



universität
wien

Statistical approach enabling technology-specific assay interference prediction from large screening data sets

04 March 2024

Vincenzo Palmacci – AIDD ESR11



About me



Problem statement

- Thanks to their sensitivity and efficiency, **fluorescence-based assays are the most widely employed technology** for the high-throughput-screening (HTS) of compounds [1, 2].
- Despite the technical advantages brought to the field, fluorescence-based assays result in a significant number of **false positive readouts caused by assay interference** [3].
- If false readouts remain undetected, they may **trigger costly follow-up studies** that may eventually turn out as futile.

Fields of application

-HITS TRIAGING:

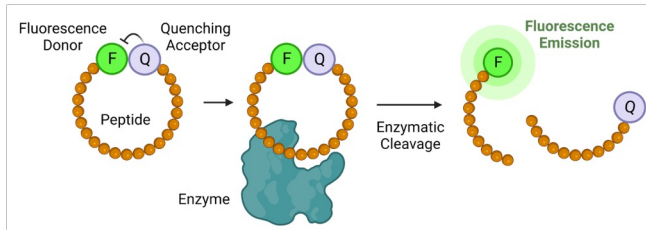
The practice of selecting a compound series with a promising efficacy profile that meets basic safety requirements and to justify investment in its optimization [4].

-NEGATIVE DESIGN:

Battery of methods that are usually employed to eliminate molecules with undesired properties [5].

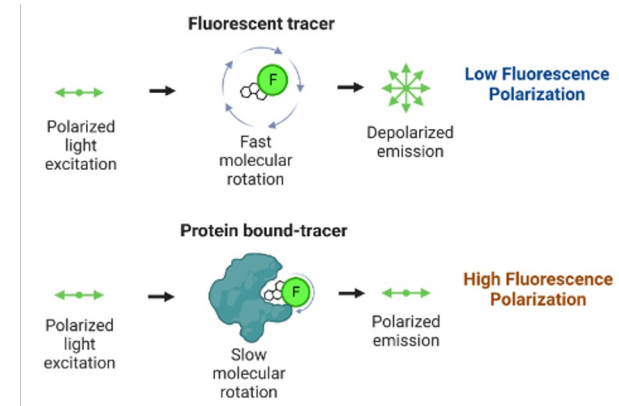
Fluorescence Intensity Assays (FLINT)

Fluorogenic assays



- Convenient for screening enzymatic inhibitors
- Fluorescent emission upon enzymatic cleavage

Fluorescence polarization (FP)



- Detect dynamic interaction between the biological target and the ligand
- Fluorescent emission upon interaction

Other popular fluorescence-based assay formats

- Fluorescence Resonance Energy Transfer (FRET)

Measures the energy transfer between a donor-acceptor pair. For the energy transfer to work donor and acceptor must be in close proximity.

- Time-Resolved FRET (TR-FRET)

Measures the time a fluorophore spends in the excited state before it reverts to its ground state by emitting a photon (FLT).

Pro of fluorescence-based assays



High specificity.



High sensitivity with low background noise.

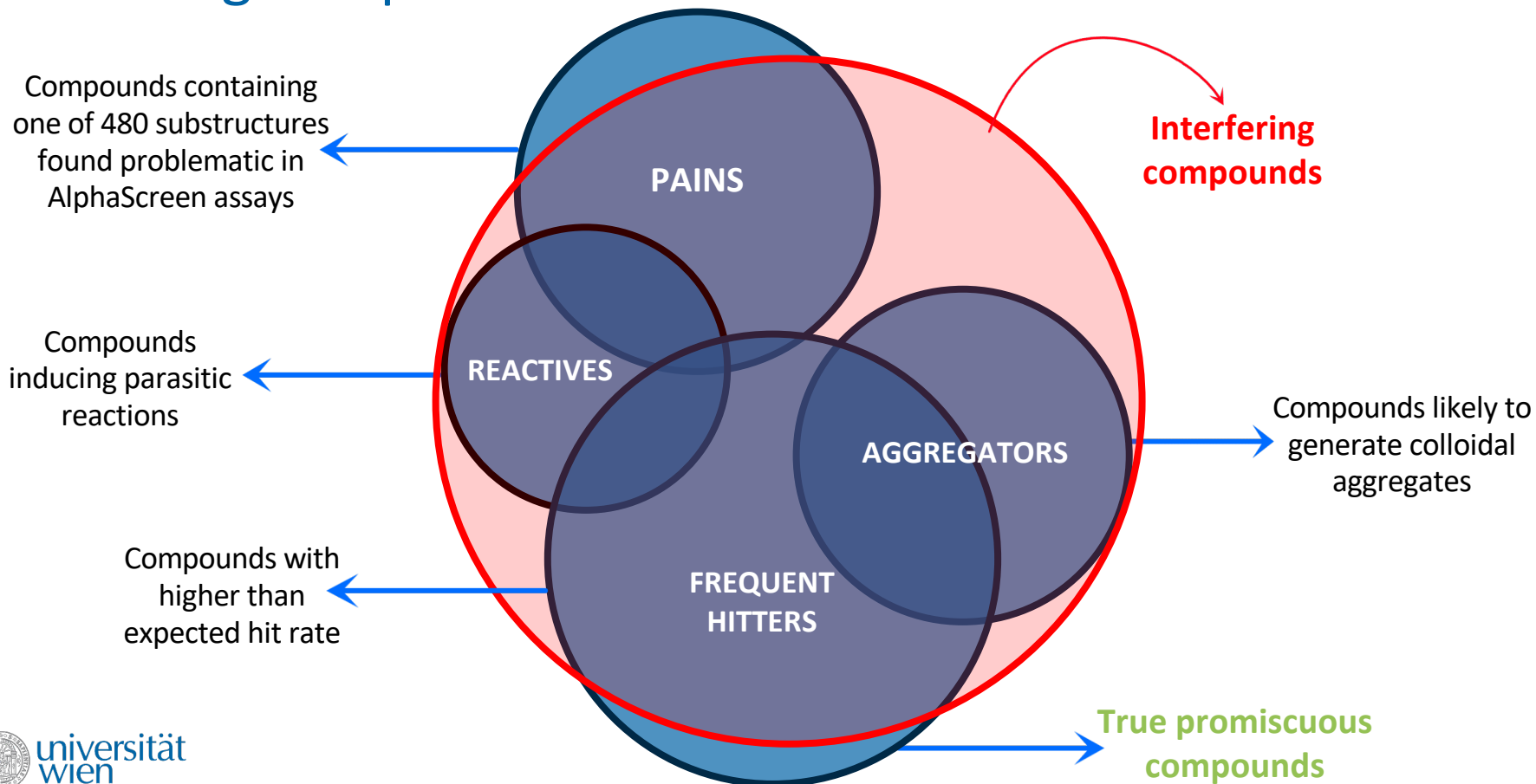


Simple operation.

No one's safe: false positive readouts

“Many hits are artefacts - their activity does not depend on a specific, drug-like interaction between molecule and protein. Artefacts have subversive reactivity that masquerades as drug-like binding and yields false signals across a variety of assays.”[6]

Interfering compounds and interference mechanisms



Dealing with assay interference: prevention measures and hit- triaging



Experimental countermeasures

- Screening with non-ionic detergents to prevent compounds aggregation
- Use of novel fluorophores emitting in a different region of the spectrum
- Use orthogonal assays to confirm the primary hits
- Implementation of counter-screen assays to identify interfering compounds

In-silico methodologies

Global methods:

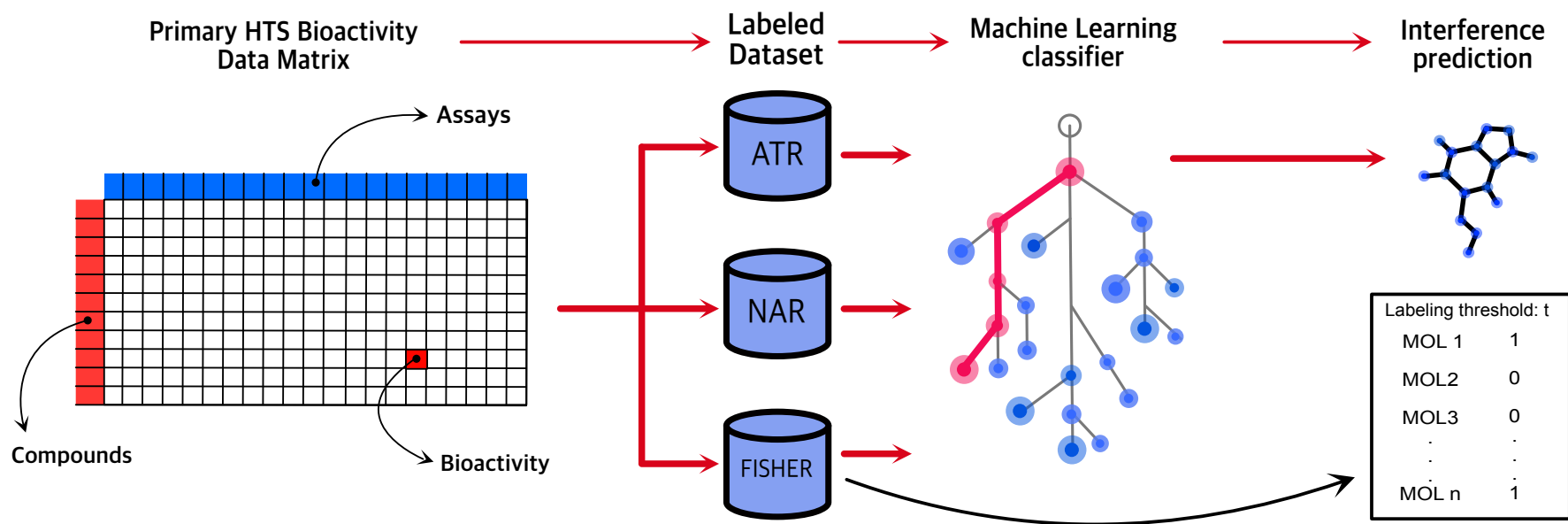
HitDexter3, Pan-Assay interference compounds (PAINS)*

Specialized methods:

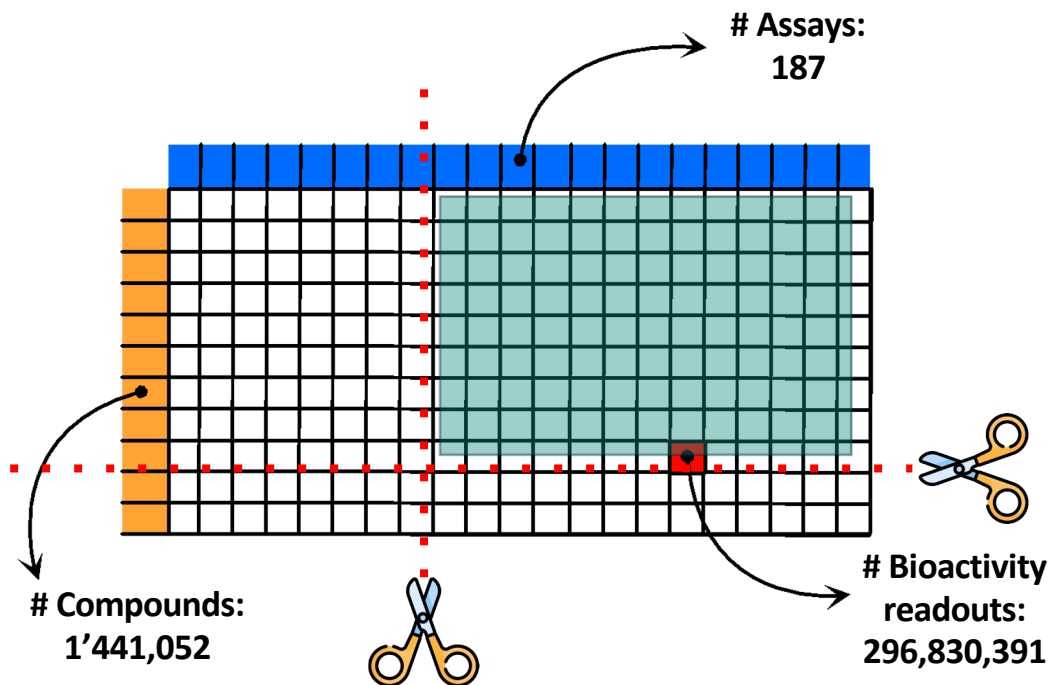
InterPred, ChemFluo, AZ (TR-)FRET interference classifiers

A new methodology to identify
compounds interfering with
specific-assay technologies

Overview



Data preprocessing: Bayer AG HTS historical data



Preprocessing pipeline

1. Assays must have bioactivity recorded for at least 80% of the compounds
2. Compounds must have bioactivity recorded for at least 80% of the assays

Dataset composition after step 1 and 2:

205 assays, 1'488'407 compounds

3. Assays must be annotated
4. Compounds must be unique (SMILES strings matching)
5. Binarize Z-scores following experimentalist indications

Dataset characterization: Bayer AG HTS historical data

Assays space:

FLINT: 56 Blue, 23 Green, 8 Red

FRET: 16

TR-FRET: 10

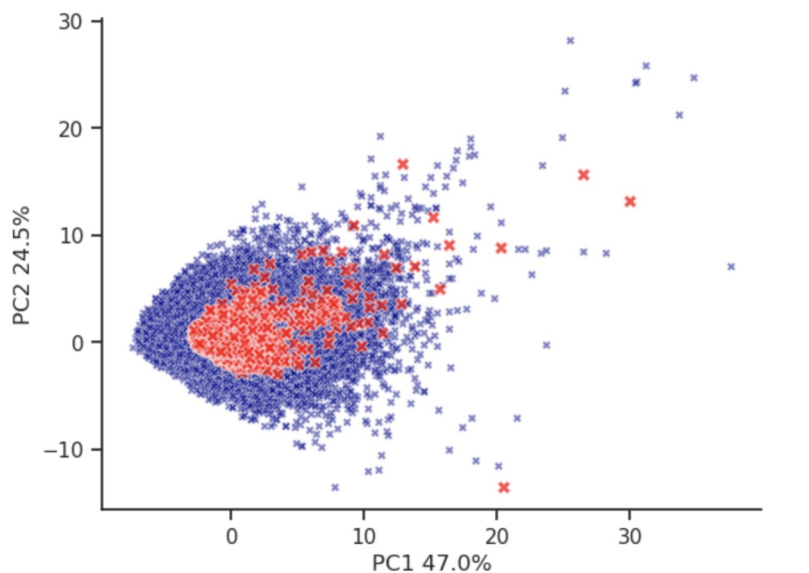
Bioluminescence: 76



Cell-based: 88

Biochemical: 99

Chemical space



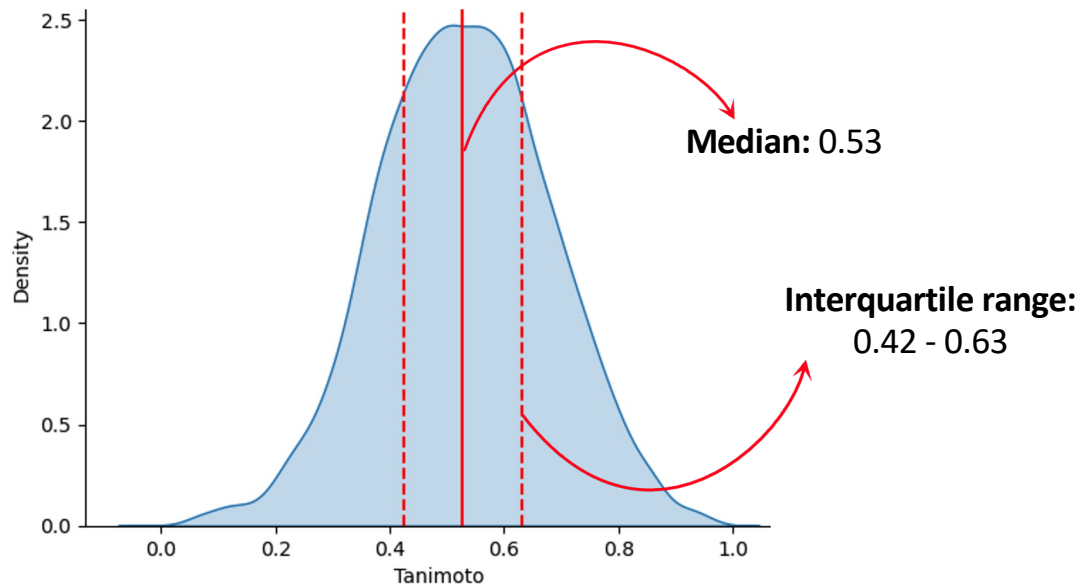
PCA comparing the training set chemical space (BLUE) and the DrugBank approved drugs space (RED)

Data collection and preprocessing: PubChem test set

BIOASSAY RECORD	
qHTS assay to test for compound auto fluorescence at 460 nm (blue) in HEK293 cells	
PubChem AID	720678
Source	
External ID	
BioAssay Type	
BIOASSAY RECORD	
Rml C and D fluorescent artifact dose-response confirmation	
PubChem AID	1696
Protein Target	dTDP-4-dehydrothymose reductase dTDP-4-dehydrothymose 3,5-epimerase
Source	PCMD
External ID	RML_FLIOR_ARTIFACT_DR
BioAssay Type	Confirmatory

10'031 unique structures

Tanimoto coefficient distribution



Distribution of compounds' maximum Tanimoto coefficient computed between the training and the PubChem test set.

Labeling compounds likely to
interfere with fluorescence-
based assays



Labelling compounds likely to interfere with the assay technology: compute interference metrics

Activity-to-tested ratio (ATR)

reloaded

Activity-to-tested ratio is computed as:

$$ATR^i = \frac{\# \text{ active readouts}}{\# \text{ times tested}}$$

Where i is the i -th compound in the dataset

1. Compute compounds ATR:
 - ATR in fluorescence assays
 - ATR in other assays
2. Apply threshold to obtain binary interference labels

Noise-to-active ratio (NAR)

Noise-to-active ratio is computed as:

$$NAR^i = \frac{\# \text{ active back} \cap \# \text{ active main}}{\# \text{ times tested}}$$

Where i is the i -th compound in the dataset

1. Compute compounds NAR considering only fluorescent assays
2. Apply threshold to obtain binary interference labels

Fisher exact test

For the compound contingency table X:

	0	1
Fluorescent assays	a	b
Other technologies	c	d

1. Compute compounds p-values applying Fisher –exact test
2. Apply threshold to obtain binary interference labels

Labelling compounds likely to interfere with the assay technology: compute binary labels

Thresholds applied to compute binary interference labels

	Percentage of likely interference compounds (thresh)			
	2%	5%	10%	20%
ATR	5.00	3.00	1.00	0.90
NAR	0.10	0.07	0.04	0.03
p-value from Fisher's exact test	0.01	0.07	0.17	0.35

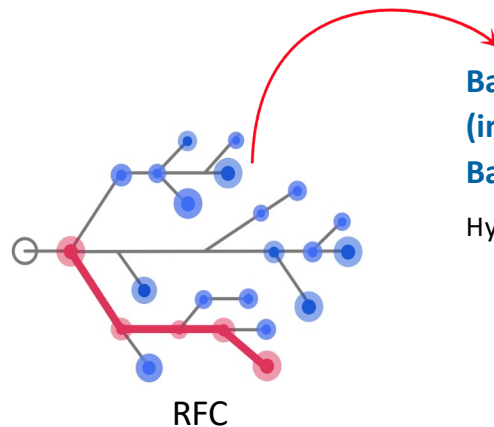
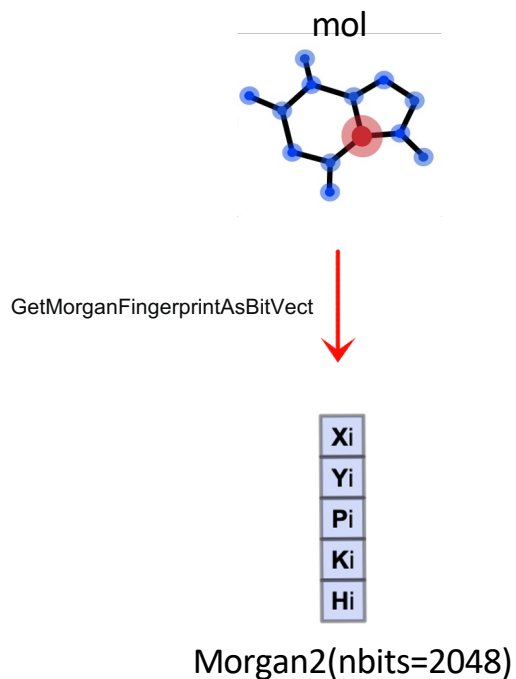
Rationale applied to interference metrics to compute binary labels

$$ATR_{Fluo}^i \geq \text{mean}(ATR_{Other}) + \text{thresh} * \text{std}(ATR_{Other})$$

$$NAR^i \geq \text{thresh}$$

$$p^i \leq \text{thresh}$$

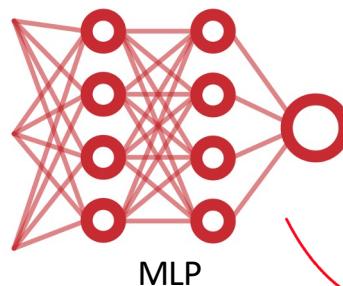
Development of machine learning classifiers for assay interference prediction



BalancedRandomForest (imbalanced-learn)
BayesOpt 50 iterations

Hyperparameters optimized:

- `n_estimators`
- `max_depth`
- `bootstrap`



MLP (PyTorchLightning)

ELU activation function
BinaryCrossEntropyLoss
WeightedRandomSampler

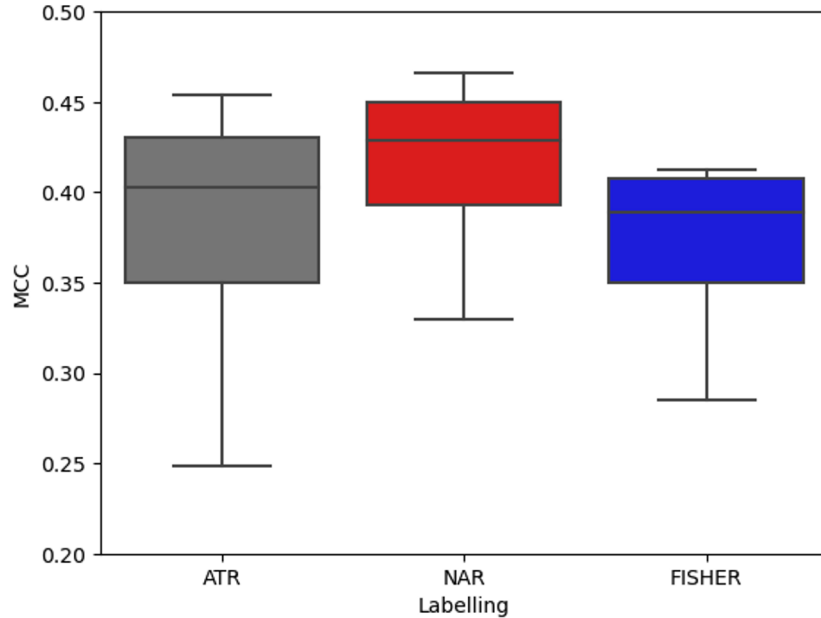
Optuna 50 iterations

Hyperparameters optimized:

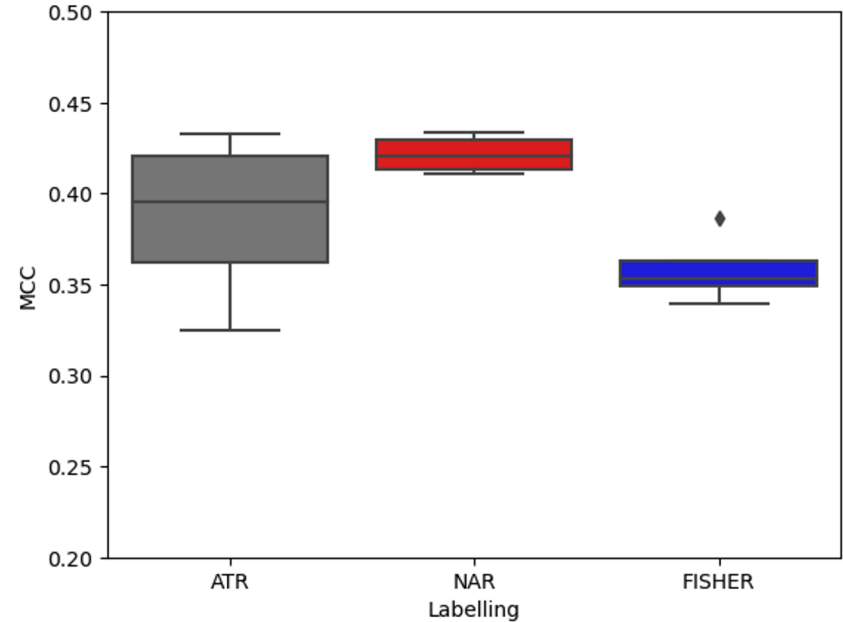
- `n_layers`
- `n_units`
- `dropout`
- `learning_rate`

Model performances on the Bayer AG test set

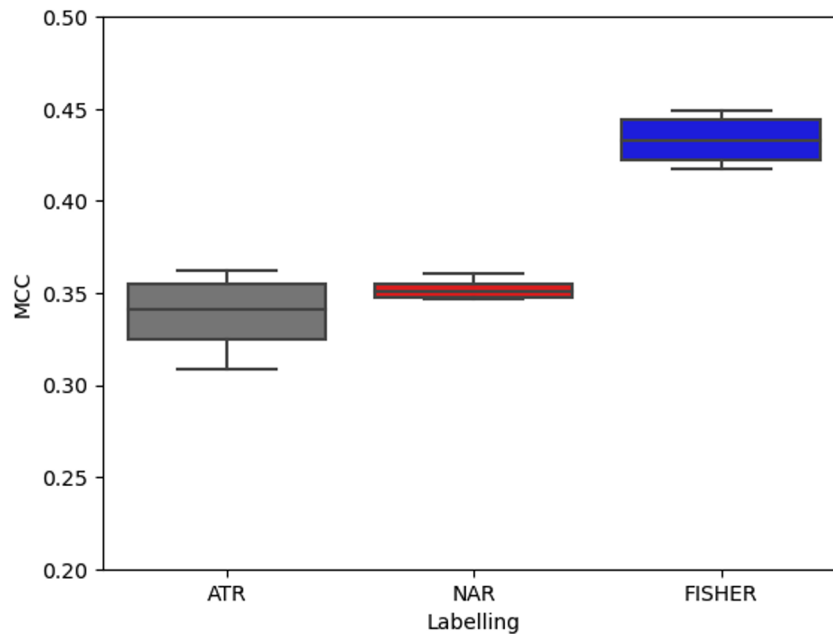
RFC MCC for different labelling methods



MLP MCC for different labelling methods

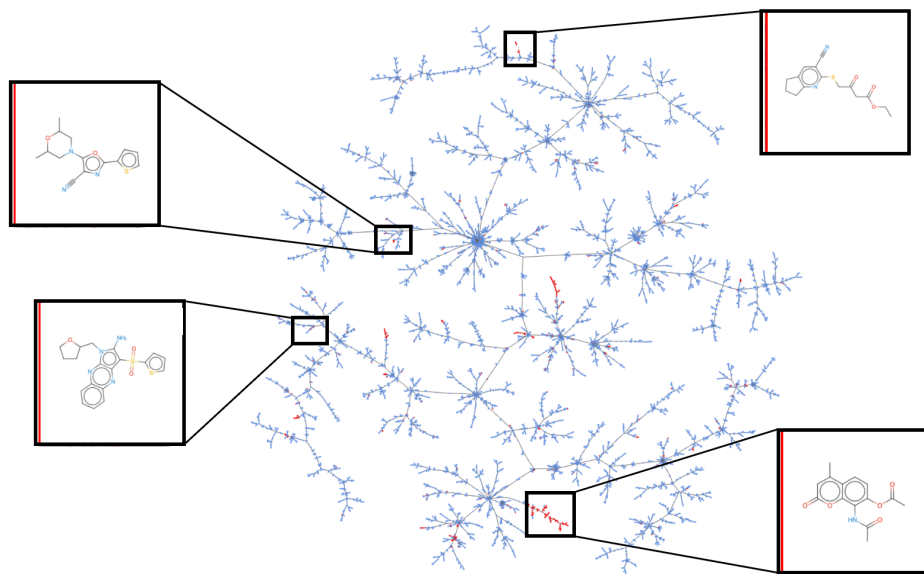


Model performances on the PubChem derived test set

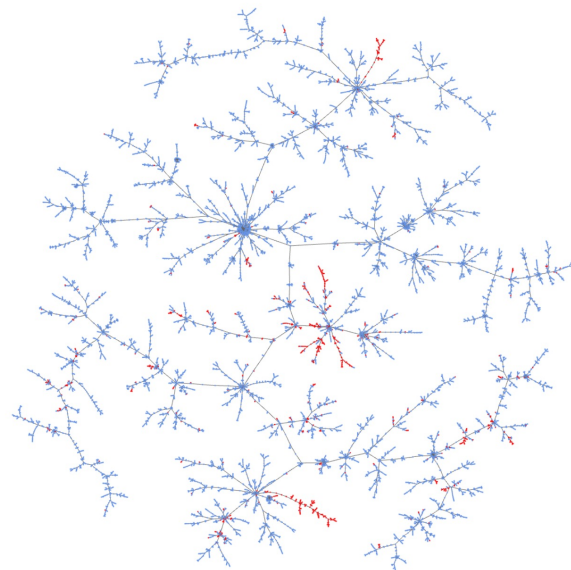


MODEL	MCC	ROC-AUC	Precision	Recall
HitDexter3.0	0.25	0.84	0.18	0.84
ChemFluo	0.34	0.82	0.23	0.75
FI-RF	0.45	0.94	0.34	0.66

Analysis of performances on the PubChem derived test set



TMAP of PubChem test set compounds colored by PubChem activity label (BLUE non-autofluorescent, RED autofluorescent)



TMAP of PubChem test set compounds colored by predicted interference label (BLUE non-interfering, RED interfering)

Conclusions

- **Single-dose HTS** data can be used with very **little preprocessing** to address assay interference
- We show that statistically derived labels can be used to train ML models for prediction of assay interference (best model reaching **MCC=0.47** on the internal test set)
- The interference labels obtained using ATR, NAR, and Fisher exact test can **approximate experimental evidence**
- Our best model **outperforms existing methods** for the prediction of autofluorescent compounds (**MCC=0.45** on the external test set)

Further experiments

- Explore if the models can predict other type of interference (e.g. aggregation)
- Extend the approach to other assay technologies