

# The NTA WebApp:

A web-based tool for rapid chemical identification from non-targeted analysis mass spectrometry data

*Alex Chao*

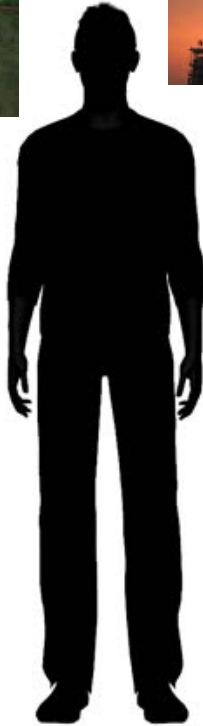
*Jeffrey M. Minucci, Matthew W. Boyce, S. Thomas Purucker, Antony J. Williams, Jon R. Sobus*



*The views expressed in this presentation are those of the author and do not necessarily reflect the views or policies of the U.S. EPA*

# Exposomics: Addressing Health Issues

## The “Exposome”



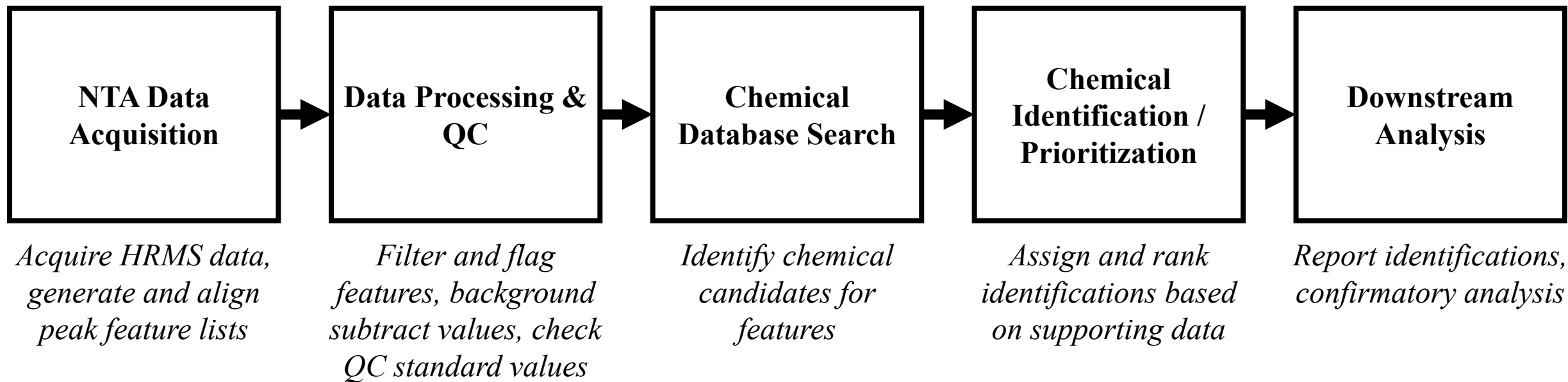
Chemical  
Monitoring  
Methods

What are the  
chemicals we  
are exposed to?

What are the  
amounts?

Do we need to  
be worried about  
it?

# Development of an NTA Workflow



# Development of a NTA Workflow



## ■ The Underlying Database: *DSSTox*

- A highly curated database of environmentally relevant chemicals (>900k)

## ■ Generating Compatibility with HRMS Data: *MS-Ready Forms*

- Mapping substance components into forms that would be observed by MS

## ■ Allowing Access into the Database: *CompTox Chemicals Dashboard*

- A web-based conduit into DSSTox allowing for batch searching of MS data

EPA's DSSTox database: History of development of a curated chemistry resource supporting computational toxicology research

Christopher M Grulke<sup>1</sup>, Antony J Williams<sup>1</sup>, Inthirany Thillanadarajah<sup>2</sup>, Ann M Richard<sup>1</sup>

"MS-Ready" structures for non-targeted high-resolution mass spectrometry screening studies

Andrew D McEachran<sup>1 2</sup>, Kamel Mansouri<sup>3 4 5</sup>, Chris Grulke<sup>4</sup>, Emma L Schymanski<sup>6</sup>, Christoph Ruttkies<sup>7</sup>, Antony J Williams<sup>8</sup>

The CompTox Chemistry Dashboard: a community data resource for environmental chemistry

Antony J Williams<sup>1</sup>, Christopher M Grulke<sup>2</sup>, Jeff Edwards<sup>2</sup>, Andrew D McEachran<sup>3</sup>, Kamel Mansouri<sup>2 3 4</sup>, Nancy C Baker<sup>5</sup>, Grace Patlewicz<sup>2</sup>, Imran Shah<sup>2</sup>, John F Wambaugh<sup>2</sup>, Richard S Judson<sup>2</sup>, Ann M Richard<sup>2</sup>

# Development of a NTA Workflow



- **Chemical ID via metadata: *Data Sources***
  - Count of chemical presence in publications, projects, data collections
- **Chemical ID via *in silico* predictions: *CFM-ID generated MS2 spectra***
  - Predicted MS2 spectra generated for all chemicals within DSSTox
- **Performance evaluations:**
  - Critical Assessment of Small Molecule Identification (**CASMI**) spectra
  - EPA's Non-Targeted Analysis Collaborative Trial (**ENTACT**) mixtures

Identifying known unknowns using the US EPA's CompTox Chemistry Dashboard

Andrew D McEachran <sup>1</sup>, Jon R Sobus <sup>2</sup>, Antony J Williams <sup>3</sup>

Linking in silico MS/MS spectra with chemistry data to improve identification of unknowns

Andrew D McEachran <sup>1 2</sup>, Ilya Balabin <sup>3</sup>, Tommy Cathey <sup>4</sup>, Thomas R Transue <sup>4</sup>, Hussein Al-Ghoul <sup>5</sup>, Chris Grulke <sup>6</sup>, Jon R Sobus <sup>7</sup>, Antony J Williams <sup>8</sup>

In silico MS/MS spectra for identifying unknowns: a critical examination using CFM-ID algorithms and ENTACT mixture samples

Alex C Tomm Anton  
Revisiting Five Years of CASMI Contests with EPA Identification Tools

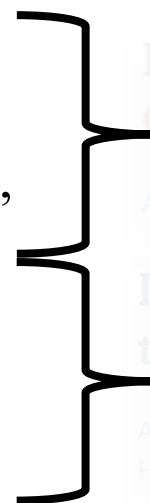
Andrew D McEachran <sup>1</sup>, Alex Chao <sup>1</sup>, Hussein Al-Ghoul <sup>1</sup>, Charles Lowe <sup>2</sup>, Christopher Grulke <sup>2</sup>, Jon R Sobus <sup>2</sup>, Antony J Williams <sup>2</sup>



# Development of a NTA Workflow



- **Chemical ID via metadata: *Data Sources***
  - Count of chemical presence in publications, projects, data collections
- **Chemical ID via *in silico* predictions: *CFM-ID generated MS2 spectra***
  - Predicted MS2 spectra generated for all chemicals within DSSTox



***DSSTox***

***Database API***

- **Performance evaluations:**
  - Critical Assessment of Small Molecule Identification (CASMI) spectra
  - EPA's Non-Targeted Analysis Collaborative Trial (ENTACT) mixtures

Identifying known unknowns using the US EPA's  
Chemistry Dashboard

Andrew D McEachran <sup>1</sup>, Jon R Sobus <sup>2</sup>, Antony J Williams <sup>3</sup>

Linking *in silico* MS/MS spectra with chemistry data  
to improve identification of unknowns

Andrew D McEachran <sup>1, 2</sup>, Ilya Balabin <sup>3</sup>, Tommy Cathey <sup>4</sup>, Thomas R Transue <sup>4</sup>,  
Hussein Al-Ghoul <sup>5</sup>, Chris Grulke <sup>6</sup>, Jon R Sobus <sup>7</sup>, Antony J Williams <sup>8</sup>

*In silico* MS/MS spectra for identifying unknowns: a  
critical examination using CFM-ID algorithms and  
ENTACT mixture samples

Revisiting Five Years of CASMI Contests with EPA

Alex Chao <sup>1, 2</sup>, Hussein Al-Ghoul <sup>1</sup>, Andrew D McEachran <sup>1, 2</sup>, Ilya Balabin <sup>3</sup>, Tom Transue <sup>4</sup>,  
Tommy Cathey <sup>4</sup>, Ilya Balabin <sup>3</sup>, Randolph R Singh <sup>3, 7</sup>, Elin M Ulrich <sup>8</sup>,  
Antony J Williams <sup>9</sup>, Jon R Sobus <sup>10</sup>

Andrew D McEachran <sup>1</sup>, Alex Chao <sup>1</sup>, Hussein Al-Ghoul <sup>1</sup>, Charles Lowe <sup>2</sup>, Christopher Grulke <sup>2</sup>,  
Jon R Sobus <sup>2</sup>, Antony J Williams <sup>2</sup>

# Development of a NTA Workflow



- **Assessing Chemical Hazard: *ToxCast***

- A program to generate chemical toxicity data via high-throughput screening assays

- **Assessing Chemical Exposure: *ExpoCast***

- A program to estimate chemical exposure potential via high-throughput models

- **Assessing Exposure via Consumer Products: *Chemical and Products Database (CPDat)***

- A database mapping chemicals with consumer product usage to assess potential exposures

The ToxCast program for prioritizing toxicity testing of environmental chemicals

David J Dix <sup>1</sup>, Keith A Houck, Matthew T Martin, Ann M Richard, R Woodrow Setzer, Robert J Kavlock

High-throughput models for exposure-based chemical prioritization in the ExpoCast project

John F Wambaugh <sup>1</sup>, R Woodrow Setzer, David M Reif, Sumit Gangwal, Jade Mitchell-Blackwood, Jon A Arnot, Olivier Joliet, Alicia Frame, James Rabinowitz, Thomas B Knudsen, Richard S Judson, Peter Egeghy, Daniel Vallero, Elaine A Cohen Hubal

The Chemical and Products Database, a resource for exposure-relevant data on chemicals in consumer products

Kathie L Dionisio <sup>1</sup>, Katherine Phillips <sup>1</sup>, Paul S Price <sup>1</sup>, Christopher M Grulke <sup>2</sup>, Antony Williams <sup>2</sup>, Derya Biryol <sup>1 3</sup>, Tao Hong <sup>4</sup>, Kristin K Isaacs <sup>1</sup>

# Development of a NTA Workflow



- **Assessing Chemical Hazard: *ToxCast***

- A program to generate chemical toxicity data via high-throughput screening assays

- **Assessing Chemical Exposure: *ExpoCast***

- A program to estimate chemical exposure potential via high-throughput models

- **Assessing Exposure via Consumer Products: *Chemical and Products Database (CPDat)***

- A database mapping chemicals with consumer product usage to assess potential exposures

The ToxCast program for prioritizing toxicity testing of environmental chemicals

David J. Dix<sup>1</sup>, Keith A. Houck, Matthew T. Martin, Ann M. Richard, R. Woodrow Setzer, Robert J. Kavlock

***DSSTox***

High-throughput models for exposure-based chemical prioritization in the ExpoCast project

John F. Wambaugh<sup>1</sup>, R. Woodrow Setzer, David M. Reif, Sumit Gangwal, Jade Mitchell-Blackwood, Jon A. Arnot, Olivier Joliet, Alicia Frame, James Rabinowitz, Thomas B. Knudsen, Richard S. Judson, Peter Egeghy, Daniel Vallerio, Elaine A. Cohen Hubal

The Chemical and Products Database, a resource for exposure-based chemical prioritization in consumer products

***Database API***

Kathie L. Dionisio<sup>1</sup>, Katherine Phillips<sup>1</sup>, Paul S. Price<sup>1</sup>, Christopher M. Grulke<sup>2</sup>, Antony Williams<sup>2</sup>, Derya Biryol<sup>1,3</sup>, Tao Hong<sup>4</sup>, Kristin K. Isaacs<sup>1</sup>

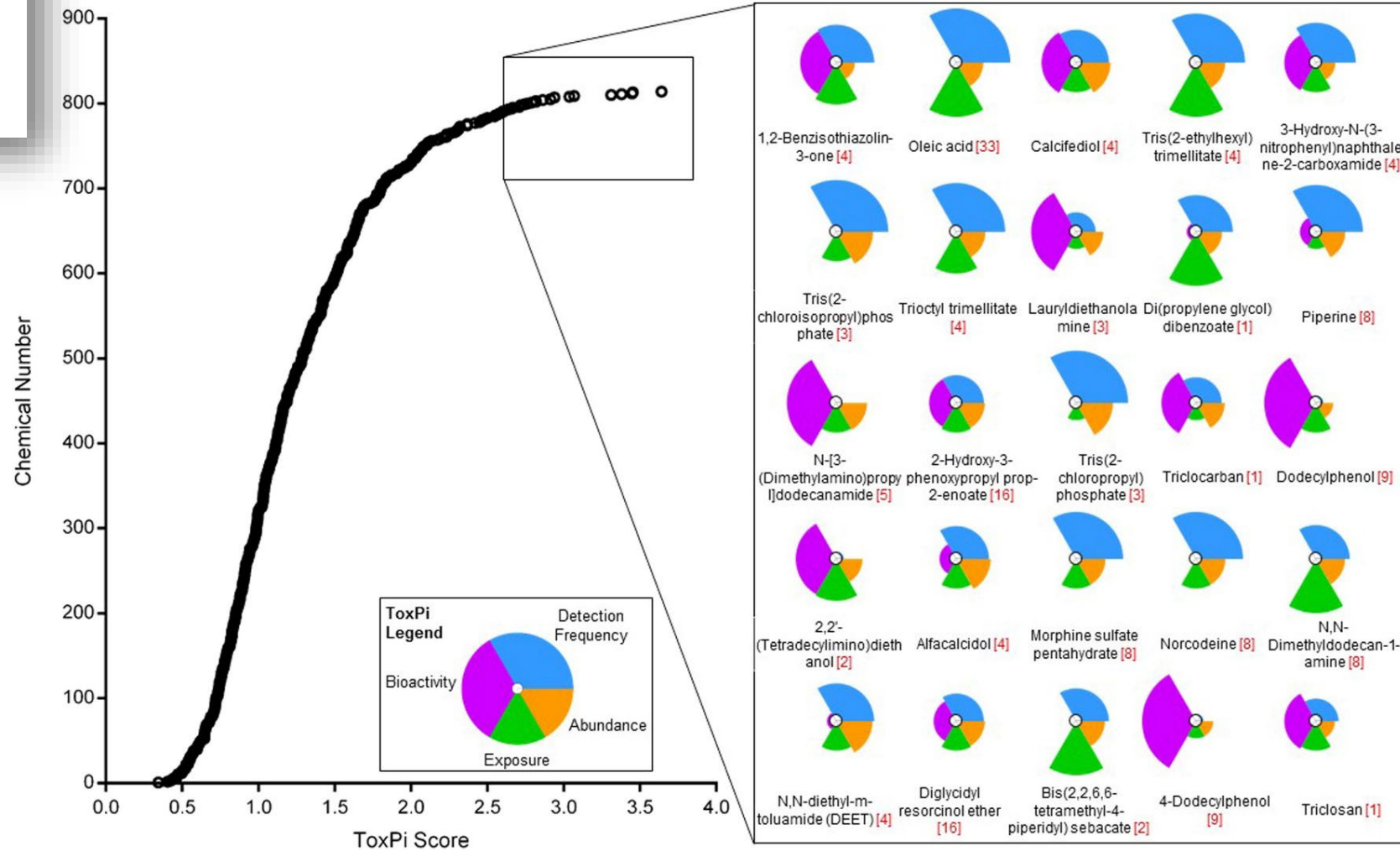


# Chemical Risk Prioritization in House Dust

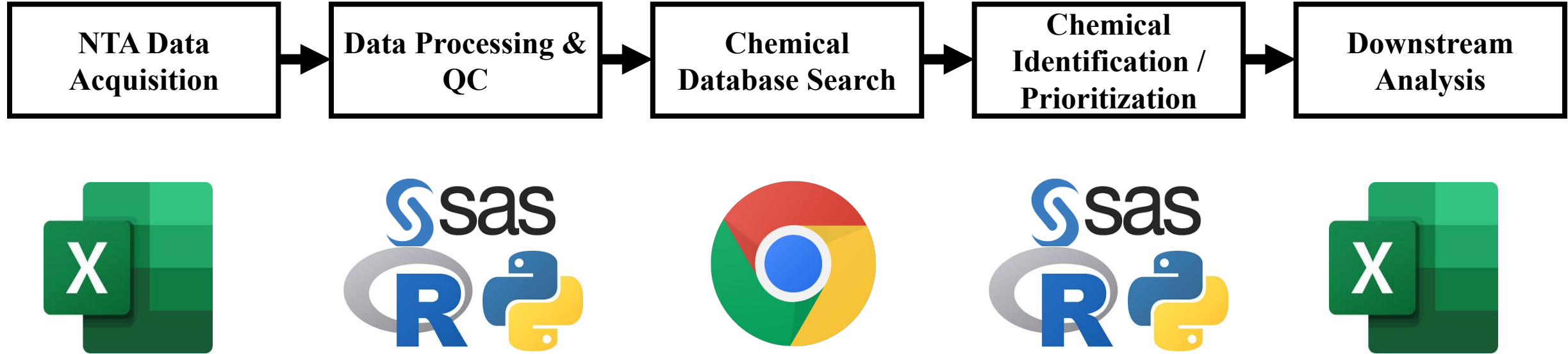
Linking high resolution mass spectrometry data with exposure and toxicity forecasts to advance high-throughput environmental monitoring

Julia E Rager<sup>1</sup>, Mark J Strynar<sup>2</sup>, Shuang Liang<sup>1</sup>, Rebecca L McMahan<sup>1</sup>, Ann M Richard<sup>3</sup>, Christopher M Grulke<sup>4</sup>, John F Wambaugh<sup>3</sup>, Kristin K Isaacs<sup>2</sup>, Richard Judson<sup>3</sup>, Antony J Williams<sup>3</sup>, Jon R Sobus<sup>5</sup>

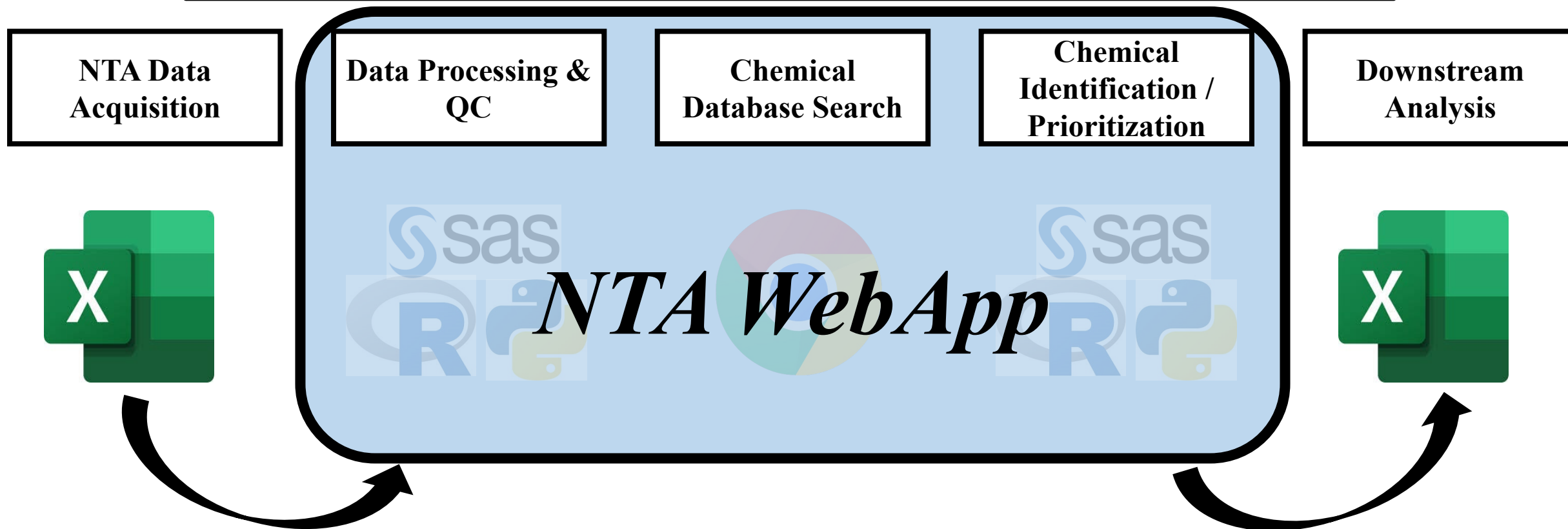
- **Household dust as a metric of exposure**
  - Vacuum dust samples collected from 56 households for NTA
- ***ToxPi Score* to prioritize chemicals for follow-up**
  - ToxCast bioactivity
  - ExpoCast exposure estimates
  - Feature abundance
  - Sample detection frequency



# Consolidation of Workflow into a WebApp



# Consolidation of Workflow into a WebApp



- ***Standardization of processes:*** Single web-accessible point for processing NTA data
- ***Reduction of processing steps:*** Once input, data are carried through whole workflow
- ***Documentation of processes:*** Full workflow tracking for reproducibility and reporting (Input files, processing / search parameters, output results, QC results)

# Data Processing Functionality

## Perform Quality Control on NTA Data Set

*Are features adducts of other features?*

*Are duplicate features present?*

*Are features reproducible across replicates?*

*Are features specific to samples?*

*Are initial annotations robust?*

*Are spiked tracers within allowable tolerances?*

***Flag and/or remove data points from the data set that do not meet QC criteria***

## Cleaning NTA Data for Reporting

- Determine median values for sample replicate groups
- Subtract blank median values for each feature from sample median values

In an NTA data set there may be **thousands** of features observed across **hundreds** of samples:

**Millions of calculations** required to clean and transform each data set

# Chemical Database Search Functionality

## Identify Potential Chemicals for Features

For each feature, search chemical database:

- Pull back all chemicals with a “matching” mass (mass within accuracy tolerance)
- Pull back associated chemical metadata

### *Distributed Structure Searchable Toxicity Database (DSSTox)*

- WebApp has direct interface with DSSTox database

## Prioritize Chemical Candidates for Features

For each candidate, compare associated data to select the most likely candidates (e.g.):

- ToxCast assay data
- ExpoCast exposure estimates
- Data source counts
- MS2 *in silico* spectra scoring

DSSTox currently contains **~906,000** chemicals to search through


A given feature may have **hundreds** of chemical candidates to compare in order to identify likely candidates




# NTA WebApp: Input Page (MS1 data)

## Tracer Input File

- Tracers: Isotopically labelled standards spiked into samples

 United States Environmental Protection Agency

Environmental Topics   Laws & Regulations   About EPA    

NTA: non-targeted analysis of MS data (beta) [Contact Us](#)

**Tools**

- MS1 Tool**
  - Run MS1 Tool
  - MS1 Tool Algorithms
  - MS1 Tool QA/QC
  - MS1 Tool References
  - MS2 CFMID Tool

**Documentation**

- Source Code

### Run NTA MS1 Tool

| Input  | Value   |
|--|---|
| Project name:                                      | <input type="text" value="Example nta"/>                  |
| Positive MPP file (csv):                           | <input type="button" value="Choose File"/> No file chosen |
| Negative MPP file (csv):                           | <input type="button" value="Choose File"/> No file chosen |
| Adduct mass accuracy units:                        | <input type="text" value="ppm"/>                          |
| Adduct mass accuracy:                              | <input type="text" value="10"/>                           |
| Adduct retention time accuracy (mins):             | <input type="text" value="0.05"/>                         |
| Tracer file (csv; optional):                       | <input type="button" value="Choose File"/> No file chosen |
| Tracer mass accuracy units:                        | <input type="text" value="ppm"/>                          |
| Tracer mass accuracy:                              | <input type="text" value="5"/>                            |
| Tracer retention time accuracy (mins):             | <input type="text" value="0.1"/>                          |
| Min sample:blank cutoff:                           | <input type="text" value="3"/>                            |
| Min replicate hits:                                | <input type="range" value="2"/>                           |
| Max replicate CV:                                  | <input type="text" value="0.8"/>                          |
| Parent ion mass accuracy (ppm):                    | <input type="range" value="5"/>                           |
| Discard features below this retention time (mins): | <input type="text" value="0.0"/>                          |
| Search dashboard by:                               | <input type="text" value="mass"/>                         |
| Save top result only?                              | <input type="text" value="no"/>                           |
| DSSTox search batch size (debugging):              | <input type="text" value="150"/>                          |

## NTA Data Input Files:

- Peak-picked, aligned MS1 data (CSV)
- Matrix of feature mass, RT, and sample abundances

## Workflow Parameters:

- Data Processing
- Chemical retrieval

# NTA WebApp: Input Page (MS2 data)

[Environmental Topics](#)[Laws & Regulations](#)[About EPA](#)

## NTA: non-targeted analysis of MS data (beta)

[Contact Us](#)

### Tools

[MS1 Tool](#)[MS2 CFMID Tool](#)

### Documentation

[Source Code](#)

## Run MS2 CFMID Tool

| Input                          | Value  |
|--------------------------------|--|
| Project name:                  | <input type="text" value="Example ms2 nta"/>               |
| Positive mode MS2 files (mgf): | <input type="button" value="Choose Files"/> No file chosen |
| Negative mode MS2 files (mgf): | <input type="button" value="Choose Files"/> No file chosen |
| Precursor mass accuracy (ppm): | <input type="text" value="10"/>                            |
| Fragment mass accuracy (Da):   | <input type="text" value="0.02"/>                          |

☐ Save Metadata?

### NTA Data Input Files:


- Exported MS2 data (MGF format)
- Precursor mass, RT, fragment/intensity pairs

### Workflow Parameters:

- Chemical retrieval
- Spectrum matching

# Submitting an NTA WebApp Job

*In progress...*

 United States  
Environmental Protection  
Agency

Environmental TopicsLaws & RegulationsAbout EPA

Search EPA.gov

NTA: non-targeted analysis of MS data (beta)[Contact Us](#)

Tools

MS1 Tool

Run MS1 Tool

MS1 Tool Algorithms

MS1 Tool QA/QC


MS1 Tool References

MS2 CFMID Tool

Documentation

Source Code

NTA Output



Job ID: XIN4V113

Processing... please wait.

# Submitting an NTA WebApp Job

*In progress...*

*Processing complete and results ready for download*

Environmental Topics

Laws & Regulations

About EPA

Search EPA.gov

NTA: non-targeted analysis of MS data (beta)

## Tools

### MS1 Tool

Run MS1 Tool

MS1 Tool Algorithms

MS1 Tool QA/QC

MS1 Tool References

MS2 CFMID Tool

## Documentation

Source Code

## NTA Output



Job ID: XIN4V113

Processing... please wait.

Environmental Topics

Laws & Regulations

About EPA

Search EPA.gov



NTA: non-targeted analysis of MS data (beta)

[Contact Us](#)

## Tools

### MS1 Tool

Run MS1 Tool

MS1 Tool Algorithms

MS1 Tool QA/QC

MS1 Tool References

MS2 CFMID Tool

## Documentation

Source Code

## NTA Output

Job ID: XIN4V113

Download results:

[Final results](#)

[All files](#)

# WebApp NTA Results Output Format

## *Feature Level Results*

| <i>Feature ID</i> | <i>Mass</i> | <i>Retention Time</i> | <i>Sample 1</i>   | <i>Sample 2</i> | <i>Sample 3</i> |
|-------------------|-------------|-----------------------|---|-----------------|-----------------|
| 1                 | 210.0876    | 6.904999              | <i>Blank-subtracted median abundance values (QC filtered)</i> |                 |                 |
| 2                 | 202.1223    | 7.808004              |   |                 |                 |
| 3                 | 670.5638    | 12.535                |   |                 |                 |
| 4                 | 706.5684    | 12.45099              |   |                 |                 |
| 5                 | 660.5236    | 12.16101              |   |                 |                 |
| 6                 | 616.4656    | 12.817                |   |                 |                 |
| 7                 | 278.147     | 9.584997              |   |                 |                 |
| 8                 | 216.1382    | 8.605996              |   |                 |                 |
| 9                 | 224.1037    | 7.854003              |   |                 |                 |

## *Chemical Level Results*

| <i>Feature ID</i> | <i>Chemical</i>      | <i>MS-Ready Formula</i> | <i>Chem. Data 1</i>   | <i>Chem. Data 2</i> | <i>Chem. Data 3</i> |
|-------------------|----------------------|-------------------------|---|---------------------|---------------------|
| 1                 | Chemical Candidate 1 | <i>MS-Ready Formula</i> | <i>Chemical-specific data and metadata values (ToxCast, ExpoCast, data sources, MS2 scores)</i> |                     |                     |
|                   | Chemical Candidate 2 |                         |   |                     |                     |
|                   | Chemical Candidate 3 |                         |   |                     |                     |
|                   | Chemical Candidate 4 |                         |   |                     |                     |
| 2                 | Chemical Candidate 1 |                         |   |                     |                     |
|                   | Chemical Candidate 2 |                         |   |                     |                     |
|                   | Chemical Candidate 3 |                         |   |                     |                     |
|                   | Chemical Candidate 4 |                         |   |                     |                     |
|                   | Chemical Candidate 5 |                         |   |                     |                     |
|                   |                      |                         |   |                     |                     |

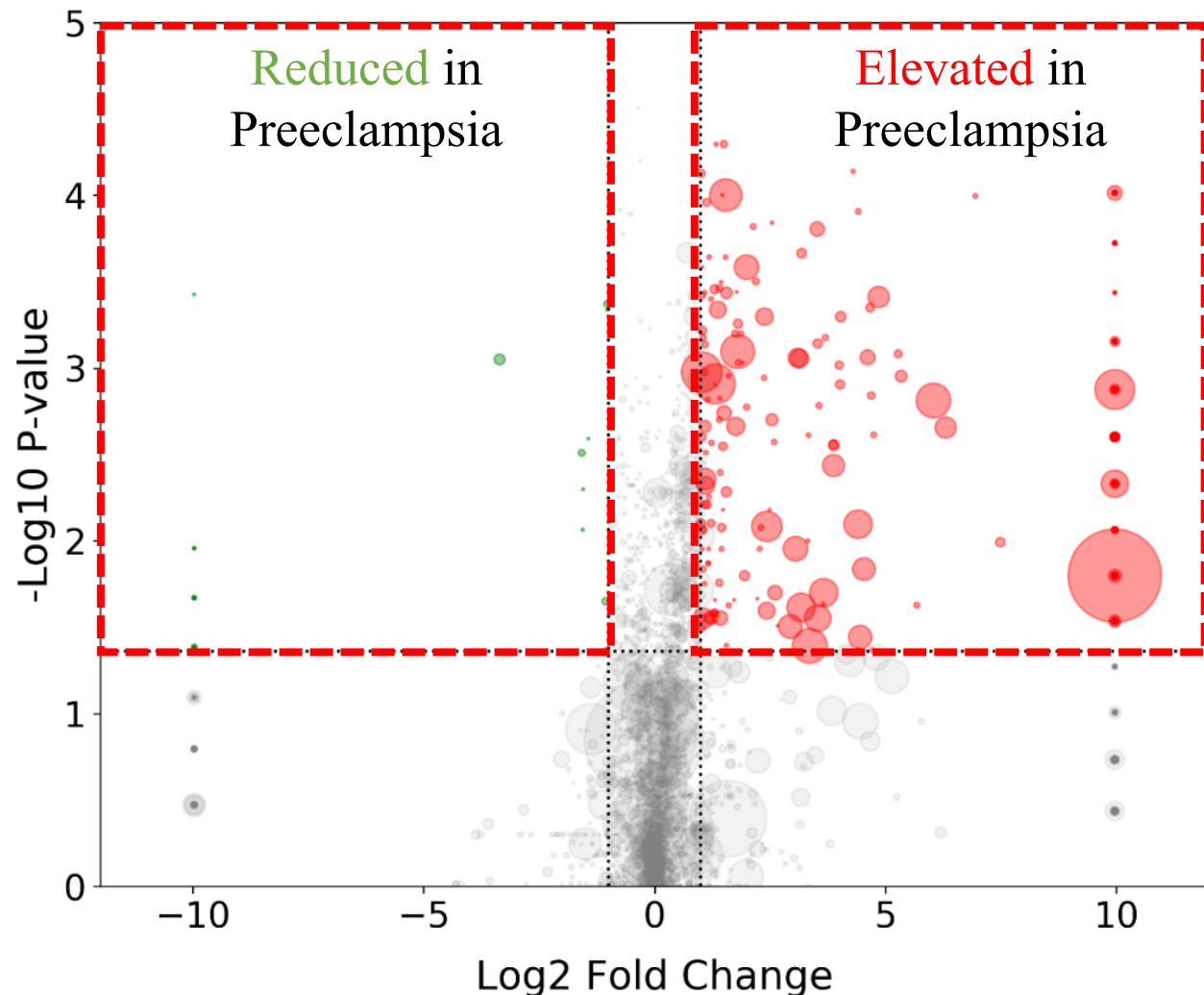


## Intensity



| <i>Tracer Compound</i> | <i>Expected Mass</i> | <i>Expected Retention Time</i> | <i>Sample 1</i>   | <i>Sample 2</i> | <i>Sample 3</i> | <i>Sample 4</i> | <i>Sample 5</i> | <i>Sample 6</i> | <i>Sample 7</i> | <i>Sample 8</i> | <i>Sample 9</i> |
|------------------------|----------------------|--------------------------------|---|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|
| Compound 1             | 210.0876             | 6.90                           | <i>Observed retention time, mass, intensity values</i><br><i>Calculated mass/retention time error, intensity CV's</i> |                 |                 |                 |                 |                 |                 |                 |                 |
| Compound 2             | 202.1223             | 7.81                           |   |                 |                 |                 |                 |                 |                 |                 |                 |
| Compound 3             | 670.5638             | 12.54                          |   |                 |                 |                 |                 |                 |                 |                 |                 |
| Compound 4             | 706.5684             | 12.45                          |   |                 |                 |                 |                 |                 |                 |                 |                 |

# Use Case for Chemical Prioritization



## NTA on placenta samples:

- Normotensive (n = 17) and preeclamptic (n = 18)
- **183 molecular features** found significantly different (~6000 potential candidates)
- Feature chemicals prioritized for targeted confirmatory work via:
  - Reference MS2 spectrum match
  - *In silico* MS2 spectrum match
  - Data source counts
  - Consumer product database presence (CPCat)
  - **46 chemicals** prioritized / acquired
- **25 chemicals** confirmed via targeted analyses



# What's Next?

- **Development of tools for improved NTA results**
  - Database incorporation of publicly available MS2 spectra
  - The Hazard Comparison Dashboard: database aggregation of publicly available toxicity data (*public version to be online shortly*)
  - Semi-quantitation methods: Generation of concentration estimates and uncertainty bounds from NTA data (*manuscript submitted*)
  - Compound method amenability predictions (*manuscript published*)
  - LC-MS retention time predictions (*manuscript published*)

## Predicting compound amenability with liquid chromatography-mass spectrometry to improve non-targeted analysis

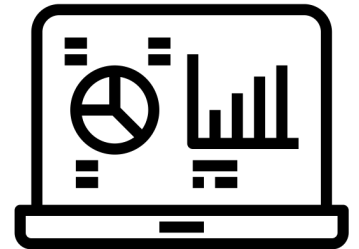
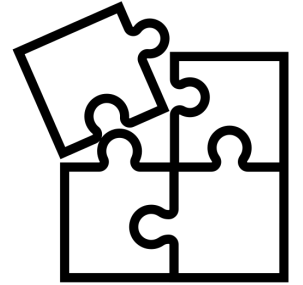
Charles N Lowe<sup>1</sup>, Kristin K Isaacs<sup>2</sup>, Andrew McEachran<sup>3</sup>, Christopher M Grulke<sup>2</sup>, Jon R Sobus<sup>2</sup>, Elin M Ulrich<sup>2</sup>, Ann Richard<sup>2</sup>, Alex Chao<sup>2</sup>, John Wambaugh<sup>2</sup>, Antony J Williams<sup>2</sup>

## Development and Application of Liquid Chromatographic Retention Time Indices in HRMS-Based Suspect and Nontarget Screening

Reza Aalizadeh<sup>1</sup>, Nikiforos A Alygizakis<sup>1 2</sup>, Emma L Schymanski<sup>3 4</sup>, Martin Krauss<sup>5</sup>, Tobias Schulze<sup>5</sup>, María Ibáñez<sup>6</sup>, Andrew D McEachran<sup>7</sup>, Alex Chao<sup>7</sup>, Antony J Williams<sup>7</sup>, Pablo Gago-Ferrero<sup>8 9</sup>, Adrian Covaci<sup>10</sup>, Christoph Moschet<sup>11</sup>, Thomas M Young<sup>11</sup>, Juliane Hollender<sup>4 12</sup>, Jaroslav Slobodnik<sup>2</sup>, Nikolaos S Thomaidis<sup>1</sup>

# The NTA WebApp: Summary

- **The NTA WebApp is a synthesis of diverse work**
  - Multiple databases are integrated directly into workflow
  - More data = greater ability for chemical assignment and prioritization
- **The NTA WebApp is a step towards standardization**
  - Reproducible work through a tool containing entire workflow
  - Transparency and tracking of workflow
- **Upcoming developments**
  - New modules / databases in development for incorporation
  - *Manuscript in draft*: Public-facing version of WebApp to accompany publication (expected Fall 2022)



# Acknowledgements



Hussein Al-Ghoul  
Kathie Dionisio  
Louis Groff  
Jarod Grossman  
Chris Grulke  
Kristin Isaacs  
Charles Lowe  
James McCord  
Andrew McEachran  
Seth Newton  
Allison Phillips  
Katherine Phillips  
Marie Russell  
John Sloop  
Elin Ulrich  
John Wambaugh

