# Improved compound identification in non-target screening

**David Cox**, Craig Butt, Janna Anichina, Adrian Taylor | SCIEX

ASMS 2022

SCIEX
The Power of Precision

**Searching chemical databases** by m/z or molecular formula often generates a **very long list** of possible compounds. The **correct answer** is often the one that is **easiest to read**.

# Compound ID difficulty: targeted, suspect, unknown

**Targeted**
Compound is in library, has been measured by this lab
- Mass error
- Retention time
- Ion ratio
- Fragment mass error
- Isotope match
- Spectral library match

**Suspect**
Compound is in library, has **never** been measured by this lab before
- Mass error
- Fragment mass error
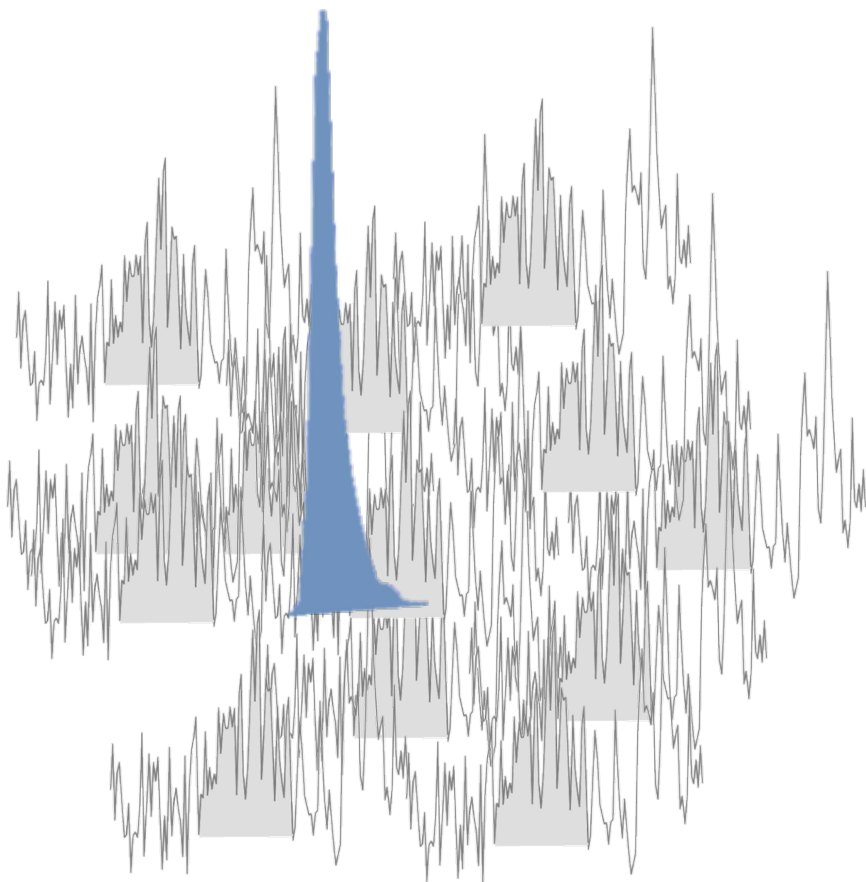- Isotope match
- Spectral library match

**Unknown**
Compound is in PubChem database, no spectrum ever measured
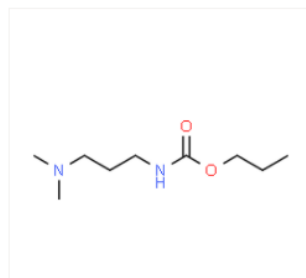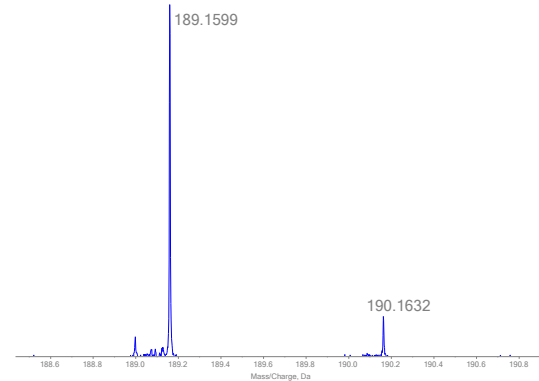- Mass error
- Isotope match
- Predict formula

**General Unknown**
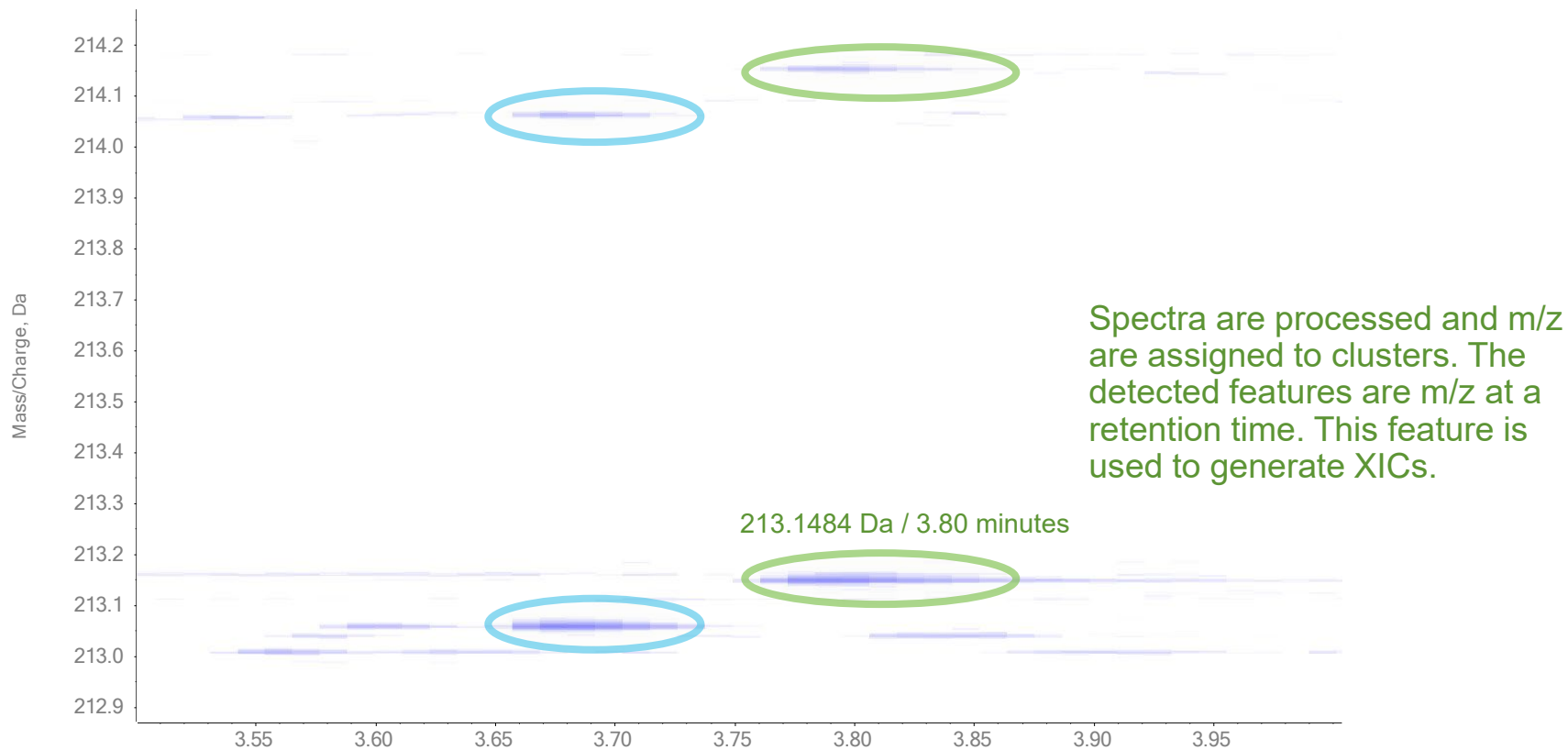Not in PubChem database

# Finding "real" features

# Identifying a compound



**Propamocarb**

| | |
|---|---|
| Molecular Formula | $C_9H_{20}N_2O_2$ |
| Average mass | 188.267 Da |
| Monoisotopic mass | 188.152481 Da |
| ChemSpider ID | 30114 |

# Feature finding



Spectra are processed and m/z are assigned to clusters. The detected features are m/z at a retention time. This feature is used to generate XICs.
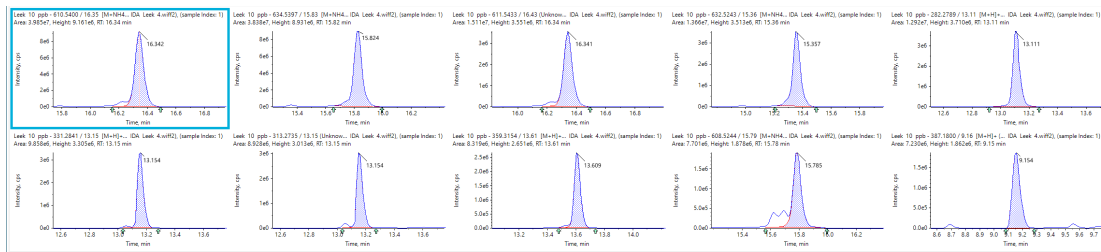
213.1484 Da / 3.80 minutes

# What is a "good" peak?

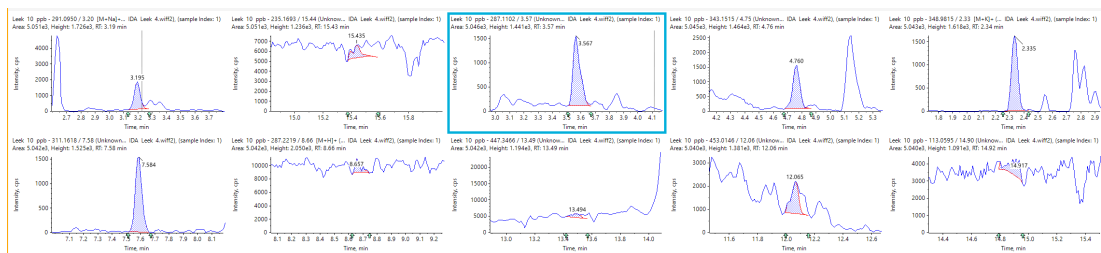I can't describe it, but I know it when I see it …

# Using intensity to sort "good" from "bad" peaks

While features with a high intensity are typically real, and features with a low intensity are often junk, where do you draw a cut-off?
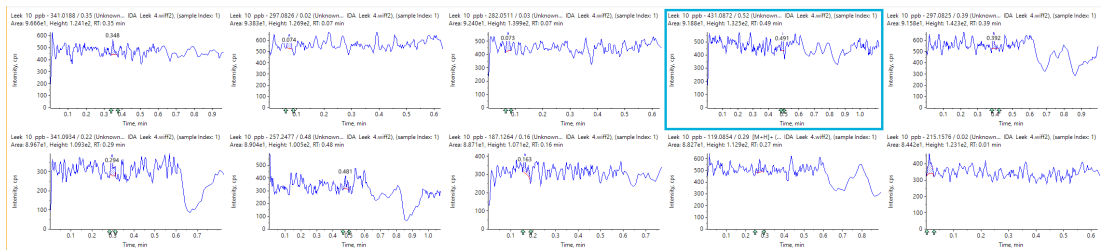
1st

8 971

13 394

# What is a "good" peak?

We can use other measures.
But how do we decide on cut-offs? How to keep track of many measures at once?

## Intensity

S:N, Area Height, Region Height, Quality, Retention Time Delta min_, Total Width, Width at 50 , Baseline Delta Height, Width at 5, Width at 10, Slope of Baseline, Tailing Factor, Asymmetry Factor, Points Across Baseline, Points Across Half Height, AutoPeak Asymmetry, AutoPeak Candidate Model Quality, AutoPeak Group Confidence, AutoPeak Integration Quality, AutoPeak Num Peaks, AutoPeak Peak Width Confidence

# Enabling machine learning with Analytics prototype in SCIEX OS software

- Modified version of processing software

- Used Microsoft's open source ML.NET framework

- Implemented a classifier that uses the
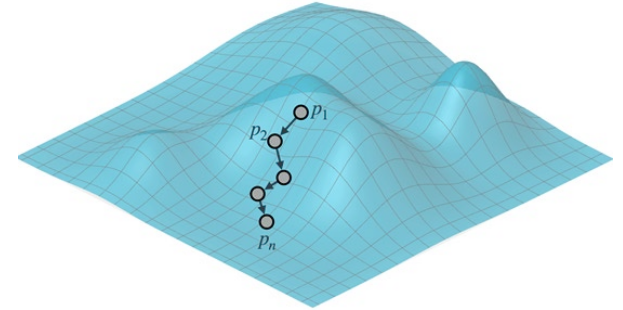  Stochastic Dual Coordinate Ascent (SDCA) method

  – Trained using a custom column in the results table, where a
    user enters classifications ("good" or "bad")

  – Results of classification are output to new columns
    (predicted classification and score)



https://livebook.manning.com/book/grokking-machine-learning/appendix-b/v-15/29

- For more details on ML.NET and machine learning in general:

  https://dotnet.microsoft.com/en-us/apps/machinelearning-ai/ml-dotnet
  https://rubikscode.net/2021/02/01/machine-learning-with-ml-net-ultimate-guide-to-classification/
  https://www.kdnuggets.com/2020/05/5-concepts-gradient-descent-cost-function.html
  https://www.jmlr.org/papers/v14/gonen13a.html

# Review a few peaks to train the model

With minimal keystrokes (↑ ↓ g b), classification of a **few hundred** peaks is fast and easy.
This set of data is used to train the model so it can be used to automatically classify the remaining ~13 000 features.

# Machine learning classification finds more "real" features

Deciding what to focus our attention on is much easier now. We can still find true features, even at low intensities.

1st

8 971

13 394

# How well does machine learning perform compared to other scoring?

Random score

True positive rate

1

0.5

0

0          1

False positive rate

A **randomly** generated number has **no predictive** power. It generates as many true positives as false positives.

Perfect score

0          1

A (mythical) perfect scoring algorithm would obtain all true positives before any false positives occur.

# Intensity sorting is better than a random score ¯\\_(ツ)_/¯

**Random score**

True positive rate

1

0.5

0

0          1

False positive rate

**Intensity**

0          1

For example, classifying peaks above 0.004% relative intensity obtained **90% of the 559 true peaks** in the test set, but generated **222 false positives**.

**Perfect score**

0          1

# Machine learning scoring performs better than intensity scoring



Random score

Intensity

ML trained on **83** examples

ML trained on **262** examples

Perfect score

True Positive Rate

False positive rate

ML score above 50% obtained **90% of the 559 true peaks** in the test set, and only generated **56 false positives**.

# Finding "real" features



# Identifying a compound



**Propamocarb**

| | |
|---|---|
| Molecular Formula | $C_9H_{20}N_2O_2$ |
| Average mass | 188.267 Da |
| Monoisotopic mass | 188.152481 Da |
| ChemSpider ID | 30114 |

# Library searching is common for identification

This only works if the MS/MS for the compound has been seen before, and has been entered into the spectral library you have.



Acquired        Library        Fit / Purity

A               A    100 / 100

                X    100 / 30

                Y    50 / 30

                Z    30 / 30

# Searching a formula in online databases

$$C_{14}H_{15}ClN_4O_3$$



Sorted by **Reference Count**

| | |
|---|---|
| 1 | N-(2-Acetamidoethyl)-4-[5-(chloromethyl)-1,2,4-oxadiazol-3-yl]benzamide |
| 2 | (4-Chloro-1-methyl-1H-pyrazol-5-yl)[4-(2-furoyl)-1-piperazinyl]methanone |
| 3 | 4-Chloro-N-[(1-ethyl-3-methyl-1H-pyrazol-4-yl)methyl]-3-nitrobenzamide |
| 4 | N-(3-Chloro-2-methylphenyl)-2-(3,5-dimethyl-4-nitro-1H-pyrazol-1-yl)acetamide |
| 5 | 2-Chloro-N-[(1-ethyl-5-methyl-1H-pyrazol-4-yl)methyl]-4-nitrobenzamide |
| 6 | N-(6-Amino-1-benzyl-2,4-dioxo-1,2,3,4-tetrahydro-5-pyrimidinyl)-2-chloro-N-methylacetamide |
| 7 | N-(2-Chlorobenzyl)-2-(3,5-dimethyl-4-nitro-1H-pyrazol-1-yl)acetamide |
| 8 | 5-Chloro-N-[(1-ethyl-5-methyl-1H-pyrazol-4-yl)methyl]-2-nitrobenzamide |
| 9 | 4-Chloro-N-[(1-ethyl-5-methyl-1H-pyrazol-4-yl)methyl]-3-nitrobenzamide |
| 10 | 4-Chloro-N-[(1-ethyl-3-methyl-1H-pyrazol-4-yl)methyl]-2-nitrobenzamide |
| ... | |

# Searching a formula in online databases

$$C_{14}H_{15}ClN_4O_3$$

Sorted by **Reference Count**

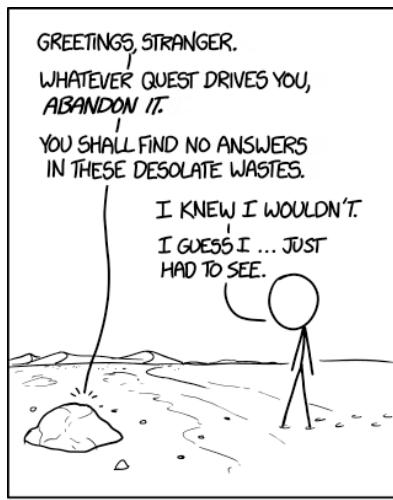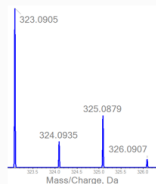| | |
|---|---|
| 1 | N-(2-Acetamidoethyl)-4-[5-(chloromethyl)-1,2,4-oxadiazol-3-yl]benzamide |
| 2 | (4-Chloro-1-methyl-1H-pyrazol-5-yl)[4-(2-furoyl)-1-piperazinyl]methanone |
| 3 | 4-Chloro-N-[(1-ethyl-3-methyl-1H-pyrazol-4-yl)methyl]-3-nitrobenzamide |
| 4 | N-(3-Chloro-2-methylphenyl)-2-(3,5-dimethyl-4-nitro-1H-pyrazol-1-yl)acetamide |
| 5 | 2-Chloro-N-[(1-ethyl-5-methyl-1H-pyrazol-4-yl)methyl]-4-nitrobenzamide |
| 6 | N-(6-Amino-1-benzyl-2,4-dioxo-1,2,3,4-tetrahydro-5-pyrimidinyl)-2-chloro-N-methylacetamide |
| 7 | N-(2-Chlorobenzyl)-2-(3,5-dimethyl-4-nitro-1H-pyrazol-1-yl)acetamide |
| 8 | 5-Chloro-N-[(1-ethyl-5-methyl-1H-pyrazol-4-yl)methyl]-2-nitrobenzamide |
| 9 | 4-Chloro-N-[(1-ethyl-5-methyl-1H-pyrazol-4-yl)methyl]-3-nitrobenzamide |
| 10 | 4-Chloro-N-[(1-ethyl-3-methyl-1H-pyrazol-4-yl)methyl]-2-nitrobenzamide |
| ... | |
| ... | |
| 301 | 4-Chloro-N'-[(Z)-(3,4-dimethoxyphenyl)methylene]-1-methyl-1H-pyrazole-5-carbohydrazide |
| 302 | 2-Oxo-2-[(1,3,5-trimethyl-1H-pyrazol-4-yl)amino]ethyl 6-chloronicotinate |
| 303 | 2-Oxo-2-[(1,3,5-trimethyl-1H-pyrazol-4-yl)amino]ethyl 2-chloronicotinate |
| 304 | **Cycloxaprid** |
| 305 | N-(4-Chlorophenyl)-2-(3,5-dimethyl-4-nitro-1H-pyrazol-1-yl)propanamide |
| 306 | N-(2-Chloro-5-nitrobenzyl)-3-ethyl-N-methyl-1H-pyrazole-5-carboxamide |
| 307 | Ethyl 2-[(5-chloro-3-methyl-1-phenyl-1H-pyrazol-4-yl)carbonyl]hydrazinecarboxylate |
| 308 | N-[3-Chloro-2-(dimethylamino)phenyl]-2-(2,4-dioxo-3,4-dihydro-1(2H)-pyrimidinyl)acetamide |
| 309 | N-[3-Chloro-4-(1H-pyrazol-1-yl)phenyl]-N'-(1-hydroxy-2-propanyl)ethanediamide |
| 310 | 4-Chloro-2-{2-[(1,3-dimethyl-1H-pyrazol-5-yl)amino]-2-oxoethoxy}benzamide |
| ... | |
| ... | |

# Searching a formula in online databases

$C_{14}H_{15}ClN_4O_3$



Mass spectrum:
323.0905
324.0935
325.0879
326.0907
Mass/Charge, Da

**Randall Munroe:**
https://xkcd.com/1334/



GREETINGS, STRANGER.

WHATEVER QUEST DRIVES YOU, ABANDON IT.

YOU SHALL FIND NO ANSWERS IN THESE DESOLATE WASTES.

I KNEW I WOULDN'T.

I GUESS I ... JUST HAD TO SEE.

I HATE FEELING DESPERATE ENOUGH TO VISIT THE SECOND PAGE OF ~~GOOGLE~~ RESULTS.

*chemical*

Sorted by **Reference Count**

| | |
|---|---|
| 1 | N-(2-Acetamidoethyl)-4-[5-(chloromethyl)-1,2,4-oxadiazol-3-yl]benzamide |
| 2 | (4-Chloro-1-methyl-1H-pyrazol-5-yl)[4-(2-furoyl)-1-piperazinyl]methanone |
| 3 | 4-Chloro-N-[(1-ethyl-3-methyl-1H-pyrazol-4-yl)methyl]-3-nitrobenzamide |
| 4 | N-(3-Chloro-2-methylphenyl)-2-(3,5-dimethyl-4-nitro-1H-pyrazol-1-yl)acetamide |
| 5 | 2-Chloro-N-[(1-ethyl-5-methyl-1H-pyrazol-4-yl)methyl]-4-nitrobenzamide |
| 6 | N-(6-Amino-1-benzyl-2,4-dioxo-1,2,3,4-tetrahydro-5-pyrimidinyl)-2-chloro-N-methylacetamide |
| 7 | N-(2-Chlorobenzyl)-2-(3,5-dimethyl-4-nitro-1H-pyrazol-1-yl)acetamide |
| 8 | 5-Chloro-N-[(1-ethyl-5-methyl-1H-pyrazol-4-yl)methyl]-2-nitrobenzamide |
| 9 | 4-Chloro-N-[(1-ethyl-5-methyl-1H-pyrazol-4-yl)methyl]-3-nitrobenzamide |
| 10 | 4-Chloro-N-[(1-ethyl-3-methyl-1H-pyrazol-4-yl)methyl]-2-nitrobenzamide |
| ... | |
| ... | |
| 301 | 4-Chloro-N'-[(Z)-(3,4-dimethoxyphenyl)methylene]-1-methyl-1H-pyrazole-5-carbohydrazide |
| 302 | 2-Oxo-2-[(1,3,5-trimethyl-1H-pyrazol-4-yl)amino]ethyl 6-chloronicotinate |
| 303 | 2-Oxo-2-[(1,3,5-trimethyl-1H-pyrazol-4-yl)amino]ethyl 2-chloronicotinate |
| 304 | **Cycloxaprid** |
| 305 | N-(4-Chlorophenyl)-2-(3,5-dimethyl-4-nitro-1H-pyrazol-1-yl)propanamide |
| 306 | N-(2-Chloro-5-nitrobenzyl)-3-ethyl-N-methyl-1H-pyrazole-5-carboxamide |
| 307 | Ethyl 2-[(5-chloro-3-methyl-1-phenyl-1H-pyrazol-4-yl)carbonyl]hydrazinecarboxylate |
| 308 | N-[3-Chloro-2-(dimethylamino)phenyl]-2-(2,4-dioxo-3,4-dihydro-1(2H)-pyrimidinyl)acetamide |
| 309 | N-[3-Chloro-4-(1H-pyrazol-1-yl)phenyl]-N'-(1-hydroxy-2-propanyl)ethanediamide |
| 310 | 4-Chloro-2-{2-[(1,3-dimethyl-1H-pyrazol-5-yl)amino]-2-oxoethoxy}benzamide |
| ... | |
| ... | |

… tell me when you see something you can **pronounce**.

# The correct answer is usually the one that is easiest to read

$C_{14}H_{15}ClN_4O_3$

**Sorted by Readability**

| | |
|---|---|
| 1 | **Cycloxaprid** |
| 2 | (5s,8r)-cycloxaprid |
| 3 | avadomide hydrochloride (usan) |
| 4 | uracil, 1-{p-[3-(2-chloroethyl)ureido]benzyl}- |
| 5 | 5-chloro-n-(oxan-2-yloxy)-2-(1h-1,2,4-triazol-1-yl)benzamide |
| 6 | 8-chloro-4-(2-methyl-1,4-oxazepan-4-yl)-6-nitroquinazoline |
| 7 | n-(1-tert-butyl-1h-pyrazol-3-yl)-4-chloro-2-nitrobenzamide |
| 8 | 1-(5-chloro-2-methoxyphenyl)3-(6-ethoxypyrimidin-4-yl)urea |
| 9 | 1-(5-chloro-2-methoxyphenyl)-3-(2-(6-oxopyridazin-1(6h)-yl)ethyl)urea |
| 10 | 2-(7-chloro-5-nitro-1h-indazol-1-yl)-1-(piperidin-1-yl)ethan-1-one |
| ... | |
| ... | |

**Sorted by Reference Count**

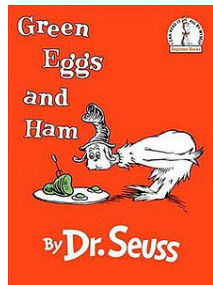| | |
|---|---|
| 1 | N-(2-Acetamidoethyl)-4-[5-(chloromethyl)-1,2,4-oxadiazol-3-yl]benzamide |
| 2 | (4-Chloro-1-methyl-1H-pyrazol-5-yl)[4-(2-furoyl)-1-piperazinyl]methanone |
| 3 | 4-Chloro-N-[(1-ethyl-3-methyl-1H-pyrazol-4-yl)methyl]-3-nitrobenzamide |
| 4 | N-(3-Chloro-2-methylphenyl)-2-(3,5-dimethyl-4-nitro-1H-pyrazol-1-yl)acetamide |
| 5 | 2-Chloro-N-[(1-ethyl-5-methyl-1H-pyrazol-4-yl)methyl]-4-nitrobenzamide |
| 6 | N-(6-Amino-1-benzyl-2,4-dioxo-1,2,3,4-tetrahydro-5-pyrimidinyl)-2-chloro-N-methylacetamide |
| 7 | N-(2-Chlorobenzyl)-2-(3,5-dimethyl-4-nitro-1H-pyrazol-1-yl)acetamide |
| 8 | 5-Chloro-N-[(1-ethyl-5-methyl-1H-pyrazol-4-yl)methyl]-2-nitrobenzamide |
| 9 | 4-Chloro-N-[(1-ethyl-5-methyl-1H-pyrazol-4-yl)methyl]-3-nitrobenzamide |
| 10 | 4-Chloro-N-[(1-ethyl-3-methyl-1H-pyrazol-4-yl)methyl]-2-nitrobenzamide |
| ... | |
| 301 | 4-Chloro-N'-[(Z)-(3,4-dimethoxyphenyl)methylene]1-methyl-1H-pyrazole-5-carbohydrazide |
| 302 | 2-Oxo-2-[(1,3,5-trimethyl-1H-pyrazol-4-yl)amino]ethyl 6chloronicotinate |
| 303 | 2-Oxo-2-[(1,3,5-trimethyl-1H-pyrazol-4-yl)amino]ethyl 2-chloronicotinate |
| 304 | **Cycloxaprid** |
| 305 | N-(4-Chlorophenyl)-2-(3,5-dimethyl-4-nitro-1H-pyrazol-1-yl)propanamide |
| 306 | N-(2-Chloro-5-nitrobenzyl)-3-ethyl-N-methyl-1H-pyrazole-5-carboxamide |
| 307 | Ethyl 2-[(5-chloro-3-methyl-1-phenyl-1H-pyrazol-4-yl)carbonyl]hydrazinecarboxylate |
| 308 | N-[3-Chloro-2-(dimethylamino)phenyl]2-(2,4-dioxo-3,4-dihydro-1(2H)-pyrimidinyl)acetamide |
| 309 | N-[3-Chloro-4-(1H-pyrazol-1-yl)phenyl]-N'-(1-hydroxy-2-propanyl)ethanediamide |
| 310 | 4-Chloro-2-{2-[(1,3-dimethyl-1H-pyrazol-5-yl)amino]-2-oxoethoxy}benzamide |
| ... | |
| ... | |

When compounds are sorted by how easy their name is to read, using a common children's book readability score, the correct answer is often in the first few rows.

While sorting by how many references a compound has will often reveal the correct answer, for newly registered pesticides, this does not work. In this example it is row 304, and it is on page 16 of a chemical compound search.
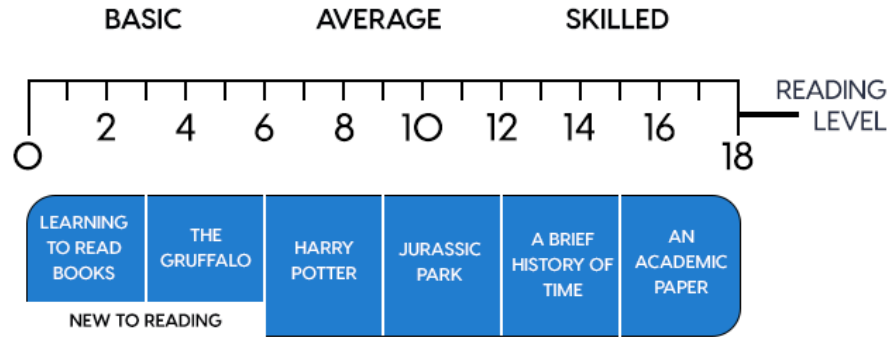
# How do we measure readability?

- Simple algorithms (some dating back to 1923) for estimating how difficult it is to read the text in a book
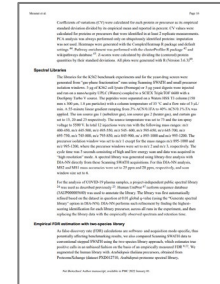


https://en.wikipedia.org/wiki/Green_Eggs_and_Ham#/media/File:Green_Eggs_and_Ham.jpg

https://readable.com/readability/flesch-reading-ease-flesch-kincaid-grade-level/

# Using readability scores on a larger data set

**Customer sample**

- 242 compounds relevant to food safety
- 53 of these not in available spectral library

**Can the compounds that are not in the spectral library be identified?**

# Automate the workflow using Python

- Pre-calculate
  - Use textstat to score readability of all compound name synonyms from PubChem database

- For a given m/z or predicted chemical formula
  - Find all possible PubChem matches using PubChemPy

- Lookup the readability of the synonyms for all matches
  - Keep anything that is readable (from infant up to my thesis supervisor)

- Search anything remaining in Google along with search terms relevant to the sample
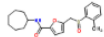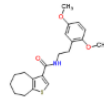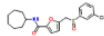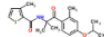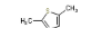  - Using Google custom search API

# Using reference counting, Isofetamid is found on page 5

## Isofetamid

On the fifth page of results, it finally shows up, buried amongst other structures with similar reference counts

Found 2513 results

Search term: **MF = 'C_{20}H_{25}NO_{3}S'**

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | ... |

| ID | Structure | Molecular Formula | Molecular Weight | # of Data Sources | # of References ▼ |
|---|---|---|---|---|---|
| 13137952<br>- 0/1 defined | | $C_{20}H_{25}NO_3S$ | 359.4824 | 9 | 10 |
| 13140800 | | $C_{20}H_{25}NO_3S$ | 359.4824 | 9 | 10 |
| 13147899<br>- 0/1 defined | | $C_{20}H_{25}NO_3S$ | 359.4824 | 9 | 10 |
| 27473807 | | $C_{20}H_{25}NO_3S$ | 359.4824 | 9 | 10 |

# Using readability, Isofetamid is successfully identified

Isofetamid

| | A | cid | name | googleHits | readability | msmsMatch1 |
|---|---|---|---|---|---|---|
| 1 | | cid | name | googleHits | readability | msmsMatch1 |
| 2 | 0 | 71657865 | isofetamid | 8350 | 12.6 | 0.980003843 |
| 3 | 1 | 58742385 | | 0 | 300 | -1 |
| 4 | 2 | 117979745 | | 0 | 300 | -1 |
| 5 | 3 | 25192457 | 1-(4-ethylphenylsulfonyl)-4-phenylazepan-4-ol | 0 | 146 | -1 |
| 6 | 4 | 25164282 | | 0 | 300 | -1 |
| 7 | 5 | 25129861 | | 0 | 300 | -1 |

Using readability on its own, or combined with Google searching or theoretical MS/MS matching, enables successful identification of the correct compound
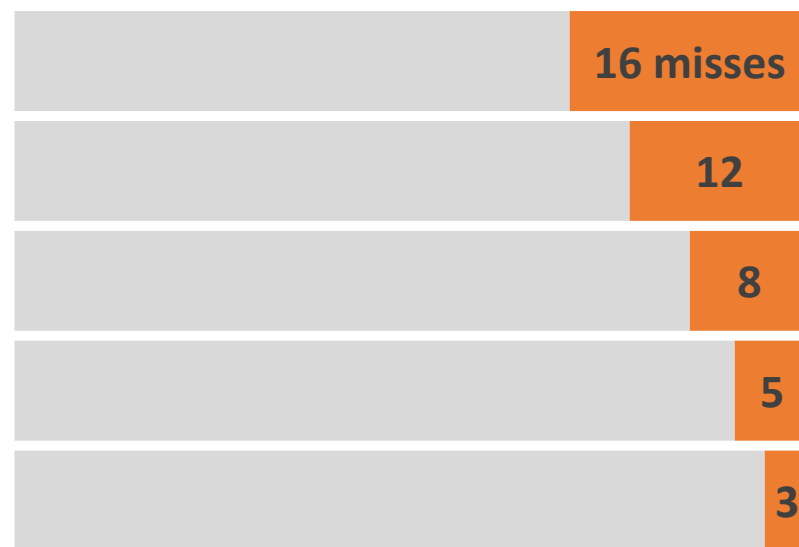
# More unknowns identified using readability

Identified using **readability** and number of **Google page hits**

| | |
|---|---|
| top 1 | **52 hits** |
| top 2 | **53** |
| top 5 | **53** |
| top 10 | **53** |
| top 25 | **53** |

Scoring the number of **Google page hits** for only those compound names that are **readable** gave the **correct** identification within the **top 2 compounds**.

Identified using **reference counts**

| | |
|---|---|
| top 1 | **16 misses** |
| top 2 | **12** |
| top 5 | **8** |
| top 10 | **5** |
| top 25 | **3** |

While using **reference counts** does work for most unknowns, this technique did **not** work for several compounds. For 3 of these examples, the correct identification was **not** in the **top 25 compounds**. Purchasing this many compounds to confirm a **single** compound would be prohibitively **expensive**.

# Trademarks / Licensing

The SCIEX clinical diagnostic portfolio is For In Vitro Diagnostic Use. Rx Only. Product(s) not available in all countries. For information on availability, please contact your local sales representative or refer to www.sciex.com/diagnostics. All other products are For Research Use Only. Not for use in Diagnostic Procedures.

Trademarks and/or registered trademarks mentioned herein, including associated logos, are the property of AB Sciex Pte. Ltd. or their respective owners in the United States and/or certain other countries (see www.sciex.com/trademarks).

© 2022 DH Tech. Dev. Pte. Ltd. RUO-MKT-11-14749-A

**Searching chemical databases** by m/z or molecular formula often generates a **very long list** of possible compounds. The **correct answer** is often the one that is **easiest to read**.
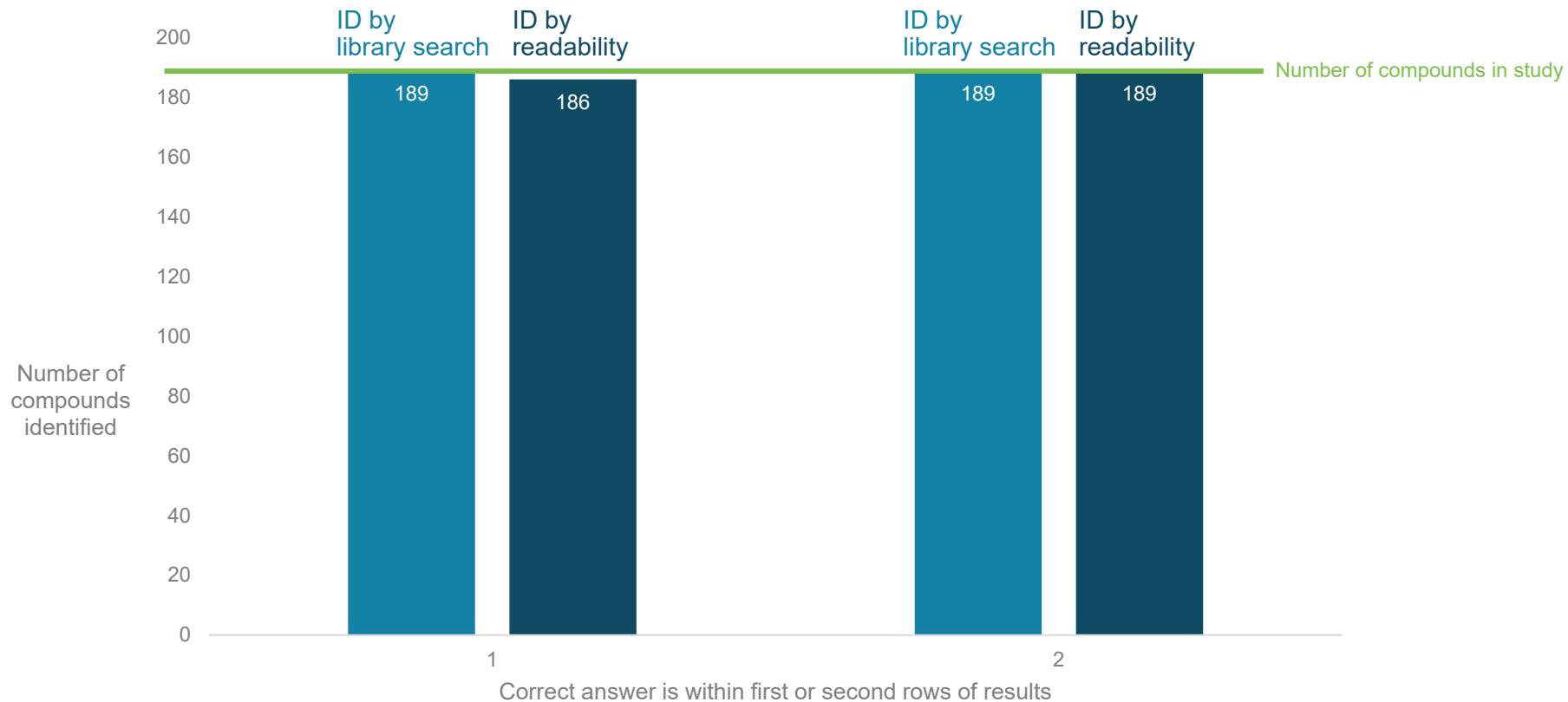
SCIEX
The Power of Precision

Questions and answers

Does **readability** have anything to do with **spectral library** matching?

# Readability scoring is almost as effective as library searching

# Trademarks / Licensing

The SCIEX clinical diagnostic portfolio is For In Vitro Diagnostic Use. Rx Only. Product(s) not available in all countries. For information on availability, please contact your local sales representative or refer to www.sciex.com/diagnostics. All other products are For Research Use Only. Not for use in Diagnostic Procedures.

Trademarks and/or registered trademarks mentioned herein, including associated logos, are the property of AB Sciex Pte. Ltd. or their respective owners in the United States and/or certain other countries (see www.sciex.com/trademarks).

© 2022 DH Tech. Dev. Pte. Ltd. RUO-MKT-11-14749-A