# Agilent LC/Q-TOF Simplified Workflow

**Experimental Acquisition**

**DB & Library Matching**

**Data Analysis**

Sample Extracts

Chemical Database

MS$^1$ Feature Table

LC/Q-TOF HRMS

DB MS-Ready Structures

Filtered Feature Table

Chemical Candidate Table

MS$^2$ Acquisition

MS$^1$ Acquisition

DB MS-Ready Formula & Monoisotopic Mass

Aggregated Match Table

MS$^2$ .d Files

Reference MS$^2$ Spectra

MS$^2$ Reference Matches

MS$^2$ .mgf Files

*in silico* MS$^2$ Spectra

MS$^2$ *in silico* Matches

# Chemical Database = DSSTox

# MS-Ready Structures

Spiked Substance: **Tamoxifen**

Spiked Substance: **Tamoxifen citrate**

Predicted Formula for Observed Molecular Feature: $C_{26}H_{29}NO$

Dashboard Search

**DTXSID1034187**

**DTXSID8021301**

1st: **DTXSID1034187**     2nd: **DTXSID8021301**

**MS-Ready Processing**

**DTXCID9014187**     **DTXCID9014187**     **DTXCID9014187**     **DTXCID9014187**

# Dashboard Access

# Agilent LC/Q-TOF Simplified Workflow

**Experimental Acquisition**

**DB & Library Matching**

**Data Analysis**

Sample Extracts → LC/Q-TOF HRMS → MS$^2$ Acquisition, MS$^1$ Acquisition

MS$^2$ Acquisition → MS$^2$ .d Files → MS$^2$ .mgf Files

Chemical Database → DB MS-Ready Structures → DB MS-Ready Formula & Monoisotopic Mass

MS$^1$ Acquisition → DB MS-Ready Formula & Monoisotopic Mass

MS$^2$ .d Files → Reference MS$^2$ Spectra

MS$^2$ .mgf Files → in silico MS$^2$ Spectra

MS$^1$ Feature Table → Filtered Feature Table → Chemical Candidate Table → Aggregated Match Table

DB MS-Ready Formula & Monoisotopic Mass → MS$^1$ Feature Table

Reference MS$^2$ Spectra → MS$^2$ Reference Matches → Aggregated Match Table

in silico MS$^2$ Spectra → MS$^2$ in silico Matches → Aggregated Match Table

# EPA NTA WebApp

**Feature Removal:**

1) Duplicate features
2) Non-reproducible features
3) Blank features (sample:blank)
4) Non-responsive features (dilutions)

**Feature Flagging:**

1) Multi-mode hits (+ and -)
2) Meas. precision (CV threshold)
3) Formula match (score ≥ threshold)
4) Negative mass defect
5) Halogenation
6) Has/is adduct
7) Has/is neutral loss
8) Has/is multimer

**Dashboard Integration:**

1) Data source & pub counts
2) Bioactivity & exposure levels
3) Presence on lists
4) Product & use categories

# Agilent LC/Q-TOF Simplified Workflow

# Generation of *in silico* Spectra

McEachran, Andrew D., et al. *Scientific data* 6.1 (2019): 1-9
Allen, Felicity, et al. *Metabolomics* 11.1 (2015): 98-110.

# CFM-ID Database Matching



1. Query database by mass

MGF file ···· Exp MS2 Spectrum (Mass = 356.119) → CFM-ID Database

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

Alex Chao[1,2] · Hussein Al-Ghoul[1,2] · Andrew D. McEachran[1,3] · Ilya Balabin[4] · Tom Transue[4] · Tommy Cathey[4] · Jarod N. Grossman[2,3] · Randolph R. Singh[1,5] · Elin M. Ulrich[6] · Antony J. Williams[7] · Jon R. Sobus[6]

Retrieve candidate compounds within mass window

Candidate 1 $C_{19}H_{20}N_2O_3S$ (Mass = 356.119) ···· In silico MS2 Spectra (CE 10, 20, 40)

Candidate 2 $C_{19}H_{20}N_2O_3S$ (Mass = 356.119) ···· In silico MS2 Spectra (CE 10, 20, 40)

Candidate 3 $C_{21}H_{21}ClO_3$ (Mass = 356.118) ···· In silico MS2 Spectra (CE 10, 20, 40)

2. Score in silico spectra

### CFM-ID Scores

|  | in silico CE 10 | in silico CE 20 | in silico CE 40 |
|---|---|---|---|
| Candidate 1 | 0.5 | 0.3 | 0.1 |
| Candidate 2 | 0.2 | 0.1 | 0.02 |
| Candidate 3 | 0.1 | 0.05 | 0.01 |

# CFM-ID Database Matching (w/ Formula Information)

*1. Query database by mass*

MGF file ···· Exp MS2 Spectrum (Mass = 356.119) → CFM-ID Database

*Retrieve candidate compounds within mass window*

→ $C_{19}H_{20}N_2O_3S$ →

*Formula Identified*

*Filter candidates by formula*

Candidate 1 $C_{19}H_{20}N_2O_3S$ (Mass = 356.119) ···· *In silico* MS2 Spectra (CE 10, 20, 40)

Candidate 2 $C_{19}H_{20}N_2O_3S$ (Mass = 356.119) ···· *In silico* MS2 Spectra (CE 10, 20, 40)

Candidate 3 $C_{21}H_{21}ClO_3$ (Mass = 356.118) ···· *In silico* MS2 Spectra (CE 10, 20, 40)

*2. Score in silico spectra*

## CFM-ID Scores

| | *in silico* CE 10 | *in silico* CE 20 | *in silico* CE 40 |
|---|---|---|---|
| Candidate 1 | 0.5 | 0.3 | 0.1 |
| Candidate 2 | 0.2 | 0.1 | 0.02 |
| Candidate 3 | 0.1 | 0.05 | 0.01 |

# CFM-ID Database Matching (w/ Multiple CE$_{experimental}$)

# CFM-ID Scoring Approaches

# EPA'S Non-Targeted Analysis Collaborative Trial



## The Trial Mixtures:



10 Mixtures ranging from 95 to 365 compounds
(Total: 1,269 unique compounds)

*"Pass" compounds = 377 with MS2 data*

## EPA Setup:



Agilent 1290 UPLC
Agilent 6530B Q-TOF with ESI source

Ulrich, Elin M., et al. *Analytical and bioanalytical chemistry* 411.4 (2019): 853-866.
Sobus, Jon R., et al. *Analytical and bioanalytical chemistry* 411.4 (2019): 835-851.

# Reference vs. *in silico* Library Coverage



PCDL | CFM-ID

88  111  77

101

"Pass" Compounds

| MS2 Library | % of "Pass" Compounds Identified |
|---|---|
| Agilent PCDL | 53% |
| CFM-ID Top Hit | 50% |
| PCDL and/or CFM-ID Top Hit | 73% |

PCDL → Agilent reference $MS^2$ library

"Pass" compounds (n=377) → ENTACT chemicals observed with $MS^2$ data

# NTA Workflows: Using CFM-ID Results as Filters

**Score**
Filter out candidates
below score cutoff

**Variability in score
distribution**

**Rank**
Filter out candidates
above rank cutoff

**Variability in number of
candidate compounds**

MS2 Spectrum 1

✕

Candidate Scores

MS2 Spectrum 1

n = 10

Candidate Scores

Filter by Top 20

MS2 Spectrum 2

✕

Candidate Scores

MS2 Spectrum 2

n = 500

✕

Candidate Scores

# Normalizing CFM-ID Results Values

**Score Quotient**
Normalize score to the highest candidate compound score

**Score Percentile**
Normalize rank to the number of candidate compounds

|  | Rank | CFM-ID Score | Maximum Score | Score Quotient | Score Percentile |
|---|---|---|---|---|---|
| **Candidate Compound 1** | 1 | 0.5 | 0.5 | 1 | 100 |
| **Candidate Compound 2** | 2 | 0.4 | 0.5 | 0.8 | 80 |
| **Candidate Compound 3** | 3 | 0.39 | 0.5 | 0.78 | 60 |
| **Candidate Compound 4** | 4 | 0.1 | 0.5 | 0.2 | 40 |
| **Candidate Compound 5** | 5 | 0.05 | 0.5 | 0.1 | 20 |

Score Quotient = Score / Maximum Score

# NTA Workflows: Using CFM-ID Normalized Results as Filters

**Score Quotient**
Filter out candidates below score quotient cutoff

**Score Percentile**
Filter out candidates below percentile cutoff

*Score quotient cutoff = 0.5*
*Keep candidates scoring at least half of max score*

*Score percentile cutoff = 0.5*
*Keep the top 50% of candidates*

MS2 Spectrum 1

Candidate Scores

MS2 Spectrum 1

Candidate Scores

MS2 Spectrum 2

Candidate Scores

MS2 Spectrum 2

Candidate Scores

# Applying Cut-off Filters to Data

| | CFM-ID Score | Maximum Score | Score Quotient |
|---|---|---|---|
| Candidate Compound 1 | 0.5 | 0.5 | 1 |
| Candidate Compound 2 | 0.4 | 0.5 | 0.8 |
| Candidate Compound 3 | 0.39 | 0.5 | 0.78 |
| Candidate Compound 4 | 0.1 | 0.5 | 0.2 |
| Candidate Compound 5 | 0.05 | 0.5 | 0.1 |

# Applying Cut-off Filters to Data

| | CFM-ID Score | Maximum Score | Score Quotient |
|---|---|---|---|
| Candidate Compound 1 | 0.5 | 0.5 | 1 |
| Candidate Compound 2 | 0.4 | 0.5 | 0.8 |
| Candidate Compound 3 | 0.39 | 0.5 | 0.78 |
| Candidate Compound 4 | 0.1 | 0.5 | 0.2 |
| Candidate Compound 5 | 0.05 | 0.5 | 0.1 |

▲ True Compound

● Other Candidate Compounds

| True Positives | |
|---|---|
| False Negatives | |
| True Negatives | |
| False Positives | |

Score Quotient

# Applying Cut-off Filters to Data

| | CFM-ID Score | Maximum Score | Score Quotient |
|---|---|---|---|
| Candidate Compound 1 | 0.5 | 0.5 | 1 |
| Candidate Compound 2 | 0.4 | 0.5 | 0.8 |
| Candidate Compound 3 | 0.39 | 0.5 | 0.78 |
| Candidate Compound 4 | 0.1 | 0.5 | 0.2 |
| Candidate Compound 5 | 0.05 | 0.5 | 0.1 |

▲ True Compound

● Other Candidate Compounds

| True Positives | 1 |
|---|---|
| False Negatives | 0 |
| True Negatives | 0 |
| False Positives | 4 |

# Applying Cut-off Filters to Data

| | CFM-ID Score | Maximum Score | Score Quotient |
|---|---|---|---|
| Candidate Compound 1 | 0.5 | 0.5 | 1 |
| Candidate Compound 2 | 0.4 | 0.5 | 0.8 |
| Candidate Compound 3 | 0.39 | 0.5 | 0.78 |
| Candidate Compound 4 | 0.1 | 0.5 | 0.2 |
| Candidate Compound 5 | 0.05 | 0.5 | 0.1 |

▲ True Compound

● Other Candidate Compounds

| True Positives | 1 |
|---|---|
| False Negatives | 0 |
| True Negatives | 2 |
| False Positives | 2 |



Score Quotient
Cut-off = 0.5

# Applying Cut-off Filters to Data

| | CFM-ID Score | Maximum Score | Score Quotient |
|---|---|---|---|
| Candidate Compound 1 | 0.5 | 0.5 | 1 |
| Candidate Compound 2 | 0.4 | 0.5 | 0.8 |
| Candidate Compound 3 | 0.39 | 0.5 | 0.78 |
| Candidate Compound 4 | 0.1 | 0.5 | 0.2 |
| Candidate Compound 5 | 0.05 | 0.5 | 0.1 |

▲ True Compound

● Other Candidate Compounds

| True Positives | 0 |
|---|---|
| False Negatives | 1 |
| True Negatives | 3 |
| False Positives | 1 |



Score Quotient Cut-off = 0.9

# Balancing Cut-offs

$$True\ Positive\ Rate\ (TPR) = \frac{TP}{TP + FN}$$

How many of the true compounds are we keeping?

$$False\ Positive\ Rate\ (FPR) = \frac{FP}{FP + TN}$$

How much of the junk are we getting rid of?

# Quotient Vs. Percentile Cutoffs



Global ROC Curves (All ENTACT Mixtures)

# Quotient Vs. Percentile Cutoffs



Global ROC Curves (All ENTACT Mixtures)

Legend:
- Quotient (by formula)
- Percentile (by formula)
- Quotient (by mass)
- Percentile (by mass)

X-axis: False Positive Rate
Y-axis: True Positive Rate

Score Quotient Cut-off

▲ True Compounds
● Other Candidate Compounds

# Quotient Vs. Percentile Cutoffs



Global ROC Curves (All ENTACT Mixtures)

Legend (ROC plot):
- Quotient (by formula) — magenta solid
- Percentile (by formula) — green solid
- Quotient (by mass) — magenta dotted
- Percentile (by mass) — green dotted

Axes: True Positive Rate (y), False Positive Rate (x)

Right plot legend:
- Score Quotient (y-axis, 0 to 1)
- Score Quotient Cut-off (dashed line)
- ▲ True Compounds
- ● Other Candidate Compounds

# Quotient Vs. Percentile Cutoffs



Global ROC Curves (All ENTACT Mixtures)

# Quotient Vs. Percentile Cutoffs



Global ROC Curves (All ENTACT Mixtures)

Cut-off Values for Global TPR = 0.9

|  | Cut-off value |
|---|---|
| Quotient (by formula) | 0.18 |
| Percentile (by formula) | 38 |
| Quotient (by mass) | 0.13 |
| Percentile (by mass) | 32 |

*Apply to individual ENTACT mixtures*

# CFM-ID Cut-off Filtering: Individual ENTACT Mixtures

Experimental Acquisition | Database/Library Matching | Data Analysis

Sample Extracts
A
LC-QTOF/MS
B → MS2 Acquisition
B → MS1 Acquisition
C → MS2 Acquisition .d Files
F → MS2 Exported .mgf Files

DSSTox Database[1] (~875,000 Substances)
G
DSSTox MS-Ready Structures[2]
G
DSSTox MS-Ready Formulae
H
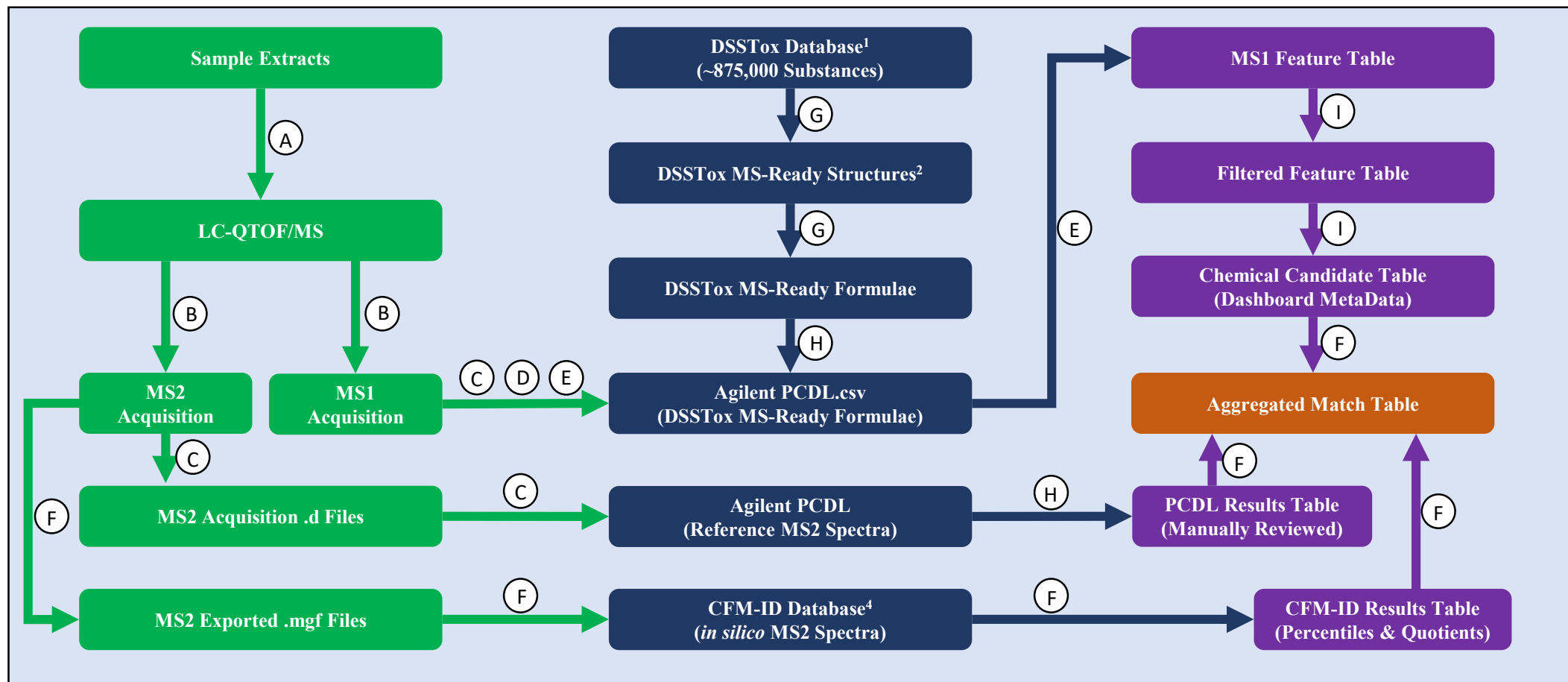Agilent PCDL.csv (DSSTox MS-Ready Formulae)
Agilent PCDL (Reference MS2 Spectra)
CFM-ID Database[4] (*in silico* MS2 Spectra)

MS1 Feature Table
I
Filtered Feature Table
I
Chemical Candidate Table (Dashboard MetaData)
F
Aggregated Match Table
PCDL Results Table (Manually Reviewed)
CFM-ID Results Table (Percentiles & Quotients)

C D E
E
F
H
F
F
F

**Software & Tools**

A) Excel Macro (naming & randomization)
B) Agilent MassHunter Data Acquisition
C) Agilent MassHunter Qualitative Analysis
D) Agilent Profinder (peak picking & alignment)
E) Agilent Mass Profiler Professional (formula matching)
F) Python Script
G) CompTox Chemicals Dashboard[3]
H) Excel
I) EPA NTA WebApp

**General Examples**

Sobus et al. https://link.springer.com/article/10.1007%2Fs00216-018-1526-4

Newton et al. https://www.sciencedirect.com/science/article/pii/S026974911732691X?via%3Dihub

Hedgespeth et al. https://www.sciencedirect.com/science/article/pii/S004896971933298X?via%3Dihub

**Specific References**

[1]Grulke et al. https://www.sciencedirect.com/science/article/pii/S2468111319300234

[2]McEachran et al. https://jcheminf.biomedcentral.com/articles/10.1186/s13321-018-0299-2

[3]Williams et al. https://jcheminf.biomedcentral.com/articles/10.1186/s13321-017-0247-6

[4]McEachran et al. https://www.nature.com/articles/s41597-019-0145-z

# Contributing Researchers


Credit: the Research Triangle Foundation

**EPA ORD**
Hussein Al-Ghoul*
Alex Chao*
Jarod Grossman*
Kristin Isaacs
Sarah Laughlin*
Charles Lowe
James McCord
Jeff Minucci
Seth Newton
Katherine Phillips
Tom Purucker
Randolph Singh*
Mark Strynar
Elin Ulrich

\* = ORISE/ORAU

**EPA ORD (cont.)**
Chris Grulke
Kamel Mansouri*
Andrew McEachran*
Ann Richard
John Wambaugh
Antony Williams

**Agilent**
Jarod Grossman
Andrew McEachran

**GDIT**
Ilya Balabin
Tom Transue
Tommy Cathey

# Questions?



sobus.jon@epa.gov

*The views expressed in this presentation are those of the authors and do not necessarily represent the views or policies of the U.S. Environmental Protection Agency.*



Credit: the Research Triangle Foundation