# HANDLING NUISANCE COMPOUNDS

## AT CZ-OPENSCREEN

**Adam Hanzlík**

UNIVERSITY OF CHEMISTRY AND TECHNOLOGY PRAGUE

CZ OPENSCREEN AT IMG

ENBIK 2025

1

# TOOLS TO HANDLE NUISANCE

## SUBSTRUCTURE FILTERS

- PAINS (Baell)
- GSK
- BMS
- LINT (Pfizer)
- and more (~2500)

## MACHINE LEARNING MODELS

- HitDexter 3
- Luciferase Advisor
- BadApple
- SCAM Detective

## BAD COMPOUND LISTS

- Aggregator Advisor
- Nuisance Compounds in Cellular Assays
- CONS (Baell)
- Obsolete Compounds

**Blacklist**

## SECONDARY SCREENS

- ALARM NMR
- Orthogonal target
- Redox assays
- Technology counter assays

**A compound that behaved badly before is likely to do so again**

- Bias for compounds set used to derive the method

- Lack of significance metric

- Redundancy when combining sources

- Lack of experimental context

- Hard to interpret

- SMARTS design sometimes not as intended

These pitfalls were carefully considered during the development of new tools to address nuisance behaviour at CZ-OPENSCREEN

- Based on a collection of known filters provided by datamol-io/medchem
- Filters merged based on matching profile against Chembl

✓ 2500 filters reduced to 1500

✓ faster

✓ number of filters matched becomes more meaningful

✓ substructure filter fingerprint as input for ML

| chembl | mol1 | mol2 | mol3 | mol4 |
|--------|------|------|------|------|
| filter1 | 0 | 0 | 1 | 0 |
| filter2 | 1 | 0 | 0 | 1 |
| filter3 | 0 | 0 | 1 | 0 |

match the same structures

| | mol1 | mol2 | mol3 | mol4 |
|--------|------|------|------|------|
| filter{1,3} | 0 | 0 | 1 | 0 |
| filter2 | 1 | 0 | 0 | 1 |

| chembl | mol1 | mol2 | mol3 | mol4 |
|--------|------|------|------|------|
| filter1 | 1 | 0 | 0 | 1 |
| filter2 | 1 | 0 | 0 | 0 |

if filter1 then filter2

- Around 4.5 million (sample X experiment X activity) measurements

- ~100K unique structures used in at least 15 experiments

- Comprehensive experiment metadata on every experiment:

    - target

    - method (Luminescence, Fluorescence intensity...)

    - assay format (cell based, biochemical)

    - reporter

- Quality Control performed on samples

Comparison of all autofluorescent-tagged against all in all primary assays. (p = 0.00e+00)

Legend:
- 202956 structure_ids in full subset, (2882094 data points)
- 711 structure_ids in full subset, (46804 data points)

## VISUAL

TWO SIDE-BY-SIDE HISTOGRAMS OF DISTRIBUTIONS

## STATISTICAL

P-VALUE OF NONPARAMETRIC TEST (H0: SAME DISTRIBUTION)

It does not matter wherever we look at structures, samples, libraries, sources, method subsets or all primary assay data.

We are always comparing two distributions of activity scores.

All readings are normalized using a modified Z-score algorithm (b-score, median instead of mean).

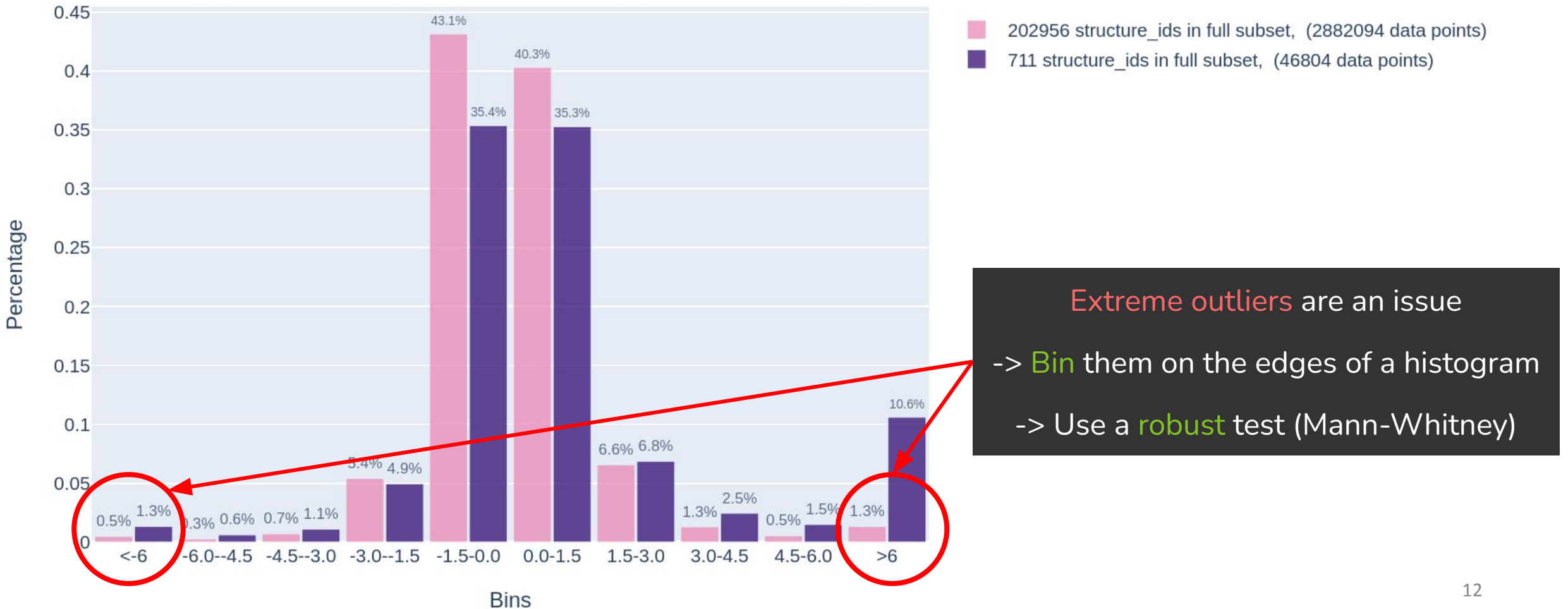Row-wise, column-wise and plate-wise median polishing

Extreme outliers are an issue
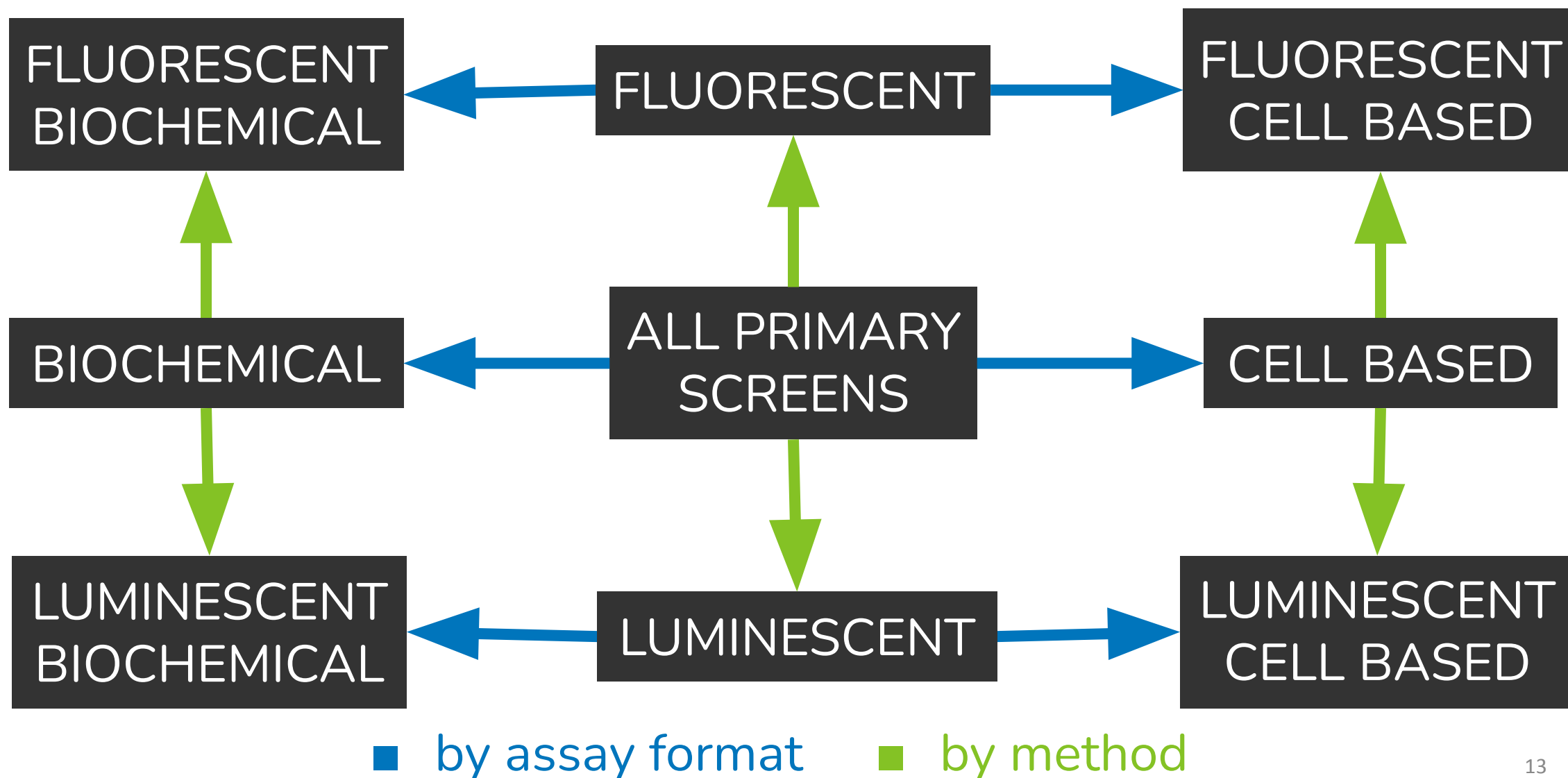
-> Bin them on the edges of a histogram

-> Use a robust test to compare the distributions

Comparison of all autofluorescent-tagged against all in all primary assays. (p = 0.00e+00)



Legend:
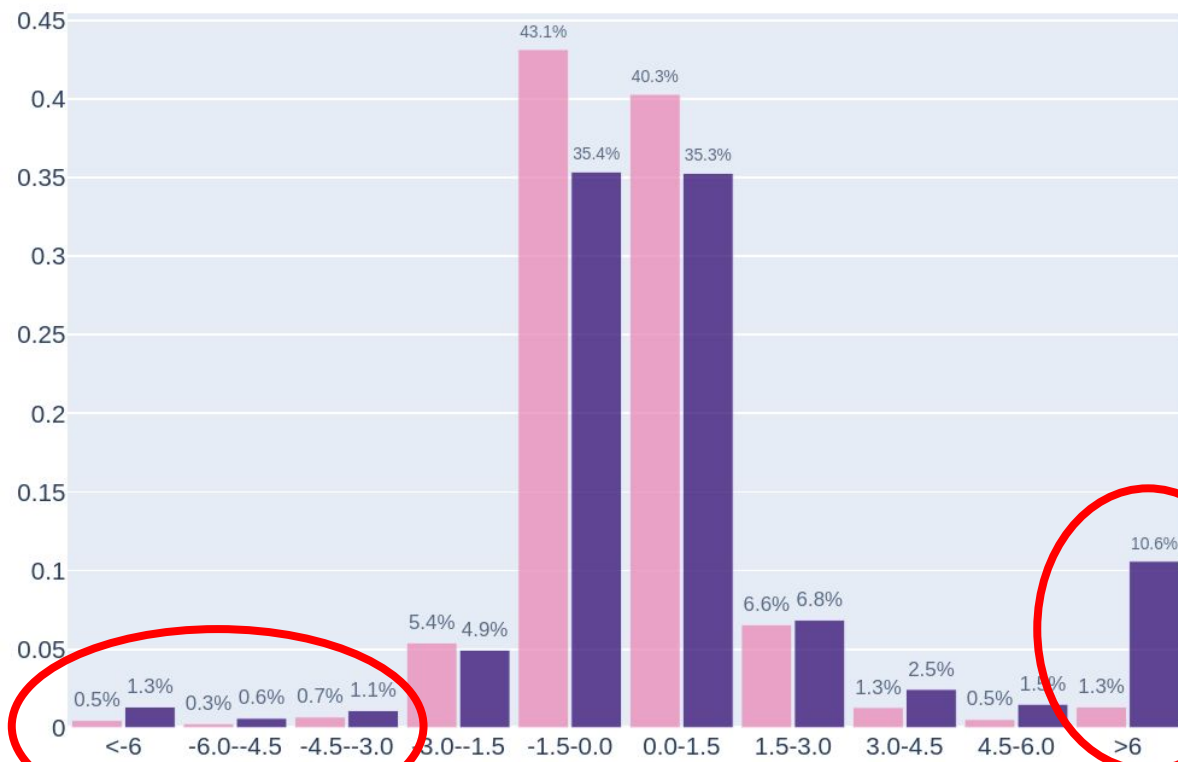- 202956 structure_ids in full subset, (2882094 data points)
- 711 structure_ids in full subset, (46804 data points)

Chart values (Percentage vs Bins):
- <-6: 0.5%, 1.3%
- -6.0--4.5: 0.3%, 0.6%
- -4.5--3.0: 0.7%, 1.1%
- -3.0--1.5: 5.4%, 4.9%
- -1.5-0.0: 43.1%, 35.4%
- 0.0-1.5: 40.3%, 35.3%
- 1.5-3.0: 6.6%, 6.8%
- 3.0-4.5: 1.3%, 2.5%
- 4.5-6.0: 0.5%, 1.5%
- >6: 1.3%, 10.6%

Extreme outliers are an issue

-> Bin them on the edges of a histogram

-> Use a robust test (Mann-Whitney)

12

✓ One method can assess ANY set of samples or structures globally or focused on a particular method and/or assay format subset

✓ Flexible and fast comparisons (slices can be preset and precalculated

✓ Visual + quantitative -> interpretation + significance

✓ Interactive - bin ranges, outlier thresholds…

UNIFIED FILTER SET

GSK  PFIZER

OTHERS  UNIFIED NON REDUNDANT SUBSTRUCTURE FILTER SET  BMS

NIH  PAINS

NUISANCE HISTOGRAM

NUISANCE FLAGS

ID 1 — AUTOFLUORESCENCE

ID 6 — CELL TOXICITY

ID 14 — FREQUENT HITTER

SUMMARY

17

# ACKNOWLEDGEMENTS

Milan Voršilák

supervisor

Ctibor Škuta

consultation

Petr Bartůněk

consultation

Tomáš Muller

consultation