# Machine Learning for Enhanced Identification in RPLC/HRMS Non-Targeted Workflows

Hiu-Lok Ngan,[†] Viktoriia Turkina,[‡] Denice van Herwerden,[‡] Hong Yan,[†] Zongwei Cai,[*,†] and Saer Samanipour[*,‡,§]

[†]State Key Laboratory of Environmental and Biological Analysis, Department of Chemistry, Hong Kong Baptist University, Hong Kong, P. R. China
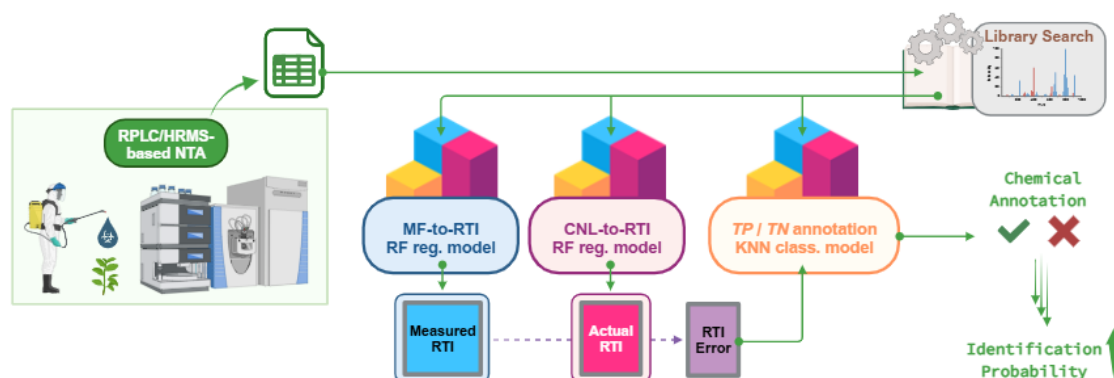
[‡]Van 't Hoff Institute for Molecular Sciences (HIMS), University of Amsterdam, Amsterdam, 1098 XH, the Netherlands

[§]UvA Data Science Center, University of Amsterdam, Amsterdam, 1012 WP, the Netherlands

**\*Corresponding authors: Drs. Zongwei Cai and Saer Samanipour**

**E-mails:** zwcai@hkbu.edu.hk; s.samanipour@uva.nl

**Abstract Graphic**

**Abstract**

16 In HRMS-based non-targeted analysis (NTA) without access to any retention information of

17 unknown compounds, spectral matching is one of the most employed approaches for the

18 assessment of chemical identification probability (IP). Recently, within the metabolomics

19 community, the use of true positive ($TP$) probability has been proposed as an alternative to the

20 conventional confidence assessment approaches. In this study, a combination of information

21 extracted from the MS/MS spectra and calibrant-free predicted retention time indices (RTIs)

22 yielded the probability of $TP$ for each chemical annotation by integrating 3 machine learning

23 (ML) models. Firstly, they include a molecular fingerprint (MF)-to-RTI model trained by 4,713

24 calibrants. Then a cumulative neutral loss (CNL)-to-RTI model was trained by 485,577

25 experimental spectra. Finally, a binary classification model was trained by 1,686,319 $TP$ and

26 true negative ($TN$) annotations. Our results demonstrated a high correlation (training: $R^2 = 0.96$;

27 testing: $R^2 = 0.88$) between MF-derived and CNL-derived RTI values, suggesting reduced RTI

28 error for $TP$ annotations. By incorporating the output parameters from a previously developed

29 library search algorithm, monoisotopic mass, and RTI error for $TP$ determination, the k-nearest

30 neighbors algorithm achieved a weighted $F1$ score of 0.65 and a Matthews correlation

31 coefficient of 0.30 for the annotations with their spectral matching scores ≥50% of total score.

32 The attained ML models were applied to RPLC/HRMS NTA of pesticide mixtures that were

spiked in solvent blank and 100× and 10× diluted black tea matrix. The chemical IPs of *TP* candidates were increased by 54.5%, 52.1%, and 46.7%, respectively. This work demonstrates the application of ML at large-scale model training to enhance chemical IP of unknown compounds.

**Keywords**

Non-targeted screening, Identification confidence, Quantitative structure-retention relationship (QSRR), MS/MS spectral reference library, Supervised learning, Model transferability

**Scientific Contribution**

Three supervised machine learning models were developed to improve the identification probabilities of unknown compounds in non-targeted RPLC/HRMS. Our work provided a higher confidence *in silico* analysis solution for the early-stage data analytics, aiding to shortlist chemicals-of-interest for further chemical identification by reference standards.

## 1. Introduction

Mass spectrometry (MS)-based non-targeted analysis (NTA) is a high-throughput technique for profiling sample analytes, distinguishing itself from targeted analysis through its discovery-oriented approach, which does not require prior knowledge of compound structures. The combination of reversed-phase liquid chromatography (RPLC) with high-resolution MS (HRMS) has become a prominent methodology for conducting NTA of the chemical exposome [1, 2], including contaminants such as pesticides [3], and per- and polyfluoroalkyl substances (PFAS) [4]. Researchers investigating human exposure to environmental chemicals not only focus on the exposome but also examine pollutant-induced alterations in the metabolome. Consequently, RPLC/HRMS-based NTA is increasingly employed for uncovering biological insights in environmental toxicology [5].

Accurate molecular annotation is crucial for compound identification. In LC/MS-based NTA, each deconvoluted two-dimensional (2D) LC/MS tensor is characterized by the information from the retention time (RT) and the mass-to-charge ratio (*m/z*) domains. To communicate confidence of compound identification, the Chemical Analysis Working Group at the Metabolomics Standards Initiative proposed a minimum metadata requirement for non-novel metabolite identification in 2007 [6]. For each annotated LC/MS tensor, its identification can only be putative if its measured *m/z* value and tandem MS (MS/MS) spectrum are just matched with its reference spectra from external libraries while the verification by chemical standard data acquired on the same platform is missed. Since 2014, similar criteria have been introduced in environmental analysis [7–9], emphasizing the necessity of further experimental effort [10], such as at least matching the retention behavior of reference standard on the same instrumental setting to achieve higher confidence match [7].

Although applying an in-house database of standard RT and MS/MS spectrum is the only efficient way to reach the highest identification confidence (IC) in RPLC/HRMS-based NTA,

- 4 -

expanding the database requires significant resources. Accurate *in silico* analysis at the early stage can aid to screen out a panel of chemicals whose annotated identities require further validation by reference standard. Practically, it involves spectral matching against the *m/z* values of molecular ions and their associated fragment ions [11, 12]. Several approaches have been introduced for determining matching tolerance, including false discovery rate (FDR) estimation [13], modified cosine, and neutral loss-based spectral alignment [14]. Additionally, novel software solutions, such as DeepMASS [15], MAW [16], and ULSA [12], have been developed to rank candidate compounds based on spectral similarity or spectral entropy similarity [17]. However, these approaches only handle information in the *m/z* domain. Considering multiple hits resulted for each MS/MS spectrum during reference library search, incorporating retention behavior into a computational process is necessary. Recently, the concept of "identification probability" (IP) was proposed by Metz et al [18], advocating for the measurement of ambiguity rather than confidence to improve transferable annotation results across analytical platforms. In our study, we aim to investigate whether incorporating chemical retention informatics with the information extracted from mass spectrum through machine learning (ML) modeling techniques enhance the compound IP in RPLC chemical space.

Model transferability is the ability of a model to accurately predict beyond the data it was originally trained on [19]. It requires diverse data sources and substantial volume during statistical modeling. Previous studies have reported on ML models for RT prediction in LC-based analyses of chemicals, such as the chemical exposome [20–24]. However, the transferability of these models is often limited by the specific chromatographic conditions [25–29], as RT can vary significantly with individual chromatographic methods. A harmonized RT index (RTI) scale is therefore necessary to offer a chromatographical method-independent alternative, although there are different schemes of RT index (RTI). It has its potentials of RTI model transferability and of integration with spectral information to enhance IC [30–33].

95    Herein, we present a computational approach to achieve calibrant-free transferable

96    prediction of RTI value on a harmonized scale and integrate it with reference spectral match

97    for compound IP enhancement in RPLC/HRMS-based NTA through statistical modeling.

98    Compounds-of-interest with high predicted IP thus are suggested to have further experimental

99    validation. Among the 3 ML models involved in our identification workflow, there are 2 models

100   built on the quantitative structure-retention relationship (QSRR) to horizontally compute

101   expected RTI values separately for the chemicals beyond the calibrants' data. Upon MS/MS

102   information being extracted through reference library search, a binary classification model that

103   incorporates information from both retention and *m/z* domains is used to determine true positive

104   (*TP*) annotations. The integrative use of our models is demonstrated by RPLC/HRMS data on

105   various pesticides-spiked blanks and black tea samples. Incorporation of the developed ML

106   models increased the IP of the chemical exposome by half on average, showing this

107   computational approach benefits from non-targeted screening in scoping chemical compounds

108   that require further identification.

109   **2.   Materials and Methods**

110   2.1.   Overall Roadmap

111   In this work, 3 ML models and the MS/MS spectral matching algorithm, ULSA [12], were

112   employed. As the roadmap illustrated in **Figs. 1A and S01**, the first model (Model 1) is a

113   random forest (RF) regression model that correlates molecular fingerprint (MF) to the true RTI

114   values of calibrants from 3 comparable scales. The second model (Model 2) is another RF

115   regression model that utilizes experimental MS/MS spectra from various external spectral

116   databases to predict the expected RTI values on a harmonized scale without knowledge about

117   the extract structure. During modeling Models 1 and 2, train/test split were performed based

118   on MF leverage rather than RTI to ensure the model was trained and tested on structurally

119   diverse compounds. The third model (Model 3) is a k-nearest neighbors (KNN) binary

- 6 -

classification model designed to predict the acceptance of matched compound annotations from ULSA or any other spectral matching algorithms with the spectral matching scores ≥50% of total score. This model integrates feature deviations from both the RTI and *m/z* domains. The computational method was validated using both internal train/test splits and external independent datasets, including RPLC/HRMS data from pesticide analyses in blank spike and tea extracts.

2.2. Model 1: Molecular Fingerprint-Based Regression Model

A previously reported RF model for predicting molecular RTI was retrained for this study with some modifications [33]. In brief, this model correlates the pre-selected MF features to calibrant RTI for transferable prediction of expected RTI value. For each canonical simplified molecular-input line-entry system (SMILES)-identified molecule, 780 AtomPairs2DFingerprintCount (APC2D) and 881 PubChem Fingerprinter (PubChem) MF features were computed using PaDEL [34]. From the features in the PubChem system, only 148 bits (i.e., from the 115[th] to the 262[nd])—specifically those describing cyclic structure—were selected and condensed as 10 features to complement with the APC2D features used to describe non-cyclic structure as previously described [33]. Co-correlation among features was examined to consider the portion of features used for each tree's training (discussed in **Supplementary Information**). As a result, each compound was represented by a 790-bit MF vector.

Different from the previous study, stratification based on the calibrants' MFs was involved in work. Transferable prediction of RTI values was achieved by training a model on 4,713 calibrants across 3 comparable scales of RTI in RPLC: the C3–14 n-alkylamide system (Amide) [35], the RTI system developed by Aalizadeh et al. from the University of Athens (UoA) [31], and the C0–23 cocamide diethanolamine homologous series (Cocamide) [36]. Details on model training and validation are discussed in **Supplementary Information**. To assess model transferability, accuracy of the predicted RTI value on the harmonized scale was tested by the

https://doi.org/10.26434/chemrxiv-2024-mdl4q-v2 ORCID: https://orcid.org/0000-0002-6002-9660 Content not peer-reviewed by ChemRxiv. **License:** CC BY-NC-ND 4.0

true values of 1,237 untrained calibrants in terms of mean absolute error. The list of compounds used to train and test this MF-to-RTI RF regression model was available online (see file 1 on [10.5281/zenodo.14752891](10.5281/zenodo.14752891)).

2.3.   Model 2: Cumulative Neutral Loss-Based Regression Model

In addition to predicting expected RTI values on a harmonized scale from 2D chemical descriptors, a second RF regression model was established to correlate empirical positive electrospray ionization (ESI+) MS/MS spectral data to the RTI values predicted from Model 1, following a strategy similar to a prior study [32]. Feature for this Model 2 was based on cumulative neutral loss (CNL), a method recently applied for matching analog molecules [14, 37, 38]. This approach has proven advantageous for chemical componentization [14, 39]. Based on the previous results [32, 40], a list of 15,961 high probability masses-of-interest (MOIs) was pre-selected as features. The list is available online (see file 2 on [10.5281/zenodo.14752891](10.5281/zenodo.14752891)). Moreover, monoisotopic mass was included as an additional feature, as it was previously reported to be one of the most discriminative to CNL-based model due to its intrinsic relationship with molecular weight [32].

Model 2's training dataset was compiled from various MS/MS reference spectral data from known compounds that have chemical identifiers, specifically InChIKeys and SMILES IDs. Databases including MassBank EU, MoNA, and NIST LC-MS/MS provided multiple spectra for specific compounds. In total, there were 27,211 distinct molecules, consisting of 693,685 MS/MS spectra used for training (the list is available online, see file 3 on [10.5281/zenodo.14752891](10.5281/zenodo.14752891)). Each spectrum contains at least 2 MOIs. To address the occurrence of multiple compounds, Hao Guo et al. implemented a random train/test split strategy, reserving only 1 spectrum for training and testing while using others for validation [41]. This approach effectively manages sample balance, ensuring each compound carries equal weight in the model, which is particularly important for regression modeling. However, the output labels of

https://doi.org/10.26434/chemrxiv-2024-mdl4q-v2 **ORCID:** https://orcid.org/0000-0002-6002-9660 Content not peer-reviewed by ChemRxiv. **License:** CC BY-NC-ND 4.0

Model 2, the predicted expected RTI values generated by Model 1, were primarily localized in the center rather than scattered in the margin regions (around RTI values of 0 and 1,500), making edge case overloading neglectable. In our study, we aimed to include as many molecules as currently available in the relevant spectral databases. This CNL-to-RTI RF regression model was trained on 485,577 query MS/MS spectra. Details on model training and validation are discussed in **Supplementary Information**. To evaluate the exchangeability between Model 1 and Model 2, exploratory data analysis (EDA) and quartile analysis were performed to study the similarity in distribution and the closeness of the predicted RTI values with the true RTI values of the compound calibrants used to train Model 1. Predictive power of Model 2 was tested on 208,104 query spectra by using $R^2$ value and root mean square error (RMSE, calculated according to **eq1**) as the performance metrics.

2.4.   Spectral Reference Library and Universal Library Search Algorithm

ULSA, a previously developed MS/MS spectral matching algorithm [12], was employed in this study to annotate compounds by matching MS/MS spectra from various public spectral databases, including MassBank EU, MoNA, and GNPS. During the execution of ULSA, 7 different parameters were summed to derive a final performance score using the formula outlined in **eq2**. The final scores ranged from 0 to 7, with higher values indicating increased reliability of the assigned annotation based on spectral matching. "FinalScoreRatio" ($S_8$) was then calculated by dividing the final score by 7. In **eq2**, $MS_{tol}$ represents the user-defined $m/z$ tolerance in Daltons for an ion to be accepted as a positive match. The 7 features ($S_{1-7}$) that were extracted from the spectra for feature selection in Model 3 included: $S_{1-2}$, the fractions of matched fragments to the reference spectra (denoted as $f_{matched}/f_{reference}$ in **eq2**) and to the total fragments in the user spectrum ($f_{matched}/f_{user}$); $S_3$, the observed mass error in the precursor ion ($MS_{err,pre}$); $S_4$, the absolute average mass error of all matched fragment ions ($|MS_{err,frag}|$); $S_5$, the standard deviation of mass error for matched fragments ($STD_{err,frag}$);

195 and $S_{6-7}$, the forward ($MF_{forward}$) and reverse ($MF_{reverse}$) match factors calculated via dot

196 product. Importantly, these output parameters are from statistical calculations that can be

197 computed directly from the spectra, making it rather independent from ULSA implementation.

198 Also, the publicly available spectral libraries are separated to ULSA. For any reference MS/MS

199 spectral data that is newly acquired, uploading it to these databases allowing them to be

200 matched. That is particularly valuable for annotating the chemicals of emerging concern.

201 2.5.  Model 3: Predicting Acceptance of Molecular Annotation

202   Model 3 is a binary classification model aimed at determining whether a molecular

203 annotation suggested by ULSA was *TP* and should be accepted (the outputs were labeled as

204 "1"), or it was indicates true negative (*TN*) and should be rejected (labeled as "0"). The features

205 for this model include the RTI difference between the RTI values derived from Models 1

206 ($RTI_{MF}$) and 2 ($RTI_{CNL}$), along with monoisotopic mass and 8 candidate features obtained from

207 ULSA. We hypothesize that a larger RTI error indicates *TN* annotations, while a smaller error

208 correlates with *TP* annotations. Some features as discussed in **Supplementary Information**

209 exhibit right-skewed distribution were log-transformed, yet all data were scaled by mean

210 normalization for model selection and optimization. Several ML algorithms, including logistic

211 regression (LR), decision tree (DT), RF, and KNN, were evaluated. Pearson correlation was

212 applied to eliminate redundant features that have less feature importance and *r* values >0.80

213 for fair comparison among models. Further feature selection focuses on the model's sensitivity

214 to excluded features, prioritizing those with the most negative contribution to predictions.

215 Details information on model selection and models' hyper-parameter tuning are discussed in

216 **Supplementary Information**. A set of semi-synthetic MS/MS spectra was prepared following

217 methodologies from prior research [12]. In short, HRMS spectra with a mass resolution >5,000,

218 acquired using an ESI+ source, were randomly selected from databases such as MassBank EU,

219 MoNA, GNPS, and NIST. To mimic false negative (*FN*) ions (i.e., loss ions) caused by ion

- 10 -

suppression or acquisition defects, 10% fragments were removed randomly. Data augmentation was applied by adding noise to simulate false positive (*FP*) ions. For instance, if the mass tolerance for MS/MS spectral matching was set to $\pm 20$ mDa, *FP* ions with random mass errors in a range of $\pm 30$ mDa were incorporated to the accurate masses of each fragment ion. These semi-synthetic sample spectra were then analyzed by ULSA, yielding 4,368,902 annotated MS/MS spectra. Since all annotated spectra with a "FinalScoreRatio" ($S_8$) <0.50 are in 100% of *TNs*, ML was deemed unnecessary for distinguishing *TP* from *TN*. Ultimately, 1,686,319 spectra were used for training and 421,381 for testing. Among the training samples, 1,535,009 spectra were *TN* (label "0"), while 151,310 instances represented *TP* (label "1"). Nine additional replicates of *TP* spectra were used to balance the training dataset. The optimal binary classification model was tested on an external independent dataset, which comprised of experimental data of pesticides spiked into blank samples at varying concentrations (1, 2.5, 5, 10, 25, 50, 100, and 1,000 ppb). Because both *TP* and *TN* are known in the blank spikes, the optimal model was primarily assessed based on the Matthews correlation coefficient (MCC) score (**eq3**) and secondarily by weighted *F*1 score (**eq4**).

2.6. Application of the Models in the Identification Workflow: NTA of Pesticides-Spiked Black Tea Samples

After developing all the predictive models, they were integrated into the NTA pipeline for aiding pesticide annotation as a real-world application. Their overall effectiveness in enhancing compound IPs was evaluated using a set of RPLC-MS/MS data from pesticides-spiked black tea extract samples, as previously reported in a CNL model development study [40]. The pesticide standards were part of the LC/MS Pesticide Comprehensive Mix Kit (Part Number 5190-0551) from Agilent Technologies. Pesticide mixtures were spiked in solvent blank and 10× and 100× diluted black tea matrix respectively, resulting 1 ppm final solutions. The details of data acquisition were outlined in the prior CNL model study [40].

- 11 -

245    - Data Pre-Processing

246    To extract HRMS information required as the model input, a user interface, namely "jHRMS

247    ToolBox", was previously developed for no-code application [42]. Detailed descriptions of its

248    component modules and their modules' input/output were summarized in **Table S1** for

249    facilitating MS community to reassemble the pipeline. To introduce, firstly, the upstream

250    processing of LC/HRMS files requiring the open MS format, ".mzXML". Raw data files

251    conversion was accomplished using the "MS_Import" module, built on MSConvertGUI (64-

252    bit, ProteoWizard [43]). Secondly, this ecosystem also incorporates the self-adjusting feature

253    detection (SAFD) algorithm for extracting LC-MS/MS features [44]. Thirdly,

254    componentization was achieved by executing the "CompCreate" module, which correlates the

255    extracted features with the probabilistic CNL model [40]. Finally, the resultant MS/MS spectra

256    were then analyzed using ULSA for compound annotation [12]. It returned a spreadsheet that

257    contains necessary information for the identification workflow, including the suggested

258    compound, $m/z$ value of the molecular ion, a list of $m/z$ values of the correlated fragments, and

259    the scores $S_{1-8}$.

260    - Models Application

261    In **Figs. 1B** and **S01**, the resultant spreadsheet was fed into Models 1 and 2 to calculate RTI

262    error (as defined by **eq5**). Subsequently, Model 3 was deployed for each annotated

263    RPLC/HRMS tensor, utilizing the necessary ULSA-returning scores, monoisotopic mass, and

264    RTI error to predict whether the suggested chemical annotation was acceptable.

265    Model 1 (the MF-to-RTI model) was deployed for obtaining the annotated compounds' RTI

266    values (considered as measured values). Leverage ($h_{ii}$) was used to assess whether an

267    annotated compound was in the applicability domain (AD, the 95% leverage threshold) of

268    Model 1 (0.275) based on the training dataset. It is calculated by the equation shown in **eq6**,

269    where $X$ represents the matrix of training data, and $x_i$ is the 790-bit vector for an individual

- 12 -

270    data query. If the leverage value of a suggested annotation is within AD, the obtained RTI value

271    has low uncertainty and can be treated as a ULSA-measured value. Similar, Model 2 (the CNL-

272    to-RTI model) was deployed for obtaining the actual RTI values if and only if a RPLC/HRMS

273    tensor's spectral data was within the AD of Model 2 ($h_{ii}$ <0.146). Upon computing the RTI

274    values from both Models 1 and 2, RTI error, the difference between measured and actual RTI

275    values, was then calculated and proceeded to Model 3 as one of the features, namely "DeltaRi".

276    Model 3 was deployed for *TP* annotation determination. A script for its application is

277    available (see Section **Computations and Code Availability**). In short, the RPLC/HRMS

278    tensors with the values of $S_8$ ≥0.50 and within the ADs of Models 1 and 2 were considered.

279    Independent variables included logarithm of $S_1$, $S_3$, $S_5$, logarithm of $S_8$, logarithm of a value

280    of monoisotopic mass/1000, and "DeltaRi". They were normalized by mean before Model 3

281    implementation. For real sample analyses, recall (as defined by **eq7**) was chosen as the primary

282    performance metric since *TN* and *FP* are always unknown in real sample.

283    -    Identification Probability Calculation

284    In ambiguity measurement of a candidate RPLC/HRMS tensor, a collection of reference

285    spectra can be matched. For each compound hit, multiple spectra of the same compound may

286    correspond. Hence, a weighted probability of *TP* (*P*(*TP*)) was calculated for each hit. A decision

287    threshold of 0.50 was then applied to determine whether one should retain or exclude that hit.

288    The shortlisted hits for each RPLC/HRMS tensor were then used to calculate the IP according

289    to **eq8**. In other words, if ULSA implementation yielded 5 hits, the probability of *TP* from

290    ULSA determination (*P*(ULSA)) was calculated as 1/5. Since different numbers of reference

291    spectra (1, 5, 10, 20, and 50) were matched for the respective hits, the weighted *P*(*TP*) for each

292    hit was calculated as the summation of *P*(*TP*) values for each spectrum-specific annotation

293    derived from Model 3, followed by a division of the number of spectra. Thus, the calculations

294    would be Σ(*P*(*TP of the 1ˢᵗ hit*))/1, Σ(*P*(*TP of the 2ⁿᵈ hit*))/5, Σ(*P*(*TP of the 3ʳᵈ hit*))/10, Σ(*P*(*TP*

- 13 -

295    *of the 4$^{th}$ hit*))/20, and Σ(*P*(*TP of the 5$^{th}$ hit*))/50. If only 2 hits remained from the 5, the IP

296    would then be 1/2.

297    2.7.   Computations and Code Availability

298         All scripts for ML modeling and data visualization were written in Julia v.1.6 using Visual

299    Studio Code (Microsoft), while some data visualization was conducted in Python v.3.10 on

300    Jupyter Notebook (Anaconda3). The scripts (with step-by-step instructions) for models

301    application and IP determination and also the developed models are available in the following

302    repositories respectively: git@github.com:TommyNHL/exposomeIDproba.git (scripts); and

303    git@bitbucket.com:hiulokngan/modelsExposomeProjUvA.git (models). Details regarding

304    computational power and packages used can be found in **Table S2**. Any updates will be

305    published on both repositories with a citation of the final DOI for this work.

306    **3.   Results and Discussion**

307    3.1.   Performance of Model 1 (the Molecular Fingerprint-to-Retention Time Index Model)

308         The first model developed in this study was utilized to predict $RTI_{MF}$ (considered as

309    measured RTI for each suggested annotation) on a harmonized scale. It was trained by 4,713

310    compound calibrants with 5,048 true TRI values across 3 similar RTI scale systems. The mean

311    absolute error of the predicted values of the trained calibrant was 89.52. Contrasting to the

312    mean absolute difference (78.56) among true RTI values from different scales of individual

313    calibrants with publicly available MS/MS spectra (the list is available online, see file 4 on

314    10.5281/zenodo.14752891), similar uncertainty might exist for the predicted values from

315    Model 1. The uncertainty (with a mean absolute error of 111.15) was slightly increased when

316    the model tested by 1,263 true values of 1,237 calibrants. In the histograms depicted in **Figs.**

317    **2A and 2B**, the true RTI values (colored in purple-blue) and $CNL_{MF}$ (in sky-blue) exhibited

318    similar distributions among the data used for training and beyond modeling respectively. Hence,

319    the retrained MF-to-RTI model has been validated to be transferable in this study.

- 14 -

3.2.  Model 2 (the Cumulative Neutral Loss-to-Retention Time Index Model)

The second model established in this study was employed to predict $RTI_{CNL}$ (considered as actual RTI for each suggested annotation) on a harmonized scale that was defined in Model 1. Model 2 was trained to correlate the pre-selected CNL masses with $RTI_{MF}$ by 485,577 query MS/MS spectra and tested by 208,104 spectra. Both Models 1 and 2 were used to compute RTI error of each chemical annotation upon universal library search. It was calculated by the difference between $RTI_{MF}$ and $RTI_{CNL}$. If a chemical annotation is *TP*, both predicted RTI values should be similar and thus the resultant error is close to 0. Wrong annotation gives a large different in $RTI_{MF}$ (measured value) from $RTI_{CNL}$ (actual value), resulting larger RTI error in a *TN* annotation.

- Model Performance

Due to QSRR, distribution of $RTI_{CNL}$ (green bars) were similar to $RTI_{MF}$ (sky-blue bars) and the true RTI values (purple-blue bars) of chemical calibrants used to train (**Fig. 2A**) and test (**Fig. 2B**) Model 1. Variances of RTI error from Model 1 (sky-blue box) was smaller than that from Model 2 (green box, **Fig. 2C**) for the data it was originally trained on Model 1. However, both models showed similar RTI error variance for the extended data (data not used to train Model 1, see **Fig. 2D**). These similar deviation patterns demonstrated the exchangeability of RTI values that were computed from Models 1 and 2.

For the model trained with a 7:3 train/test split, the RMSEs were 51.8 and 92.8, with $R^2$ values of 0.96 and 0.88 for the training and testing datasets (**Table 1** and **Figs. 2E and 2F**). These results were comparable to the predictive power of the CatBoost model that was trained on NORMAN dataset by different descriptors by Boelrijk et al. [32]. Our retrained RF model demonstrated the predicted values had better correlation with MF-derived RTIs (train: $R^2 =$ 0.94; test: 0.85 in Boelrijk et al.'s work). Larger RMSEs were observed in our study (train: 44.0; test: 67.0 in Boelrijk et al.'s work) might be due to the RTI scale adopted in this study

- 15 -

was across 3 comparable RTI scaling systems. The comparable predictive performance with the previous work confirming that this newly constructed Model 2 was transferable and could effectively predict expected RTI values based on QSRR.

- Feature Interpretation

In this study, the spectral data collected showed mass distribution of the CNL masses during ESI+ MS/MS analysis lying into range between 0–500 Da (see **Fig. S02**), consistent with the findings from van Herwerden et al. [40]. CNL masses with lower weight showed higher frequency than the higher masses (the list of counts for each CNL mass is available online, see file 2 on [10.5281/zenodo.14752891](10.5281/zenodo.14752891)). Similar distribution also can be seemed in the semi-synthetic data (see **Fig. S03**).

3.3. Model 3: Binary Classification Model

- Model Performance

Model 3 was a KNN model that incorporated 6 features, namely "RefMatchFragRatio", "MS1Error", "MS2ErrorStd", "FinalScoreRatio", monoisotopic mass, and RTI error. This model only applied to the individual spectral matching results with the scores of ≥50% for *TP*/*TN* annotation determination. Our optimized model was a RTI error inclusive model. The area under the receiver operating characteristics curve (AUROC) and the area under the precision-recall curve (AUPRC) were 0.95 and 0.93 in training and testing datasets, respectively (**Fig. 3A**). RTI error inclusive model exhibited slightly greater AUROC and AUPRC values than the values of RTI error exclusive model (train: AUROC = 0.94, AUPRC = 0.92; test: AUROC = 0.94, AUPRC = 0.91). If only considered positive results, the recalls were 0.99 for training and testing datasets (see **Table 1**). Our KNN model achieved a weighted *F*1 score of 0.65 and a MCC score of 0.30 for pesticides-spiked blank samples, with a recall of 0.66 (**Table 1**). Notably, the predictive power of RTI error inclusive model resulted in better recall (0.60) for determining acceptance of the chemical exposome annotation when the

samples presented with 10× and 100× diluted black tea matrix than the RTI error exclusive model (recall = 0.54). We demonstrated that ML modeling by adopting semi-synthetic data was feasible to achieve acceptable model transferability, annotating RPLC/HRMS tensors in higher IC.

By deploying Model 3, the default probability threshold for accepting an annotation as *TP* was set at 0.50, resulting in *TP* and *FP* rates of 98.8% and 23.4%, respectively (**Fig. 3C**), with a FDR of 19.1%. FDR-controlled cut-offs for the probability of accepting an annotation as *TP* were 0.89 and 0.78, respectively (**Fig. 3D**). These decision thresholds could be applied to calculate recalls for individual samples in different levels of IC. For 1 ppm samples, their recalls increased as the matrix effect decreased (**Fig. 3E**). When the matrix effect was absent, as solution concentrations increased, so did the number of deconvoluted RPLC/HRMS tensor and the number of positive candidates to be identified by our identification workflow. However, Model 3 exhibited no predictive power for the analyses of 1 ppb solution, while it functioned stably for the analyses of chemical contaminants with concentrations of 2.5 ppb or above.

- Feature Interpretation

Since the developed binary classification model indicated discriminative power to *TP*/*TN* chemical annotations, we conducted permutation analysis to evaluate feature importance. All KNN models with MCC scores >0.30 on validation dataset (the experimental pesticides-spiked blanks) during model selection and hyperparameter tuning were examined. Our results confirmed that the excluded features, logarithm of $S_2$ ("UsrMatchFragRatio") and the difference between $S_6$ ("DirectMatch") and $S_7$ ("ReverseMatch"), were one of the factors contributing to the models' predictive power on experimental data. The most significant feature was $S_3$ , "MS1Error" (**Fig. 3F**), while the remaining features indicating importances occasionally with positive and negative contributions that were close to 0.00. These findings were only specific to the binary classification models we attained, while they did not reflect

the intrinsic predictive values of features. Hence, Model 3 was significantly sensitive to the differences between the measured molecular ion mass and its actual *m/z* values reported from public reference spectra. Incorporation of other least discriminative features was still critical for Model 3 to exhibit its predictive power on experimental data.

3.4.    Performance Comparison with Solely MS/MS Spectral Library Search

To examine the extent of improvement achieved by deploying our models (Models 1, 2, and 3) accompanying with MS/MS spectral reference library search, we calculated the IP of molecular features based on spectral matching and ML-assisted determination of annotation acceptance or rejection. Performance comparison was based on 2 different definitions of "hit".

The first definition of hit refers to a collection of reference spectra matched for measuring the ambiguity of a candidate compound. We averaged the IPs of all positive results to facilitate overall performance comparison. As shown in **Fig. 3G**, for 1 ppm sample solutions, integration of our 6-component KNN model (Model 3) enhanced average IPs from 44% to 68% (24% difference) for the blank spike and from 48% to 73% (25% difference) in samples presented with 100× diluted tea matrix, representing increments of 54.5% and 52.1%, respectively. The presence of 10× diluted tea matrix slightly affected the improvement from ML incorporation, which revealed a 21% higher IP than singly ULSA alone (from 45% to 66%), equivalent to 46.7% increase.

Molecular annotation via spectral library search typically ranks based on score of matched reference spectra to measure annotation confidence [12, 16]. We also compared the performance by defining a hit as a single reference spectrum match upon ranking by matching score. Ranking was performed for library search-suggested annotations and ML-assisted annotations using "FinalScoreRatio" and by "$P(TP)$" without a defined acceptance threshold. An alternative IP, as defined in **eq8** for ranking analysis, represented the occurrence frequency of *TP* annotations in the top-ranked hits. As shown in **Fig. 4**, IPs were improved slightly (1–

420   3%) on average for samples with tea matrices due to ML incorporation, while remaining

421   comparable for spiked blanks. This finding emphasized the importance of considering

422   identification ambiguity of a candidate compound from a collection of reference spectral and

423   multiple platforms rather than a single reference spectrum hit [18]. *TP* hit rates were calculated

424   as the ratio of RPLC/HRMS tensors with any *TP* hit at least once in the top ranks to the total

425   identified spectra. Library search achieved *TP* hit rates of 87–88%, 83–87%, and 76–81% for

426   top-5, top-3, and top-1 hits, respectively. Our assessments found that ML incorporation

427   improved these rates to 89–95%, 87–90%, and 78–83%.

428   Collectively, chemical IP could be enhanced by applying our developed models. The

429   improvement would be greater if the number of reference spectra increased. This observation

430   was supported by our results of top-1, top-3, top-5 analyses and the ambiguity analyses that

431   included a collection of matched reference spectra (i.e., top-many analyses).

432   3.5.   Potentials and Limitations

433   In this work, we demonstrated a novel approach of integrative ML models to enhance

434   chemical IP in RPLC/HRMS-based NTA. Our results demonstrated the effectiveness of our

435   developed models, particularly Model 3, which integrates RTI error computed by the RTI

436   values on a harmonized scale from 2 QSRR-based models to enhance molecular annotation IPs.

437   Model 1 (the MF-to-RTI model) exhibited transferable RTI prediction and its exchangeability

438   with Model 2 (the CNL-to-RTI model). Model 2 showed strong predictive power for the

439   expected RTIs, capable to predict RTI values of chemical compounds beyond the data used to

440   modeling. The integration of ML incorporation with reference spectral library searches in our

441   identification workflow for the measurement of ambiguity not only improved recall in real

442   samples with and without tea matrices but also increased the confidence of annotation

443   acceptance compared to single library searches. Application of Model 3 resulted in acceptable

444   weighted *F*1 and MCC scores for the blank spike samples, alongside notable increases in IPs

445 with diluted tea matrix. Greater IP improvements were observed as the number of reference

446 spectra increased. In our work, for both tea matrix-containing and matrix-free samples,

447 computing a chemical IP from its corresponding collection of all available reference MS/MS

448 spectra improved half in annotation confidence. Our work provided a higher confidence *in*

449 *silico* analysis solution for the early-stage data analytics, aiding to shortlist chemicals-of-

450 interest for further chemical identification by reference standards.

451 However, a primary limitation inherent in our approach is the reliance on the quality and

452 diversity of the spectral reference libraries. Incomplete or biased libraries may hinder the

453 identification of highly structurally diverse compounds or those with limited accessible

454 reference spectral data. To address this challenge, it is essential to continuously update and

455 expand the spectral libraries, ensuring they encompass a broad spectrum of compounds

456 encountered in real-world samples. Furthermore, ongoing validation with various real-world

457 samples of *TP* and *TN* MS/MS spectra will be essential to fine-tune the models and ensure their

458 reliability across various analytical conditions.

459 **Abbreviations**

| | |
|---|---|
| 2D | Two-dimensional |
| AD | Applicability domain |
| Amide | C3–14 n-Alkylamide system |
| APC2D | AtomPairs2DFingerprintCount |
| AUPRC | Area under the precision-recall curve |
| AUROC | Area under the receiver operating characteristics curve |
| CNL | Cumulative neutral loss |
| Cocamide | C0–23 Cocamide diethanolamine homologous series |
| CV | Cross-validation |
| DT | Decision tree |
| EDA | Exploratory data analysis |
| ESI+ | Positive electrospray ionization |
| FDR | False discovery rate |
| *FN* | False negative |
| *FP* | False positive |
| HRMS | High-resolution mass spectrometry |

| | |
|---|---|
| IC | Identification confidence |
| IP | Identification probability |
| KNN | k-Nearest neighbors |
| LR | Logistic regression |
| *m/z* | Mass-to-charge ratio |
| MCC | Matthews correlation coefficient |
| MF | Molecular fingerprint |
| ML | Machine learning |
| MOIs | Masses-of-interest |
| MoNA | MassBank of North America |
| MS/MS | Tandem mass spectrometry |
| NTA | Non-targeted analysis |
| $P(TP)$ | Probability of true positive |
| $P(ULSA)$ | Probability of true positive from ULSA determination |
| PFAS | Per- and polyfluoroalkyl substances |
| PubChem | PubChem Fingerprinter |
| QSRR | Quantitative structure-retention relationship |
| $R^2$ | Coefficient of correlation |
| RF | Random forest |
| RI | Retention index |
| RMSE | Root mean square error |
| RPLC | Reversed-phase lipid chromatography |
| RT | Retention time |
| RTI | Retention time index |
| SAFD | Self-adjusting feature detection |
| SMILES | Simplified molecular-input line-entry system |
| *TN* | True negative |
| *TP* | True positive |
| ULSA | Universal Library Search Algorithm |
| UoA | the RI system developed at the University of Athens |

460

**Supplementary Information**

A word file with details regarding modeling, identification workflow, computational

power, libraries used, and the additional figures and tables for the supporting text.

**Author Information**

▪ **Disclosure**

- 21 -

The authors declare that they have no conflicts of interest or personal relationships that could influence the work reported in this paper.

- ▪ **Author Contributions**

Conceptualization, H.-L.N., V.T., and S.S.; resources, H.-L.N., V.T., D.v.H., Z.C., and S.S.; data curation, H.-L.N., V.T., and D.v.H.; writing-original draft preparation, H.-L.N.; writing-review and editing, H.-L.N., V.T., Y.H., and S.S.; visualization, H.-L.N. and Y.H.; supervision, Y.H., Z.C., and S.S.; project administration, Z.C. and S.S.; funding acquisition, Z.C. and S.S. All authors have read and agreed to the published version of the manuscript.

## Acknowledgements

## References

1. Manz KE, Feerick A, Braun JM, et al (2023) Non-targeted analysis (NTA) and suspect screening analysis (SSA): A review of examining the chemical exposome. J Expo Sci Environ Epidemiol 33:524–536.
2. Samanipour S, Barron LP, van Herwerden D, et al (2024) Exploring the chemical space of the exposome: How far have we gone? JACS Au 4:2412–2425.
3. Guo Z, Zhu Z, Huang S, Wang J (2020) Non-targeted screening of pesticides for food analysis using liquid chromatography high-resolution mass spectrometry-a review. Food Additives and Contaminants - Part A 37:1180–1201.
4. Zweigle J, Bugsel B, Zwiener C (2022) FindPFΔS: Non-target screening for PFAS-comprehensive data mining for MS2 fragment mass differences. Anal Chem 94:10788–10796.
5. Chen Y, Ngan HL, Song Y, et al (2024) Chronic real-ambient PM2.5 exposure exacerbates cardiovascular risk via amplifying liver injury in mice fed with a high-fat and high-cholesterol diet.

498         Environment & Health 2:221–232.

6. Sumner LW, Amberg A, Barrett D, et al (2007) Proposed minimum reporting standards for chemical analysis: Chemical Analysis Working Group (CAWG) Metabolomics Standards Initiative (MSI). Metabolomics 3:211–221.

7. Schymanski EL, Jeon J, Gulde R, et al (2014) Identifying small molecules via high resolution mass spectrometry: Communicating confidence. Environ Sci Technol 48:2097–2098.

8. Hollender J, Schymanski EL, Ahrens L, et al (2023) NORMAN guidance on suspect and non-target screening in environmental monitoring. Environ Sci Eur 35:75.

9. Alygizakis N, Lestremau F, Gago-Ferrero P, et al (2023) Towards a harmonized identification scoring system in LC-HRMS/MS based non-target screening (NTS) of emerging contaminants. TrAC - Trends in Analytical Chemistry 159:116944.

10. Ciccarelli D, Samanipour S, Rapp-Wright H, et al (2024) Bridging knowledge gaps in human chemical exposure via drinking water with non-target screening. Crit Rev Environ Sci Technol 55:190.

11. Hulleman T, Turkina V, O'Brien JW, et al (2023) Critical assessment of the chemical space covered by LC-HRMS non-targeted analysis. Environ Sci Technol 57:14101–14112.

12. Samanipour S, Reid MJ, Bæk K, Thomas K V. (2018) Combining a deconvolution and a universal library search algorithm for the nontarget analysis of data-independent acquisition mode liquid chromatography-high-resolution mass spectrometry results. Environ Sci Technol 52:4694–4701.

13. Scheubert K, Hufsky F, Petras D, et al (2017) Significance estimation for large scale metabolomics annotations by spectral matching. Nat Commun 8:1494.

14. Bittremieux W, Schmid R, Huber F, et al (2022) Comparison of cosine, modified cosine, and neutral loss based spectrum alignment for discovery of structurally related molecules. J Am Soc Mass Spectrom 33:1733–1744.

15. Ji H, Xu Y, Lu H, Zhang Z (2019) Deep MS/MS-aided structural-similarity scoring for unknown metabolite identification. Anal Chem 91:5629–5637.

16. Zulfiqar M, Gadelha L, Steinbeck C, et al (2023) MAW: The reproducible Metabolome Annotation Workflow for untargeted tandem mass spectrometry. J Cheminform 15:32.

17. Li Y, Kind T, Folz J, et al (2021) Spectral entropy outperforms MS/MS dot product similarity for small-molecule compound identification. Nat Methods 18:1524–1531.

18. Metz TO, Chang CH, Gautam V, et al (2024) Introducing "identification probability" for automated and transferable assessment of metabolite identification confidence in metabolomics and related studies. Anal Chem 97:1–11.

19. Assis J, Serrão EA, Fragkopoulou E, et al (2024) Misconception of model transferability precludes estimates of seagrass community reorganization in a changing climate. Nat Plants 10:1071–1074.

20. Munro K, Miller TH, Martins CPB, et al (2015) Artificial neural network modelling of pharmaceutical residue retention times in wastewater extracts using gradient liquid chromatography-high resolution mass spectrometry data. J Chromatogr A 1396:34–44.

21. McEachran AD, Mansouri K, Newton SR, et al (2018) A comparison of three liquid chromatography (LC) retention time prediction models. Talanta 182:371–379.

22. Aalizadeh R, Nika MC, Thomaidis NS (2019) Development and application of retention time prediction models in the suspect and non-target screening of emerging contaminants. J Hazard Mater 363:277–285.

23. Feng C, Xu Q, Qiu X, et al (2021) Evaluation and application of machine learning-based retention time prediction for suspect screening of pesticides and pesticide transformation products in LC-HRMS. Chemosphere 271:129447.

24. Song D, Tang T, Wang R, et al (2024) Enhancing compound confidence in suspect and non-target screening through machine learning-based retention time prediction. Environmental Pollution 347:123763.

25. Bonini P, Kind T, Tsugawa H, et al (2020) Retip: Retention time prediction for compound annotation in untargeted metabolomics. Anal Chem 92:7515–7522.

26. Stanstrup J, Neumann S, Vrhovšek U (2015) PredRet: Prediction of retention time by direct mapping between multiple chromatographic systems. Anal Chem 87:9421–9428.

27. Ruttkies C, Schymanski EL, Wolf S, et al (2016) MetFrag relaunched: Incorporating strategies beyond in silico fragmentation. J Cheminform 8:3.

28. Bach E, Szedmak S, Brouard C, et al (2018) Liquid-chromatography retention order prediction for metabolite identification. Bioinformatics 34:i875–i883.

- 23 -

555   29.   Kretschmer F, Harrieder EM, Hoffmann MA, et al (2024) RepoRT: A comprehensive repository for
556         small molecule retention times. Nat Methods 21:153–155.
557   30.   Rigano F, Arigò A, Oteri M, et al (2021) The retention index approach in liquid chromatography: An
558         historical review and recent advances. J Chromatogr A 1640:461963.
559   31.   Aalizadeh R, Alygizakis NA, Schymanski EL, et al (2021) Development and application of liquid
560         chromatographic retention time indices in HRMS-based suspect and nontarget screening. Anal Chem
561         93:11601–11611.
562   32.   Boelrijk J, van Herwerden D, Ensing B, et al (2023) Predicting RP-LC retention indices of
563         structurally unknown chemicals from mass spectrometry data. J Cheminform 15:28.
564   33.   van Herwerden D, Nikolopoulos A, Barron LP, et al (2024) Exploring the chemical subspace of
565         RPLC: A data driven approach. Anal Chim Acta 1317:342869.
566   34.   Yap CW (2011) PaDEL-descriptor: An open source software to calculate molecular descriptors and
567         fingerprints. J Comput Chem 32:1466–1474.
568   35.   Hall LM, Hill DW, Menikarachchi LC, et al (2015) Optimizing artificial neural network models for
569         metabolomics and systems biology: An example using HPLC retention index data. Bioanalysis
570         7:939–955.
571   36.   Aalizadeh R, Nikolopoulou V, Thomaidis NS (2022) Development of liquid chromatographic
572         retention index based on cocamide diethanolamine homologous series (C(n)-DEA). Anal Chem
573         94:15987–15996.
574   37.   Xue J, Guijas C, Benton HP, et al (2020) METLIN MS2 molecular standards database: A broad
575         chemical and biological resource. Nat Methods 17:953–954.
576   38.   Aisporna A, Benton HP, Chen A, et al (2022) Neutral loss mass spectral data enhances molecular
577         similarity analysis in METLIN. J Am Soc Mass Spectrom 33:530–534.
578   39.   Mohanty I, Mannochio-Russo H, Schweer J V., et al (2024) The underappreciated diversity of bile
579         acid modifications. Cell 187:1801–1818.
580   40.   van Herwerden D, O'Brien JW, Lege S, et al (2023) Cumulative neutral loss model for fragment
581         deconvolution in electrospray ionization high-resolution mass spectrometry data. Anal Chem
582         95:12247–12255.
583   41.   Guo H, Xue K, Sun H, et al (2023) Contrastive learning-based embedder for the representation of
584         tandem mass spectra. Anal Chem 95:7888–7896.
585   42.   Herwerden D van, Kant E, Jackson M, et al (2024) Modular open-access and open-source julia
586         language toolbox for processing of HRMS data: jHRMSToolBox
587   43.   Chambers MC, MacLean B, Burke R, et al (2012) A cross-platform toolkit for mass spectrometry
588         and proteomics. Nat Biotechnol 30:918–920.
589   44.   Samanipour S, O'Brien JW, Reid MJ, Thomas K V. (2019) Self adjusting algorithm for the
590         nontargeted feature detection of high resolution mass spectrometry coupled with liquid
591         chromatography profile data. Anal Chem 91:10800–10807.
592
593

## Formulae and Equations

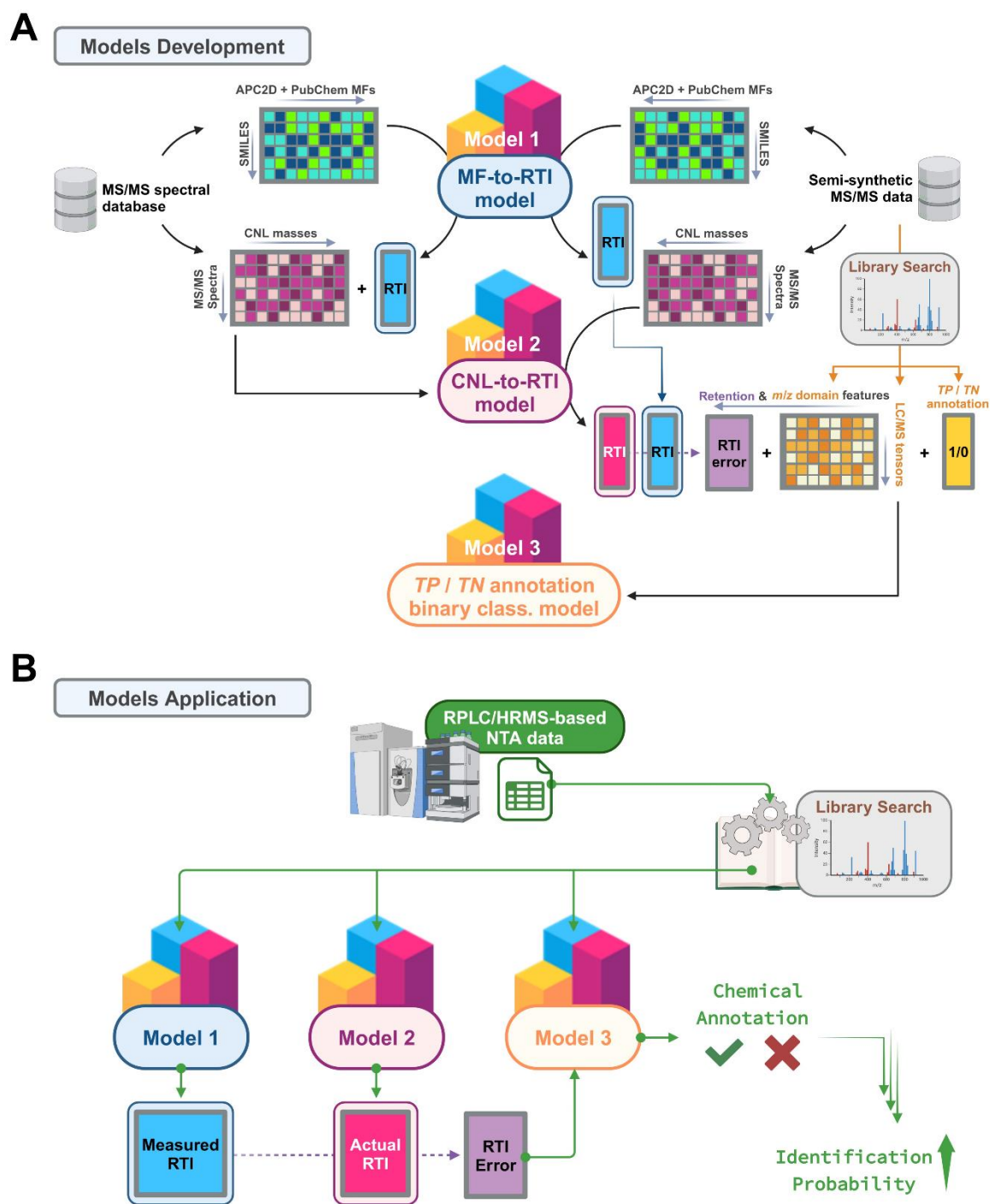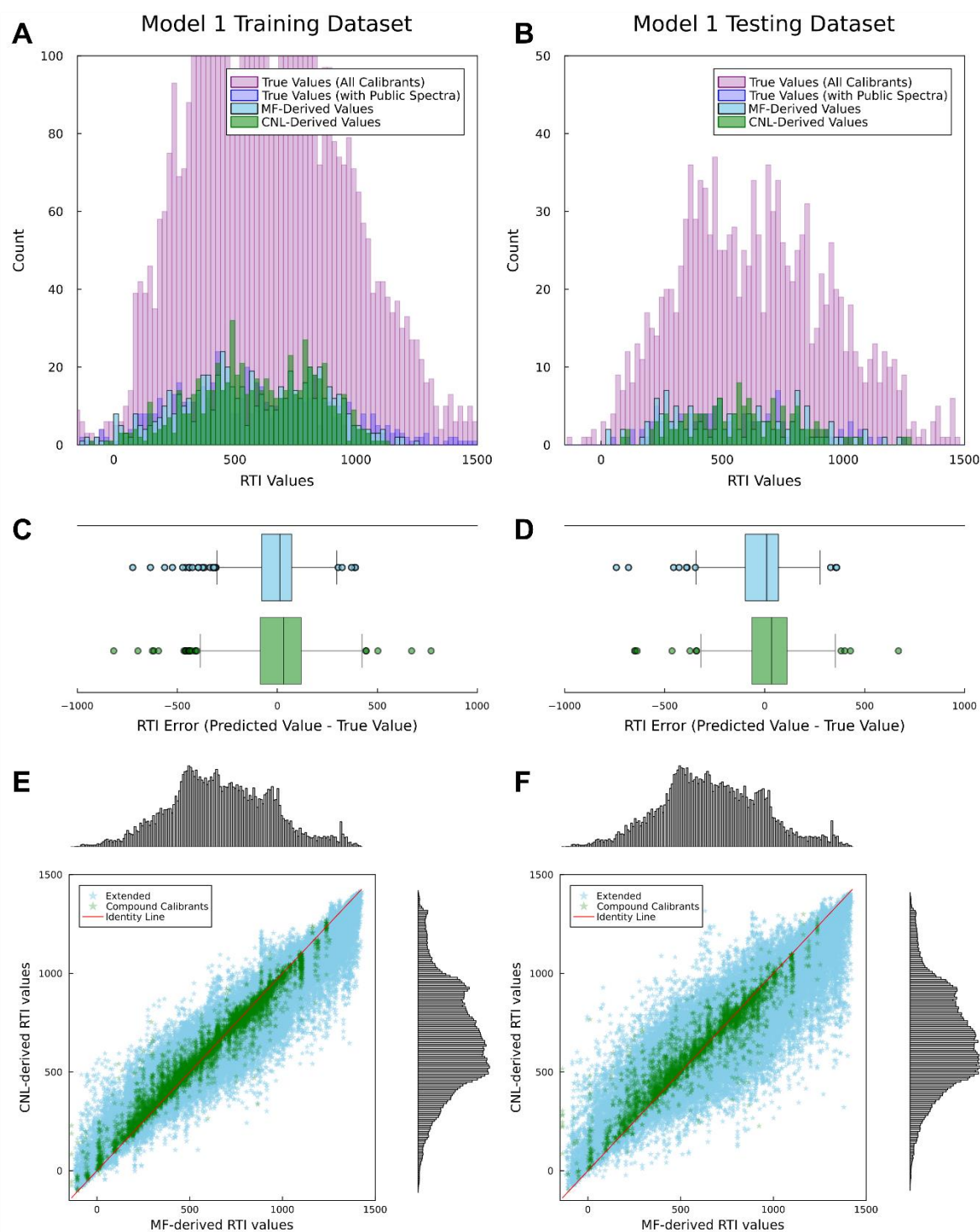| Eq1 | $$RMSE = \sqrt{\frac{1}{N}\sum_{i=1}^{N}(y_i - \hat{y}_i)^2}$$ |
|---|---|
| Eq2 | $$S_8 = \frac{f_{matched}}{f_{user}} + \frac{f_{matched}}{f_{reference}} + \frac{MS_{tol} - MS_{err,pre}}{MS_{tol}}$$ $$+ \frac{f_{matched}}{f_{reference}} \cdot \frac{MS_{tol} - \left|MS_{err,frag}\right|}{MS_{tol}} + \frac{f_{matched}}{f_{reference}} \cdot \frac{2 \cdot MS_{tol} - STD_{err,frag}}{2 \cdot MS_{tol}}$$ $$+ MF_{forward} + MF_{reverse}$$ $$= S_1 + S_2 + \frac{MS_{tol} - S_3}{MS_{tol}} + S_2 \cdot \frac{MS_{tol} - S_4}{MS_{tol}} + S_2 \cdot \frac{2 \cdot MS_{tol} - S_5}{2 \cdot MS_{tol}} + S_6 + S_7$$ |
| Eq3 | $$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$ |
| Eq4 | $$F1 = \frac{TP}{TP + \frac{1}{2}(FP + FN)}$$ |
| Eq5 | $$RTI\ error = RTI_{MF} - RTI_{CNL}$$ |
| Eq6 | $$h_{ii} = x_i^T(X^TX)^{-1}x_i$$ |
| Eq7 | $$Recall = \frac{TP}{TP + FN}$$ |
| Eq8 | $$P(Exposome\ Identification\ for\ each\ RPLC/HRMS\ tensor) = \frac{1}{N}$$ |
| Eq9 | $$P(Exposome\ Identification\ for\ each\ RPLC/HRMS\ tensor\ in\ Top$$ $$- X\ of\ Hits\ of\ reference\ spectra) = \frac{N}{X}$$ |

595

- 25 -

**Figure Legends:**

**Figure 1: A diagram of the roadmap in this study**. **(A)** Three models were developed in this study. Model 1 was a random forest (RF) regression (reg.) model for predicting molecular fingerprint (MF)-derived retention time index (RTI). Model 2 was a RF reg. model for predicting cumulative neutral loss (CNL)-derived RTI. A set of semi-synthetic MS/MS data was prepared. Models 1 and 2 were deployed to predict the expected RTI values for the semi-synthetic data. These data were processed by library search, outputting spectral matching scores for Model 3 training. Model 3 was a k-Nearest neighbors (KNN) binary classification (class.) model was trained by incorporating features from retention and *m/z* domains to predict true positive (*TP*) true negative (*TN*) chemical annotation. **(B)** The 3 developed models were applied to analyze reversed-phase liquid chromatography/high-resolution mass spectrometry (RPLC/HRMS)-based non-targeted analysis (NTA) data. By incorporating the models, the chemical identification probability could be enhanced. *Abbreviations: APC2D, AtomPairs2DFingerprintCount; PubChem, PubChem Fingerprinter.*
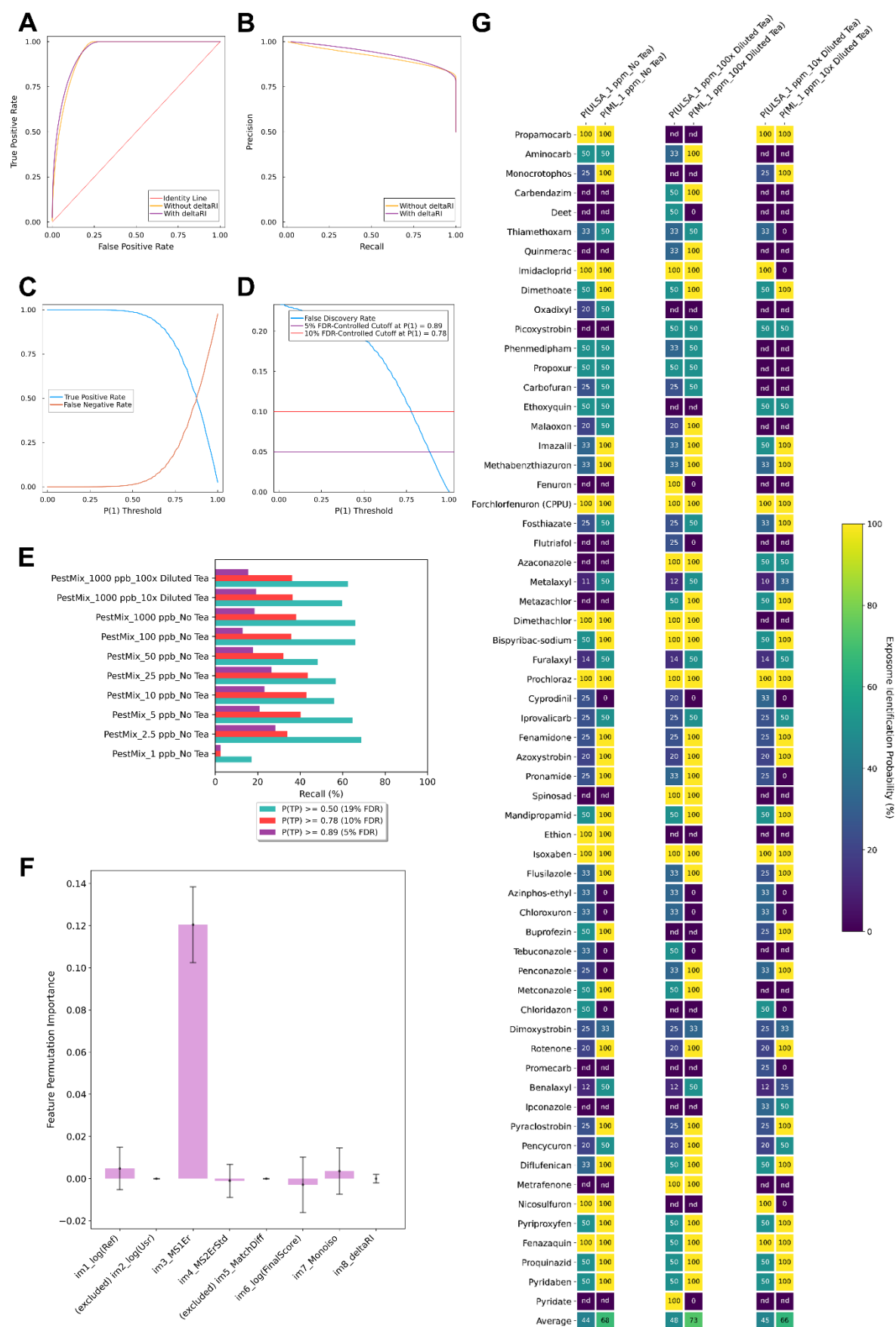
**Figure 2: Performance assessment of the quantitative structure-retention relationship (QSRR)-based models.** Histograms illustrating the distributions of true retention time indices (RTIs) of calibrants used to **(A)** train and **(B)** test Model 1 and the values predicted from Models 1 and 2. Boxplots presenting the variances of the predicted values to the true values of calibrants used to **(C)** train and **(D)** test Model 1. Correlation plots showing correspondence between molecular fingerprint (MF)-derived and cumulative neutral loss (CNL)-derived RTIs for **(E)** training and **(F)** testing datasets of Model 2. Green dots representing compounds used in Model 1's training, yet sky-blue dots indicating new compounds within the RPLC chemical space after applicability domain filtering.

- 26 -

**Figure 3: Performance assessment of the true positive/true negative annotation binary classification model. (A)** Receiver operating characteristics curves and **(B)** precision-recall curves for the retention time error (RTI) inclusive and exclusive models. **(C)** True positive (*TP*) and false negative rates and **(D)** false discovery rates (FDR) across various probability of true positive ($P(1)$ or $P(TP)$) thresholds. **(E)** Recalls of individual LC-MS/MS samples at different $P(TP)$ cut-off thresholds. Pesticide standards were part of the LC/MS Pesticide Comprehensive Mix Kit (Part Number 5190-0551) from Agilent Technologies (PestMix). Pesticide mixtures that were spiked in solvent blank and 10× and 100× diluted black tea matrix respectively to result 1 ppm final solutions were analyzed. For blank spike (No Tea) samples, various final solution concentrations (1, 2.5, 5, 10, 25, 50, 100, and 1,000 ppb) were analyzed. **(F)** Feature permutation importance analysis for all k-nearest neighbors models with Matthews correlation coefficient scores >0.30 on blank spike dataset during model selection and hyperparameter tuning. **(G)** Integration of machine learning (ML) with Universal Library Search Algorithm (ULSA) enhanced the average identification probabilities by 24% and 25% for the blank spike and the samples with the presence of 100× diluted tea matrix, equivalent to 54.5% and 52.1% increments than solely ULSA implementation, yet the presence of 10× tea matrix reduced this improvement to and 21%, equivalent to 46.7% increment. *Abbreviations: nd, not detected.*

**Figure 4: Performance assessment by ranking.** Integration of machine learning (ML) with Universal Library Search Algorithm (ULSA) slightly enhanced the average identification probability by 1–3% in top-5, top-3, and top-1 hits. True positive (*TP*) hit rates were the ratio of molecular annotations with any *TP* hit at least once in the top ranks to the total identified spectra. *Abbreviations: nd, not detected.*
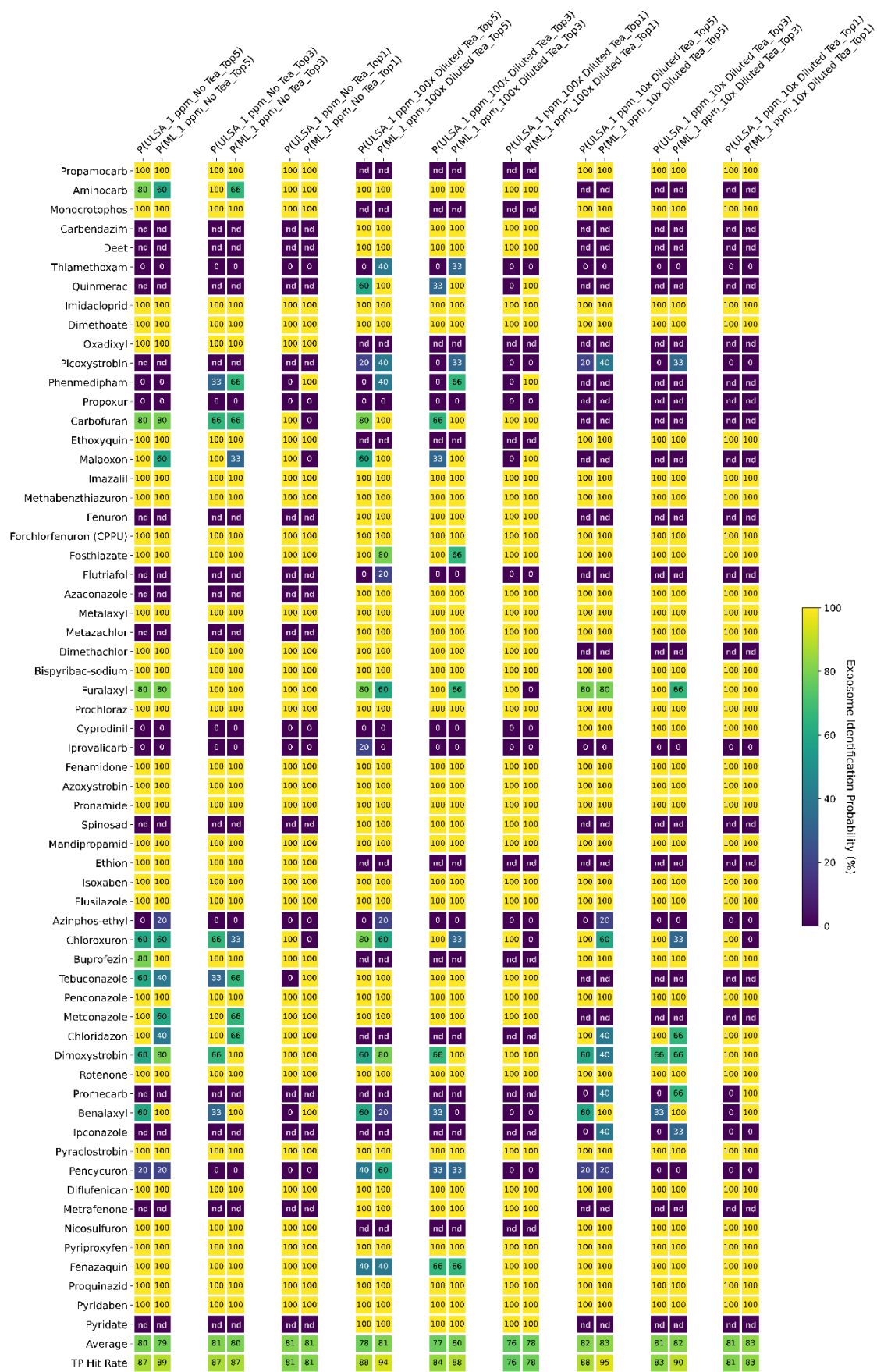
**Figure 1**

647    **Figure 2**

648

649     **Figure 3**

650

**Figure 4**

- 31 -

**Table 1: Summary of the Machine Learning Models' Predictive Power.**

| ML[a] model | Model 1 | Model 2 | | Model 3 (with feature "DeltaRi"; without "DeltaRi") | | |
|---|---|---|---|---|---|---|
| Function | to predict RTI[b] from MF[c] | to predict RTI[b] from CNL[d] | | to predict acceptance/rejection decision of a ULSA[e]-suggested molecular annotation from a MS/MS spectrum | | |
| Performance metric | | $R^2$ | RMSE[f] | Weighted $F1$ score | MCC[g] | Recall |
| Training | | 0.9634 | 51.8 | 0.89; 0.90 | 0.77; 0.79 | 0.99; 1.00 |
| Testing | | 0.8824 | 92.8 | 0.89; 0.89 | 0.77; 0.79 | 0.99; 0.99 |
| Validation | | | | 0.65; 0.65 | 0.30; 0.30 | 0.66; 0.66 |
| Real Sample | | | | | | 0.60; 0.54 |

[a] Machine learning.

[b] Retention time index.

[c] Molecular fingerprint.

[d] Cumulative neutral loss.

[e] Universal Library Search Algorithm.

[f] Root mean square error.

[g] Matthews correlation coefficient.