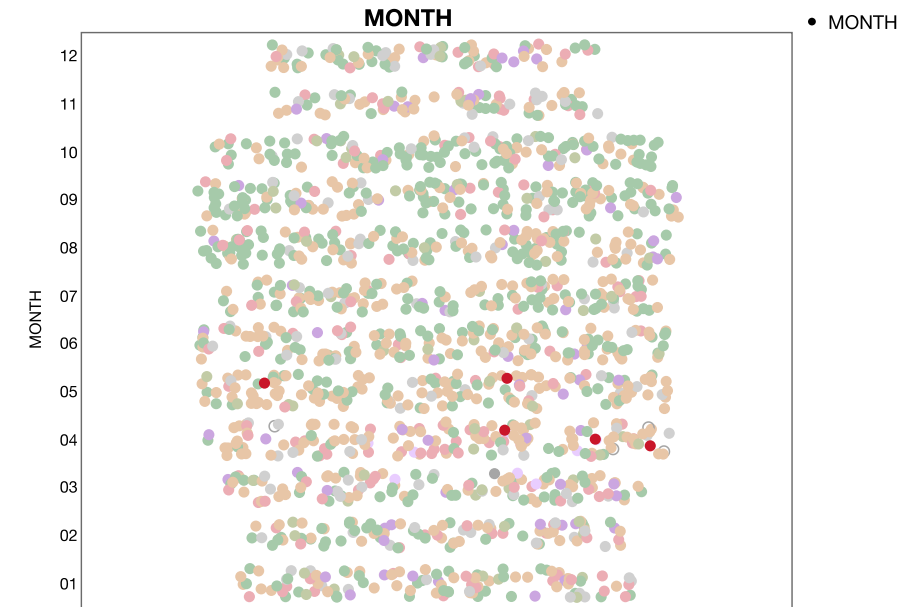
- A sampling location provides monitoring data over time. Examine trends with concentrations and patterns.
  - High concentrations seem to be annual events



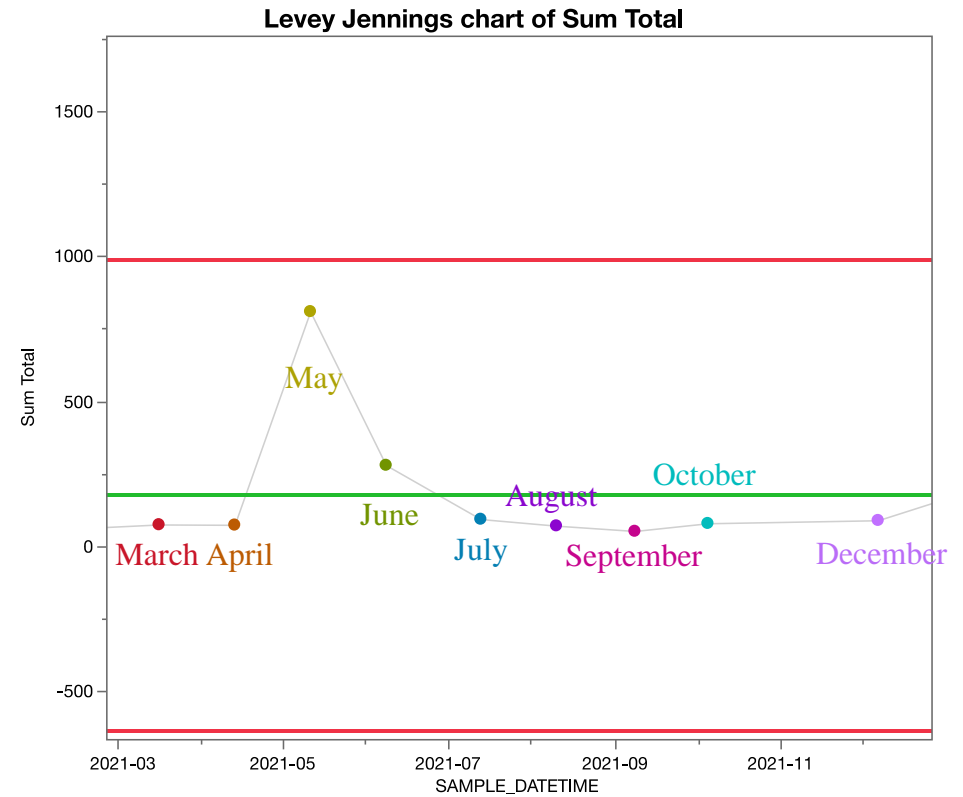
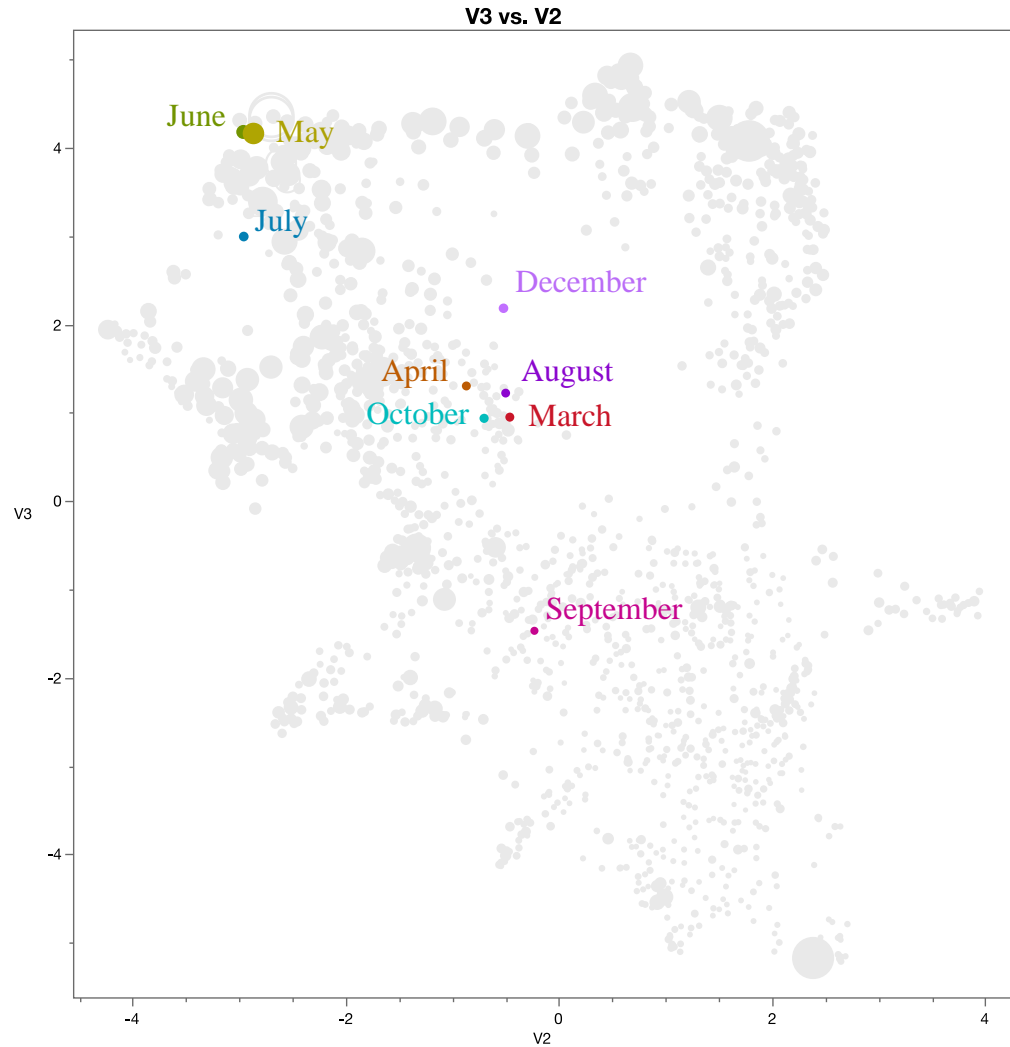
# Combining Monitoring with UMAP



- Spring melt (freshet) causes pulse of PAHs
  - Natural erosion?
  - Washing of source into river with melt?

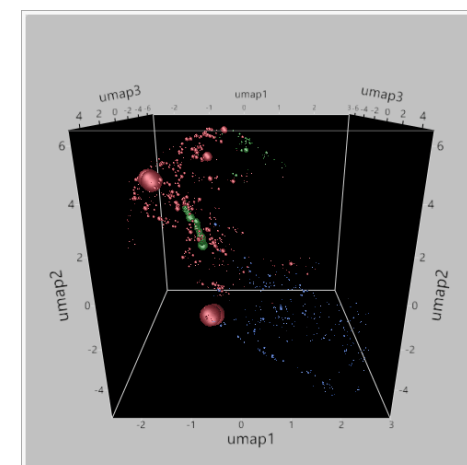


# Monthly Patterns Emerge



# Concluding Remarks

- Its okay to play with your data, just need to have a playroom that allows it
- Unlock learning about your data with interactive data exploration
- Seeing is believing
  - For you as you are exploring what your data is telling you
  - For client/public communication



# Questions?

Contact Info:

Court Sandau

Chemistry Matters Inc.

Emails: [csandau@chemistry-matters.com](mailto:csandau@chemistry-matters.com)

[court@statvis.com](mailto:court@statvis.com)

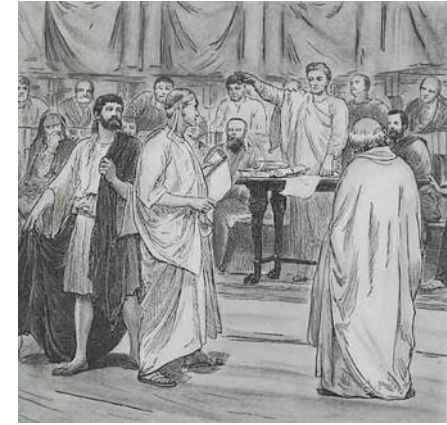
URLs: [chemistry-matters.com](http://chemistry-matters.com)

[statvis.com](http://statvis.com)

Follow us on: 

# Forensics

- Ancient roman times
- Forensics derived from latin “*forens*”:
  - belonging to the forum, public
  
- **Environmental forensics** –  
*determining identity, source, or timing of contaminant release for the purpose of public communication and/or litigation.*



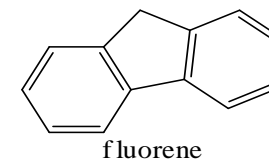
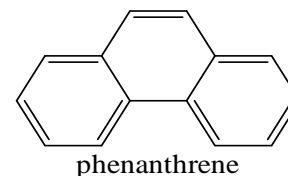
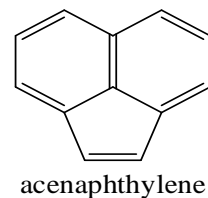
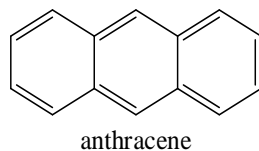
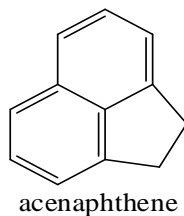
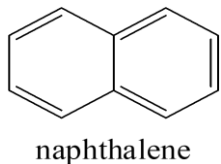
<http://www.crystalinks.com/romelaw.html>



# US EPA Priority Pollutants (PAH)

- In 1976, the US Environmental Protection Agency (USEPA) selected 16 PAHs as priority organic pollutants

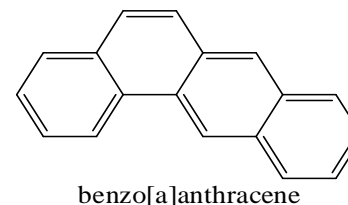
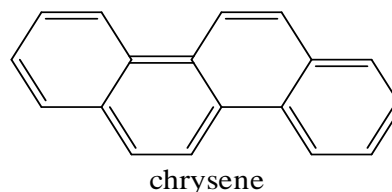
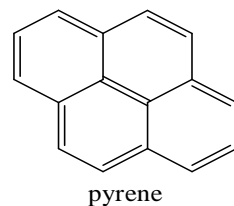
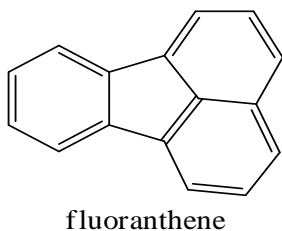
**2-ring**



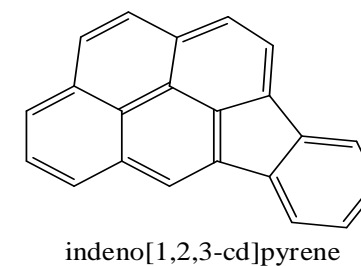
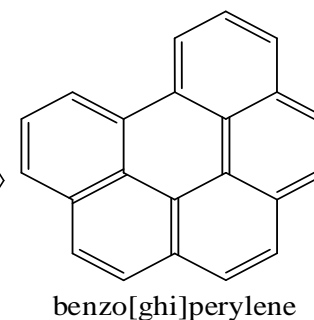
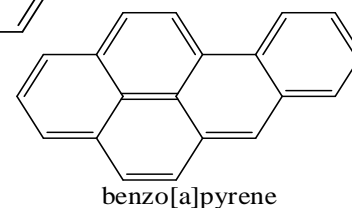
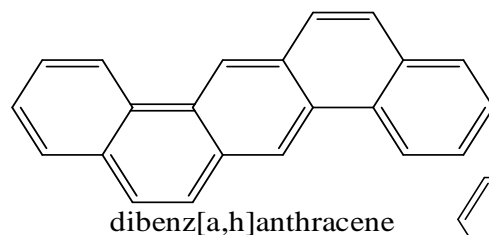
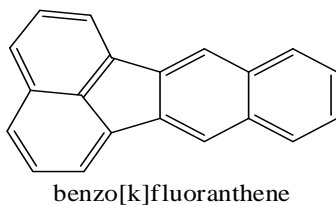
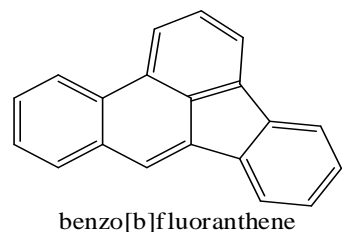
**LMW**

**3-ring**

**HMW**



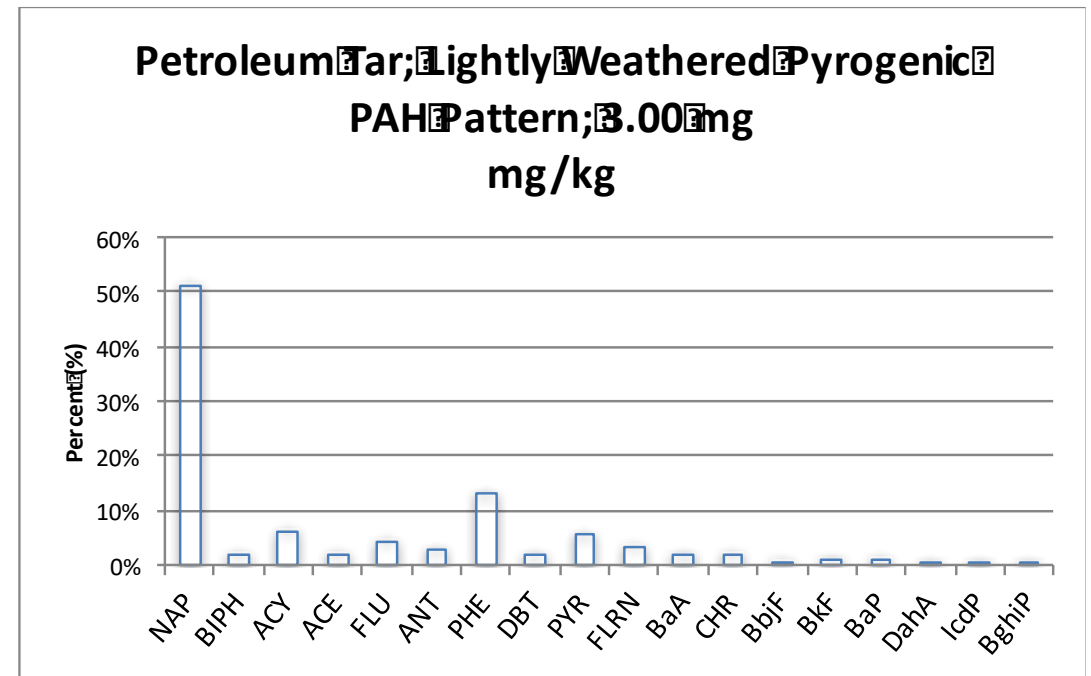
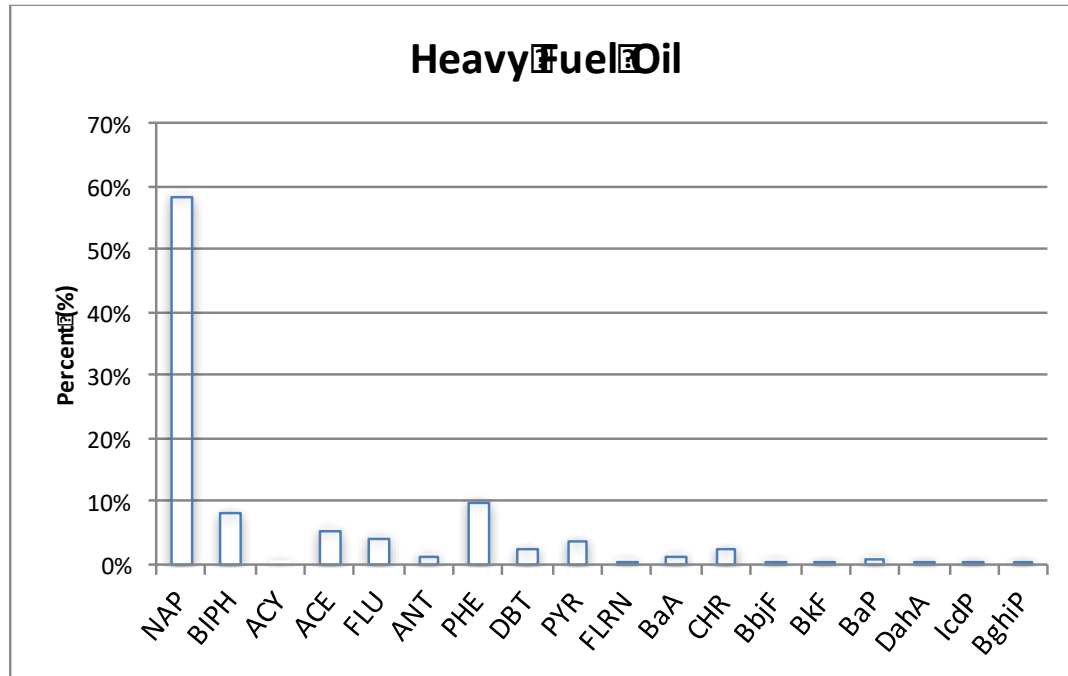
**4-ring**



**5-ring**

**6-ring**

# Forensic Fingerprints



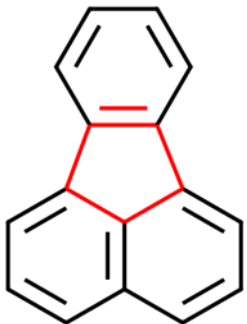
Samples from the same source will have same patterns.

With only 16 compounds, cannot distinguish fuel oil from gas plant tar. Not very forensic! But, need to use the data we have.

# Reliable Diagnostic Ratios

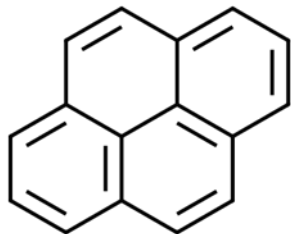
High Temperature Formation

Low Temperature Formation



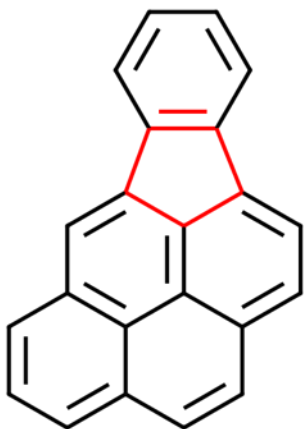
Fluoranthene

291 kJ/mol



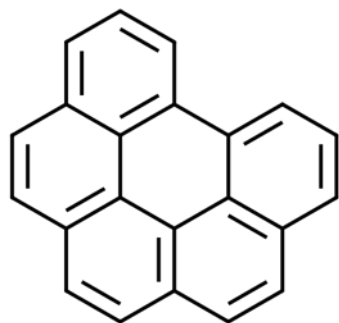
Pyrene

225 kJ/mol



Indeno[123,cd]pyrene

430 kJ/mol



Benzo[ghi]perylene

300 kJ/mol

- Uses fundamental chemistry properties (heats of formation) which drives their formation ( $\Delta_f H^\circ_m$ )
- Uses 4 of the 16 priority PAHs
- Uses similar size structures for ratio (behave the same in environment)



# Double Ratio Plots to Distinguish Source

