1   **IDSL.UFA assigns high confidence molecular formula annotations for untargeted LC/HRMS**

2   **datasets in metabolomics and exposomics**

3   Sadjad Fakouri Baygi[1], Sanjay K Banerjee[2], Praloy Chakraborty[2], Yashwant Kumar[2], Dinesh Kumar

4   Barupal[1]*

5   [1] Department of Environmental Medicine and Public Health, Icahn School of Medicine at Mount Sinai,

6   New York, NY, 10029, USA

7   [2] Non-communicable Diseases Division, Translational Health Science and Technology Institute,

8   Faridabad, Haryana, 121001, India

9   * Corresponding author: Address: CAM Building, 3rd floor, 17 E 102nd St, New York, NY 10029 Email:

10  dinesh.kumar@mssm.edu, phone: +1-530-979-4354

11

12  Abstract

13  Untargeted LC/HRMS assays in metabolomics and exposomics aim to characterize the small molecule
14  chemical space in a biospecimen. To gain maximum biological insights from these datasets, LC/HRMS
15  peaks should be annotated with chemical and functional information including molecular formula, structure,
16  chemical class and metabolic pathways. Among these, molecular formulas may be assigned to LC/HRMS
17  peaks through matching theoretical and observed isotopic profiles (MS1) of the underlying ionized
18  compound. For this, we have developed the Integrated Data Science Laboratory for Metabolomics and
19  Exposomics – United Formula Annotation (IDSL.UFA) R package. In the untargeted metabolomics
20  validation tests, IDSL.UFA assigned 54.31%-85.51% molecular formula for true positive annotations as the
21  top hit, and 90.58%-100% within the top five hits. Molecular formula annotations were also supported by
22  MS/MS data. We have implemented new strategies to 1) generate formula sources and their theoretical
23  isotopic profiles 2) optimize the formula hits ranking for the individual and the aligned peak lists and 3) scale
24  IDSL.UFA-based workflows for studies with larger sample sizes. Annotating the raw data for a publicly
25  available pregnancy metabolome study using IDSL.UFA highlighted hundreds of new pregnancy related
26  compounds, and also suggested presence of chlorinated perfluorotriether alcohols (Cl-PFTrEAs) in human
27  specimens. IDSL.UFA is useful for human metabolomics and exposomics studies where we need to
28  minimize the loss of biological insights in untargeted LC/HRMS datasets. The IDSL.UFA package is
29  available in the R CRAN repository https://cran.r-project.org/package=IDSL.UFA. Detailed documentation
30  and tutorials are also provided at www.ufa.idsl.me.

31  Introduction

32  Untargeted LC/HRMS analyses of human specimens enable studying the metabolome and exposome in
33  an unbiased manner[1, 2].They have delivered many novel biomarkers and mechanisms for diseases and
34  have improved our understanding of basic metabolic pathways[3-5]. These assays are unique in nature since
35  they record all the mass to charge (m/z) ratio signals above the limit of detection of an instrument for ionized
36  compounds in a sample[6]. This makes the collected data a rich source of information with great opportunities
37  to generate novel hypotheses about metabolome and exposome. It is critical for the promises that
38  untargeted assay offers, that the data are utilized in an inclusive way to not miss any discovery
39  opportunities.

40  A key post-data acquisition step in the untargeted LC/HRMS assays is to annotate the detected peaks
41  with a range of structural and functional information which can enable biological interpretations[3, 7, 8]. This
42  information includes a chemical structure, molecular formula, chemical class and metabolic pathway[1, 9, 10].
43  These annotations may help in understanding the nature, origin and function of the chemical structure
44  underlying a peak. Among these information, molecular formula can be assigned to a LC/HRMS peak
45  using the observed and theoretical isotopic profiles for a chemical compound.[11] Isotopic profiles are
46  distinguishable mass spectral signature that represent atomic masses and their natural abundances in
47  the molecular formulas of a compound.[12] Despite the known limitations of high-resolution mass
48  spectrometry instruments, observed experimental isotopic profiles for an ionized compound may
49  sufficiently match the theoretical counterpart within instrument errors in many instances[11, 13], allowing to
50  annotate LC/HRMS peaks with molecular formula[14]. Peak annotation by isotopic profile matching should
51  be performed using efficient computational strategies to account for instrumental errors, multi-sample
52  studies, biological plausibility and chemical diversity.[15]

53  There has been a great deal of efforts to develop computational tools for annotating peaks in a LC/HRMS
54  dataset with MS1 only data. In a MS1 peak list, a series of m/z values representing different isotopes, ESI
55  adducts, and in-source fragments can belong to one compound. Grouping these m/z values are normally
56  performed by retention time and elution profile similarities within a single file, for example by *xcms*-
57  CAMERA[16, 17], and peak intensity correlations across multiple samples such as MS-FLO[18] and CliqueMS[9]
58  tools. Clustered isotopologues from these tools can be used by the Rdisop R package[19] to assign
59  molecular formulas in a 'database independent' manner. But this approach may miss expected
60  compounds for a sample due to MS instrument's sensitivity and specificity. MetDNA[7] can search for
61  theoretical isotope profiles for a list of molecular formulas from a metabolic reaction network database in

62   the MS1 peak list. However, their 'database dependent' approach is prone to miss 1) exposure
63   compounds that are poorly represented in such biochemical databases and 2) compounds which may not
64   have any transformation products because of their bioaccumulative nature and 3) compounds that were
65   filtered out by the detection frequency and intensity thresholds while generating MS1 peak table for a
66   study. Moreover, MetDNA[7] and other tools including SIRIUS[20] and ZODIAC[21], NetID[22] are mainly
67   designed for assigning molecular formulas to peaks having MS/MS fragmentation data. Furthermore,
68   implementing these tools for larger studies where only MS1 data are available for every sample remains
69   to be challenging due to the ranking of formula hits on individual and aligned peak tables, scalable
70   computation and various sources for formulas which need to be covered for exposomics projects.

71   There is a need to develop new tools to compute and to compare theoretical and experimental isotopic
72   profiles for chemical lists from larger databases and chemical spaces for molecular formula annotations.
73   Here, we have developed a scalable, user-friendly, thoroughly tested R package, the IDSL.UFA to assign
74   molecular formulas with high confidence to peaks in untargeted LC/HRMS datasets from large-scale
75   studies. IDSL.UFA covers major possible situations in which a molecular formula can be assigned to
76   LC/HRMS peaks. We propose that processing LC/HRMS data with IDSL.UFA can find new opportunities
77   for hypotheses and biomarker discoveries for studying the role of metabolism and exposome in human
78   diseases.

79   Methods
80   **Publicly available LC/HRMS test datasets:** To test and develop the IDSL.UFA R package, we have
81   utilized the raw LC/HRMS data for human and mouse biospecimen studies (MTBLS1684[23],
82   MTBLS2542[24], ST001683[25], ST001430[26], ST001154[27], ST002044 and reference authentic standards
83   (MSV000088661) available from Metabolomics WorkBench (https://www.metabolomicsworkbench.org/),
84   MassIVE (https://www.massive.ucsd.edu), and MetaboLights (https://www.ebi.ac.uk/metabolights)
85   repositories. Data processing results that we have generated for these studies have been submitted to
86   the Zenodo.org repository and corresponding entry pages are provided in Table S.1. Sample preparation
87   and data collection procedures are available at entry pages for these studies in the repositories.

88   **Data analysis setup:** IDSL.UFA R package is available in the R-CRAN repository (https://cran.r-
89   project.org/package=IDSL.UFA). The package was installed using the '*install.packages("IDSL.UFA")*' R
90   command. The IDSL.MXP package (https://cran.r-project.org/package=IDSL.MXP) was used to read
91   mzML/mzXML/netCDF mass spectrometry data in the centroid mode. mzML files were generated from
92   the vendor specific data format using the ProteoWizard MSConvert utility[28] when needed. All data files
93   related to only one type of analysis such as "reverse phase - electrospray ionization negative mode" were
94   stored in a single file folder. Figure 1 (simplified) and Figure S.1 (detailed) show the workflow steps to
95   assign molecular formulas for a study. Data processing parameters for IDSL.UFA were provided in a
96   Microsoft excel file (https://zenodo.org/record/6466688) which was created for individual test studies. We
97   have provided the parameters files used in this manuscript in the Zenodo.org repository at
98   (https://zenodo.org/record/6466684). To run the IDSL.UFA workflow, only a single R command
99   '*UFA_workflow*(*spreadsheet* = *"address of the parameter xlsx file"*)' was needed. Tutorials to create the
100  parameter files for different scenarios are available at (https://ufa.idsl.me). For each individual peak list in
101  a study, a formula annotation list with rank, score and other peak properties was generated and exported
102  to a csv file. Likewise, for each peak in the aligned peak table, top 5-20 formulas with detection frequency
103  and median ranks across all samples were exported to a csv file for each test study.

104  **Generating the isotopic profile database (IPDB):** An IPDB is a digital collection of theoretical isotopic
105  profiles computed by the IDSL.UFA R package for a list of candidate molecular formulas. IDSL.UFA
106  queries and matches the experimental isotopic profile against this collection to annotate a LC/HRMS
107  peak. To compute the isotopic profile for a molecular formula, we have utilized the reference stable
108  isotope masses and abundances for elements in the periodic table from the PubChem database entries[29]
109  which have been sourced from International Union of Pure and Applied Chemistry (IUPAC)[30]. We have
110  also provided an online tool (https://.ipc.idsl.me) to compute an isotopic profile for a single molecular

111    formula. IDSL.UFA generates centroid isotopic profiles using a dynamic intensity threshold and a peak-
112    spacing criterion to merge adjacent isotopologues within a mass accuracy window. In this work, we have
113    covered two sources of molecular formulas.

114    **Source A (databases**): Chemical compound lists for four key databases in metabolomics and
115    exposomics including the blood exposome (chemicals expected in a mammalian blood specimen),
116    RefMet (measured and expected small molecules in biological organisms)[31], Lipid Maps (known lipid
117    molecules)[32] and the US-Food and Drug Administration substance registry[33] were obtained from their
118    online web addresses. These four databases were combined into a single compound list referenced as
119    IDSL.ExposomeDB in this manuscript and also provided at the Zenodo repository
120    (https://zenodo.org/record/5823455). Charged compounds, isotope-labeled compounds and multi-
121    components were excluded. Unique molecular formulas from this consolidated database were used for
122    computing IPDB. IPDBs for these four databases and the environmental protection agency (EPA)
123    CompTox Chemicals Dashboard[34] are available at the Zenodo repository
124    (https://zenodo.org/record/5823455).

125    **Source B (enumerated chemical space with constraints**): Molecular formulas were enumerated using
126    a set of combinatorial and filtering rules using C, H, As, B, Br, Cl, F, I, K, N, Na, O, P, S, Se, and Si
127    elements. These 16 elements were able to cover 93.76% of carbon-containing compounds ($50 \leq$ mass $\leq$
128    2000) in the IDSL.ExposomeDB combined with EPA chemistry Dashboard[34]. An enumerated chemical
129    space (ECS) can be represented using equation (1).

$$C_cH_hAs_{as}B_bBr_{br}Cl_{cl}F_fI_iK_kN_nNa_{na}O_oP_pS_sSe_{se}Si_{si} \tag{1}$$

130    where the subscripts of elements represent the number of atoms. A fully combinatorial chemical space
131    from above-mentioned 16 elements is impractical to be managed by current computational resources.
132    Therefore, we derived and coded in R a set of four rules which were inspired from the seven golden rules
133    approach[35] to constrain ECSs. These rules included **1) C/N chemical space rule** '((c/2-n-1) $\leq$
134    (h+cl+br+f+i) $\leq$ (2c+3n+6))' was used to set elemental boundaries for the organic compounds to ensure
135    entire moieties are bond to carbon and nitrogen atoms. **2) Extended SENIOR rule** was used to ensure
136    that the molecular formulas completely filled s- and p- valence electron shells.[35] **3) Maximum number of
137    halogens thresholds** was used to constrain halogenated compounds. For example, we have used the
138    maximum number of (br+cl) $\leq$ 8 and the maximum number of ((br+cl+f+i) $\leq$ 31) thresholds to cover
139    halogenated compounds in the blood exposome database. **4) Maximum number of elements rule** was
140    used to skip unrealistically complex molecular formulas generated through molecular formula
141    enumeration. For example, the maximum number of elements for glucose ($C_6H_{12}O_6$) is three (C, H, and
142    O). The ECS boundaries and rules for the MTBLS1684 study are provided in the Zenodo repository
143    (https://zenodo.org/record/5838603).

144    **MS1 peak detection and alignment:** IDSL.IPA[36] R package (https://cran.r-
145    project.org/package=IDSL.IPA) was used to generate individual peak lists for each sample and the
146    aligned peak table (m/z-RT pairs across all samples) for each study. Data processing parameter files and
147    IDSL.IPA results for each test study are provided in the Zenodo repository (see Table S.1). Details and a
148    tutorial for IDSL.IPA data processing can be found at (https://ipa.idsl.me) site.

149    **Isotopic profile matching for individual sample:** First, IDSL.UFA software accessed the peak
150    boundaries, $^{12}C$ m/z, $^{13}C$ m/z and ratio of cumulated intensity of $^{12}C$ to $^{13}C$ ($R^{13}C$) for each peak in an
151    IDSL.IPA generated peak list for a sample. Next, it finds all the theoretical isotopic profiles in an IPDB that
152    matches the $^{12}C$ and $^{13}C$ m/z for a peak. Then, for each matched theoretical profile, experimental profiles
153    are retrieved from raw data using a mass accuracy threshold within the peak boundaries for a peak. If a
154    compound formula has three isotopologues in the IPDB and only two were observed in the raw data, the
155    formula will not be annotated. IDSL.UFA requires that a minimum one MS1 scan across the peak should
156    have the full isotope profile for a formula in the IPDB.

157  For the experimental isotopic profiles, the IDSL.UFA software calculates cumulated intensities and
158  intensity-weighted average masses for each isotopologue using equations (2) and (3) across the
159  chromatographic peak to minimize the effect of fluctuations such as peak saturation.

$$\overline{Int} = \sum_{t=t_0}^{t=t_{end}} Int_t \tag{2}$$

$$\overline{m/z} = \frac{\sum_{t=t_0}^{t=t_{end}} m/z_t * Int_t}{\overline{Int}} \tag{3}$$

160  where $m/z_t$ and $Int_t$ represent mass and intensity of the matched isotopologue in individual scans across
161  the chromatographic peak from $t_0$ to $t_{end}$.

162  We have used the Profile cosine similarity ($\overline{PCS}$) to quantify profile similarity between experimental and
163  theoretical isotopic profiles using equation (4). To assess mass accuracy error for whole isotopic profile,
164  Normalized Euclidean mass error ($\overline{NEME}$) was calculated using the equation (5).[11]

$$\overline{PCS} = \sum_{i=1}^{S} \frac{I_i^{theor} I_i^{exptl}}{\sqrt{\sum_{i=1}^{S} (I_i^{theor})^2} \sqrt{\sum_{i=1}^{S} (I_i^{exptl})^2}} \tag{4}$$

$$\overline{NEME} = \sqrt{\frac{\sum_{i=1}^{S} (M_i^{theor} - M_i^{exptl})^2}{S}} \tag{5}$$

165  where $I_i$, $M_i$, and $S$ represent the intensity of the isotopologue, mass of the isotopologues, and number of
166  isotopologues in the isotopic profile, respectively. Superscripts of *theor* and *exptl* also represent
167  theoretical and experimental isotopic profiles, respectively.

168  Candidate formulas were then filtered using thresholds for 1) $\overline{PCS}$ 2) $\overline{NEME}$ 3) the top 80% of number of
169  scans with the confirmed whole isotopic profile (NDCS) and 4) minimum percentage of NDCS within a
170  chromatography peak (RCS (%)). These linear cutoffs allow eliminating false positives; however, they can
171  reject true positive peaks with poor isotopic profiles.

172  Next, a matching score for each candidate filtered formula was computed using equation (6).

$$Score = \left( \frac{S^{coeff[1]} * \left(\frac{\overline{PCS}}{100}\right)^{coeff[2]} * \left(\frac{RCS}{100}\right)^{coeff[3]}}{\left(\frac{\overline{NEME}}{maxNEME}\right)^{coeff[4]} * \left(\exp\left(\left|\ln\left(\frac{\overline{R^{13}C_{PL}}}{R^{13}C_{IP}}\right)\right|\right)\right)^{coeff[5]}} \right) \tag{6}$$

173  where $\overline{R^{13}C_{PL}}$ and $R^{13}C_{IP}$ indicate experimental and theoretical $R^{13}C$ values, respectively. $R^{13}C$ values
174  represent the ratio of the general $^{13}C$ isotopologue [M+1] relative to $^{12}C$ isotopologue [M] on the most
175  abundant mass. *coeff[1-5]* are powers of the parameters to apply different magnitudes of each variable in
176  different studies. Using this score, a ranking for candidate formula was determined. By default, IDSL.UFA
177  utilized a value of 1 for *coeff[1-5]* to rank candidate molecular formulas in the equation (6). However, we
178  have provided a score coefficient optimization strategy in the section S.1 which can be helpful for
179  improving the ranking when larger size IPDB are utilized.

180  **Summary statistics of molecular formulas annotation in the aligned peak table:** It is quite common
181  to have more than 50 samples in metabolomics and exposomics projects, which can be leveraged to
182  compute a statistic for formula annotations across all the samples. For each peak (*m/z*-RT pair) in the

183 aligned peak table, corresponding molecular formula lists across all the samples were retrieved using the
184 peak indices provided by the IDSL.IPA data processing. We then aggregated these formula lists and
185 computed two properties 1) the detection frequency and 2) median rank for each formula assigned for a
186 peak across all the samples (individual peak list). Then we generated a new sort order for each molecular
187 formula at the aligned peak table level using the following formula: $\frac{\sqrt{frequency}}{median\ rank}$. For each peak in the
188 aligned peak table, top 5-20 formulas with detection frequency and median ranks across all samples were
189 exported to a csv file for each test study.

190 **Molecular formula class detection:** Many compounds belong to a chemical class with a distinct sub-
191 structure pattern such as polychlorinated biphenyl (PCBs), polybrominated diphenyl ethers (PBDEs),
192 polycyclic aromatic hydrocarbons (PAHs), perfluoroalkyl substances (PFAS), lipids and phthalates etc.
193 The formula annotations generated via the enumerated chemical space (ECS) approach were processed
194 to detect such classes within a list of formulas. The IDSL.UFA function '*detect_formula_sets*' was used to
195 detect 1) constant ΔH/ΔC ratios for polymeric (ΔH/ΔC = 2) and cyclic (ΔH/ΔC = 1/2) chain progressions
196 within polymeric and cyclic classes (Table S.2- S.4) and 2) a constant number of carbons and fixed
197 summation of hydrogens and halogens (Σ(H+Br+Cl+F+I)) representing classes similar to PCBs, PBDEs
198 (Table S.5).

199 **Correlation analysis for gestational age:** The ST001430 study[26] includes weekly blood samples of 30
200 pregnancies. The study has 781 total samples each processed in positive and negative modes to predict
201 gestational age. To reduce batch effects, the peak heights were adjusted by raw total ion chromatograms
202 (TICs) in each sample, and then the positive and negative aligned peak height tables were stacked to
203 generate a comprehensive list of peaks. We computed a Spearman correlation coefficient between
204 gestational age and peak height data for each pregnancy. A schematic of this workflow is presented in
205 Figure S.2.

206 **Results and discussion**
207 We have engineered a new software, IDSL.UFA, to annotate LC/HRMS peaks with molecular formulas for
208 an untargeted metabolomics or exposomics study. In this approach, IDSL.UFA computes theoretical
209 isotopic profiles for molecular formulas, matches theoretical isotopic profiles against experimental
210 LC/HRMS data in individual data file using a set of matching parameters and then summarizes the
211 formula annotations using detection frequency and median ranks in multiple samples (aligned annotated
212 peak table) in a study. The IDSL.UFA software has been implemented as an R package and made
213 publicly available via the R-CRAN repository and www.ufa.idsl.me site.

214 **Section 1) Development and validation of IDSL.UFA results:** To demonstrate the validity of our
215 approach to assign molecular formulas, we have utilized datasets with true positive annotations and show
216 their ranks in the IDSL.UFA result matrices.

217 **Analysis of authentic reference standards:** First, we evaluated performance of the IDSL.UFA software
218 to detect molecular formulas in LC/HRMS data for authentic reference standards. We found that the
219 average $\overline{NEME}$ (indicator of mass difference) was 0.70 mDa and $\overline{PCS}$ (indicator of isotope profile
220 similarity) were 99.968% between experimental and theoretical isotopic profiles for 367 authentic
221 standard compounds of common metabolites. This indicated that the observed isotopic profiles were very
222 similar to the theoretical counterparts for these reference standards and suggested that molecular
223 formulas can be reliably assigned to untargeted data generated by the commonly used LC/HRMS
224 instruments. The theoretical and experimental integrated isotopic profile spectra across chromatography
225 for these standards are provided at Zenodo repository accession (https://zenodo.org/record/5803968) and
226 an example compound (Kynurenine ion [$C_{10}H_{13}N_2O_3$]⁺) is shown in Figure 2 and Figure S.3 ($\overline{NEME}$ ≤ 0.61
227 mDa and $\overline{PCS}$ = 100.000%).

228 **Analysis of untargeted LC/HRMS data with structurally annotated peaks:**

229  We selected four publicly available studies (ST001154[27], ST001683[25], MTBLS1684[23], and
230  MTBLS2542[24]). These studies have reported annotations with MSI 1-3 confidence levels
231  (https://zenodo.org/record/5838709) that were obtained using retention time, accurate mass and MS/MS
232  spectra matching. For these studies, the IDSL.UFA software assigned 61.85%, 54.31%, 70.58% and
233  85.51% molecular formula as the top hit, and 96.90%, 90.58%, 100% and 99.29% molecular formulas in
234  the top five hits in the aligned table. These results were generated using the IPDB of the
235  IDSL.ExposomeDB with 209,592 and 129,122 ion formulas in positive and negative modes from multiple
236  ionization pathways, respectively representing 83,951 unique intact molecular formulas
237  (http://zenodo.org/deposit/5838709).

238  For each selected study, an ECS IPDB was generated using the element boundaries that covered the
239  formula list of true positive annotations for the study. When IDSL.UFA software was used for each study
240  using those specific ECS IPDBs, the assignment rates were – 52.74%, 53.36%, 79.41% and 51.08%
241  molecular formula as the top hit, and 95.60%, 84.45%, 100% and 91.66% molecular formula in the top 5
242  hits in the aligned table (Figure S.4). Generally, the IDSL.UFA software annotated 924 (90.14%) and 877
243  (85.56%) molecular formulas across all four studies using IDSL.ExposomeDB and ECS IPDBs,
244  respectively. These findings demonstrate that IDSL.UFA is a sensitive approach to cover the majority of
245  formulas for chemicals detectable in a biospecimen.

246  There is a tradeoff of coverage and the confidence in annotation while choosing chemical space for
247  molecular formula annotation. We have noticed that the rank of true positive hits degrades when we have
248  used a larger chemical space (Figure S.5). However, when compounds that are known and expected to
249  be found in a blood specimen are used, we have observed that formulas for true positives are often
250  ranked top hits. Therefore, we recommend a chemical prioritization strategy by sample type and to first
251  match the compounds that are expected for that sample type and then expand the chemical space to
252  cover additional peaks.

253  **Summary of the formula annotations in the aligned peak table:** Our raw data processing generates
254  both a separate list of m/z-RT pairs for each sample (individual peak list) and a single combine list
255  (aligned-table) of m/z-RT pairs for all samples. IDSL.UFA annotates molecular formulas only to individual
256  peak lists, then, it computes the detection frequency and median rank for all formulas annotated for the
257  same peak across all samples using the aligned peak table (See methods). Our hypothesis is that the
258  most probable formula of the underlying ionized compound will have a higher detection frequency and
259  median rank across all the samples. For example, for the MTLS1684 study, 24/35 (69%) of the reported
260  annotations had a median rank of 1 and 8/35 (23%) had a median rank of 2 across all 499 samples
261  (https://zenodo.org/record/5838709). We propose that the summary of detection frequencies and ranks
262  across individual data files can be helpful in boosting the confidence for formula assignments in multi-
263  sample studies. It should be noted that IDSL.UFA does not group related peaks to flag them as potential
264  ESI adducts or in-source fragments. Such grouping of peaks can be achieved by existing solutions such
265  as MS-FLO[18] online tool or CliqueMS[9] R package.

266  **Additional validation of molecular formula assignment by MS/MS:** To further ensure that IDSL.UFA
267  can assign high confidence molecular formulas for untargeted LC/HRMS data, we utilized data from
268  ST002044 study which has high quality MS/MS data collected in the data dependent mode. A total 73 hits
269  were confirmed by matching their spectra to the NIST 2020 MS/MS library (https://chemdata.nist.gov) and
270  public mass spectral libraries (https://zenodo.org/record/6416108). IDSL.UFA assigned 78.75% of these
271  hits within a median rank of ≤ 2 in the aligned peak table generated using the IDSL.ExposomeDB formula
272  IPDB (Table S.6 and Figure S.6). These results provided additional supports to confidence in the
273  molecular formula assignment by the IDSL.UFA software using the IDSL.ExposomDB IPDB.

274  **Rank score optimization:** IDSL.UFA utilized a number of chromatographic-mass spectrometry
275  parameters to compute the rank of a molecular formula for a peak in the individual peak list. By default, a
276  score coefficient of 1 is used which works sufficiently in most situations. However, the rank can be further
277  improved by an optimization strategy that utilizes the true positive, curated and high-quality structure

278  annotations for each data file as input. This can be achieved by running a mixture of reference standards
279  using the same analytical method or by annotating peaks using MS/MS, RT and isotopic profile matching
280  using stringent criteria. For metabolite standards (MSV000088661) and blood specimens (ST002044)
281  studies, we have observed a significant improvement in the ranking of molecular formulas when
282  optimized score coefficients were utilized in the IDSL.UFA software (Table S.7).

283  **Section 2) Application of IDSL.UFA for a pregnancy study**

284  To demonstrate an application of IDSL.UFA software to characterize the metabolome and exposome for
285  blood specimens, we have re-processed a publicly available study ST001430[26] (n=781) which has weekly
286  blood samples analyzed for 30 pregnancies to accurately predict gestational age (GA in weeks). Raw
287  data were processed using the IDSL.IPA software to generate the individual peak lists and the aligned
288  peak table (https://zenodo.org/record/5804527). On average, (3,416 ESI$^-$ and 6,978 ESI$^+$) peaks were
289  detected across individual peak lists for this study and a total of (89,174 ESI$^-$ and 143,712 ESI$^+$) peaks
290  were reported in the aligned peak table. The IDSL.UFA software using the IDSL.ExposomeDB IPDB
291  annotated (80,957 ESI$^-$ and 124,647 ESI$^+$) peaks in the aligned peak table with at least one molecular
292  formula having a median rank of ≤ 5.

293  We identify the peaks that were associated with GA by computing a spearman correlation coefficient
294  between normalized peak-height for each peak and GA. On a spearman cutoff of ($p$-value ≤ 0.05, $|\rho|$ ≥
295  0.65, "two.sided" alternative), 274 peaks with a detection frequency of ≥ 5 within each subject were found
296  to be significantly associated with GA (only ≤ 36 weeks). We observed 242 (red) and 32 (blue) ascending
297  and descending correlations patterns with GA, which were consistent with the patterns reported in the
298  original paper[26] and corresponded to chemicals related to steroid hormone biosynthesis and long-chain
299  fatty acids. These results show the potential the IDSL.UFA approach to characterize the pregnancy
300  related metabolic changes (Figure 3).

301  To flag the potential peaks related to chemical exposures in the pregnancy study (ST001430), we first
302  assigned a molecular formula using an ECS that may cover diverse halogenated compounds that were
303  not found in the IDSL.ExposomDB formula list. IDSL.UFA resulted with 199,837 unique molecular
304  formulas on the aligned table (top rank ≤ 30 and number of hits ≤ 30) in the ST001430 study. Grouping
305  these formulas by a class detection approach (see method) highlighted that 7,615, 18,452, and 32,107
306  distinct formula classes. For instance, a class of heavily halogenated compounds, $C_nHClF_{2n}O_4$ (n=10-12),
307  known as chlorinated perfluorotriether alcohols (Cl-PFTrEAs) was detected for human specimens in this
308  study. Cl-PFTrEAs was previously only reported in air samples from eastern China[37] and may represent a
309  new ubiquitous global contaminant class. IDSL.UFA can only confirm isotopic profiles match (Figure S.6);
310  however, a confirmatory in-source fragment ($[M-C_3F_6O]^-$) was  consistent with the published MS/MS
311  fragmentation (Figure S.8).[37] Authentic standards for Cl-PFTrEAs are not readily available; therefore a
312  confidence level 3b (isotopic profile match combined with fragmentation-based candidate) is suggested
313  for these annotations  according to a recently proposed PFAS identification confidence level by
314  Charbonnet et al.[38] Levels of Cl-PFTrEAs were similar to the commonly known legacy halogenated
315  compounds[14] for human serum samples (Figure 4). These findings also show that IDSL.UFA software can
316  potentially detect chemicals of public health concerns in a human biospecimen and can be helpful in
317  expanding the existing database of exposome chemicals.[39]

318  **Section 3) Performance benchmarking and comparison with existing tools**

319  IDSL.UFA processed one file (D115_NEG.mzml from the ST2044 study) in ~10 minutes on a computer
320  with 6 cores, indicating the pipeline can be used in normally available computing resources.

321  To check how IDSL.UFA performed for low abundant signals, we utilized data from the MTBLS1040 study
322  which has a seven-point calibration curve for the analyzed compounds. For the hippuric acid standard in
323  the MTBLS1040 study, IDSL.UFA correctly assigned the molecular formula to the corresponding peak in

324    samples analyzed at up to 8 fmol concentration level (second-lowest point)
325    (https://zenodo.org/record/6466668).

326    IDSL.UFA software is designed to cover commonly used LC-HRMS instruments for human biospecimens
327    studies in the EBI MetaboLights and Metabolomics Workbench repositories. A mass resolution of 20,000
328    and mass accuracy of 5 ppm is often found for these instruments. We compared the results for publicly
329    available two raw data files for a BioRec human plasma sample analyzed for a lipidomic assay by
330    QToF(ST001843) and Orbitrap instruments(ST001264) using the same chromatography method in the
331    same lab. Our workflow generated 1752 peaks with 2855 formulas for the QToF data file and 1328 peaks
332    with 2209 formulas for the Orbitrap data file. A list of 151 true positive annotations from the ST110054[27]
333    study (MS/MS matches were inspected by an expert user from the same lab and chromatography
334    method) was utilized for these test data files (https://zenodo.org/record/6621138). For these true
335    positives, 35% were found to be top hits in the QToF data file and 53% in the Orbitrap data file. It seems
336    our approach works slightly better for Orbitrap data. However, an even higher resolution and better mass
337    accuracy can be helpful in removing several false positive annotations, and in improving the ranking of
338    the true positive annotations.

339    When we imported a MS1 only data file in the SIRIUS[20] tool, it did not process the file, which was
340    expected since SIRIUS only processes data files with MS/MS spectra. For a data file (D115_NEG.mzml
341    from the ST002044 study) with MS/MS spectra in the Data Dependent Acquisition (DDA) mode, SIRIUS
342    processed 885 MS/MS spectra and suggested formula annotations for 221 spectra, whereas IDSL.UFA
343    assigned molecular formula to 9303 peaks in this data file.

344    IDSL.UFA natively uses IUPAC isotope table data[29] to calculate theoretical isotopic profiles and
345    calculated almost identical isotopic profiles to that obtained from the *envi*Pat package[40] (Table S.8).
346    Negligible mass and profile similarity  differences (NEME ≤ 0.69 mDa and PCS ≥ 99.999%) were
347    observed for formula $[C_8F_{17}O_3S]^-$ between IDSL.UFA and *envi*Pat[40].

348    Next, we compared the IDSL.UFA against Rdisop[19] R package to show the advantages of a database-
349    dependent approach (IDSL.UFA) over a database independent approach (Rdisop) for molecular formula
350    annotation. For kynurenine authentic standard (MSV000088661), both IDSL.UFA and Rdisop[19] ranked
351    the M+H adduct formula as the top hit(Section S.2 and Table S.9). But Rdisop's ranking for PFOS
352    isomers were >20 in the studies ST001430 and ST002044 (both human blood samples). Whereas
353    IDSL.UFA annotated both isomers of PFOS as top hit for these studies (Table S.10-11 and Figure S.9-
354    10). This suggests that Rdisop may miss important expected compounds when a complex chemical
355    space (CHBrClFNOPS) is targeted, but IDSL.UFA will be able to annotate them for human blood
356    specimens. Next, we extended the comparison to the lipidomics analysis with 151 true positive
357    annotations. Rdisop annotated 12%, whereas IDSL.UFA reported 53% of true annotations as top hits for
358    the Orbitrap data file (Figure S.11). These comparisons suggest that a database dependent approach for
359    formula annotation, such as IDSL.UFA should be used first to screen for expected compounds in HRMS
360    data before looking for unknown-unknowns. We also provide a comparison (Table S.12) between
361    IDSL.UFA and Rdisop[19] R packages, highlighting new features that IDSL.UFA is introducing into R
362    computing workflows for metabolomics and exposomics studies.

363    Our approach to obtain homologous series with polymeric chain increment from a list of input molecular
364    formulas is different from the prior approaches[41-43-40] in which molecular formulas are enumerated only for
365    a known series or chain increment rule. Therefore, our approach has the flexibility to discover new types
366    of homologous series among a collection of formulas.

367    Conclusion
368    IDSL.UFA enabled a comprehensive characterization of the chemical space that was detected by an
369    untargeted LC/HRMS assay to study the metabolome and exposome and its role in human health. The
370    unique feature of the IDSL.UFA software is to utilize the summary statistics for the rank and frequency of
371    detected molecular formulas in the aligned annotated molecular formula table. It can complement the

372     other peak annotation efforts that use mainly MS/MS data to annotate peaks, and thus lower the number
373     of false negative reporting of peaks and minimize the under-utilization of the untargeted LC/HRMS
374     datasets. We provided various scenarios to obtain molecular formulas from a known database and
375     enumeration strategies to assign a formula to peaks in a LC/HRMS dataset. These new computational
376     strategies for molecular formula assignment can greatly expand the quality of untargeted LC/HRMS data
377     matrices and their analyses especially when MS/MS data are not available.

380     **Author's contribution:** SFB and DKB planned the study, prepared the results and drafted the
381     manuscript. SFB and DKB coded the IDSL.UFA package. YK, SB and PC provided test LC/HRMS data
382     for authentic standards of metabolites and human biospecimens. All authors have reviewed the
383     manuscript content.
384

## References

(1) Needham, B. D.; Adame, M. D.; Serena, G.; Rose, D. R.; Preston, G. M.; Conrad, M. C.; Campbell, A. S.; Donabedian, D. H.; Fasano, A.; Ashwood, P.; et al. Plasma and Fecal Metabolite Profiles in Autism Spectrum Disorder. *Biol Psychiatry* **2021**, *89* (5), 451-462. DOI: 10.1016/j.biopsych.2020.09.025

(2) Gonzalez-Dominguez, R.; Jauregui, O.; Queipo-Ortuno, M. I.; Andres-Lacueva, C. Characterization of the Human Exposome by a Comprehensive and Quantitative Large-Scale Multianalyte Metabolomics Platform. *Anal Chem* **2020**, *92* (20), 13767-13775. DOI: 10.1021/acs.analchem.0c02008

(3) Shen, B.; Yi, X.; Sun, Y.; Bi, X.; Du, J.; Zhang, C.; Quan, S.; Zhang, F.; Sun, R.; Qian, L.; et al. Proteomic and Metabolomic Characterization of COVID-19 Patient Sera. *Cell* **2020**, *182* (1), 59-72 e15. DOI: 10.1016/j.cell.2020.05.032

(4) Wozniak, J. M.; Mills, R. H.; Olson, J.; Caldera, J. R.; Sepich-Poore, G. D.; Carrillo-Terrazas, M.; Tsai, C. M.; Vargas, F.; Knight, R.; Dorrestein, P. C.; et al. Mortality Risk Profiling of Staphylococcus aureus Bacteremia by Multi-omic Serum Analysis Reveals Early Predictive and Pathogenic Signatures. *Cell* **2020**, *182* (5), 1311-1327 e1314. DOI: 10.1016/j.cell.2020.07.040

(5) Wang, L. B.; Karpova, A.; Gritsenko, M. A.; Kyle, J. E.; Cao, S.; Li, Y.; Rykunov, D.; Colaprico, A.; Rothstein, J. H.; Hong, R.; et al. Proteogenomic and metabolomic characterization of human glioblastoma. *Cancer Cell* **2021**, *39* (4), 509-528 e520. DOI: 10.1016/j.ccell.2021.01.006

(6) Franzosa, E. A.; Sirota-Madi, A.; Avila-Pacheco, J.; Fornelos, N.; Haiser, H. J.; Reinker, S.; Vatanen, T.; Hall, A. B.; Mallick, H.; McIver, L. J.; et al. Gut microbiome structure and metabolic activity in inflammatory bowel disease. *Nat Microbiol* **2019**, *4* (2), 293-305. DOI: 10.1038/s41564-018-0306-4

(7) Shen, X.; Wang, R.; Xiong, X.; Yin, Y.; Cai, Y.; Ma, Z.; Liu, N.; Zhu, Z. J. Metabolic reaction network-based recursive metabolite annotation for untargeted metabolomics. *Nat Commun* **2019**, *10* (1), 1516. DOI: 10.1038/s41467-019-09550-x

(8) Wang, L.; Xing, X.; Chen, L.; Yang, L.; Su, X.; Rabitz, H.; Lu, W.; Rabinowitz, J. D. Peak Annotation and Verification Engine for Untargeted LC-MS Metabolomics. *Anal Chem* **2019**, *91* (3), 1838-1846. DOI: 10.1021/acs.analchem.8b03132

(9) Senan, O.; Aguilar-Mogas, A.; Navarro, M.; Capellades, J.; Noon, L.; Burks, D.; Yanes, O.; Guimera, R.; Sales-Pardo, M. CliqueMS: a computational tool for annotating in-source metabolite ions from LC-MS untargeted metabolomics data based on a coelution similarity network. *Bioinformatics* **2019**, *35* (20), 4089-4097. DOI: 10.1093/bioinformatics/btz207

(10) Uppal, K.; Walker, D. I.; Jones, D. P. xMSannotator: An R Package for Network-Based Annotation of High-Resolution Metabolomics Data. *Anal Chem* **2017**, *89* (2), 1063-1067. DOI: 10.1021/acs.analchem.6b01214

(11) Fakouri Baygi, S.; Fernando, S.; Hopke, P. K.; Holsen, T. M.; Crimmins, B. S. Automated Isotopic Profile Deconvolution for High Resolution Mass Spectrometric Data (APGC-QToF) from Biological Matrices. *Anal Chem* **2019**, *91* (24), 15509-15517. DOI: 10.1021/acs.analchem.9b03335

(12) Fakouri Baygi, S.; Crimmins, B. S.; Hopke, P. K.; Holsen, T. M. Comprehensive Emerging Chemical Discovery: Novel Polyfluorinated Compounds in Lake Michigan Trout. *Environ Sci Technol* **2016**, *50* (17), 9460-9468. DOI: 10.1021/acs.est.6b01349

(13) Fakouri Baygi, S.; Fernando, S.; Hopke, P. K.; Holsen, T. M.; Crimmins, B. S. Decadal Differences in Emerging Halogenated Contaminant Profiles in Great Lakes Top Predator Fish. *Environ Sci Technol* **2020**, *54* (22), 14352-14360. DOI: 10.1021/acs.est.0c03825

(14) Fakouri Baygi, S.; Fernando, S.; Hopke, P. K.; Holsen, T. M.; Crimmins, B. S. Nontargeted Discovery of Novel Contaminants in the Great Lakes Region: A Comparison of Fish Fillets and Fish Consumers. *Environ Sci Technol* **2021**, *55* (6), 3765-3774. DOI: 10.1021/acs.est.0c08507

(15) Fakouri Baygi, S.; Hutinet, S.; Cariou, R.; Fernando, S.; Hopke, P. K.; Holsen, T. M.; Crimmins, B. S. Comparison between Automated and User-Interactive Non-Targeted Screening Tools: Isotopic Profile Deconvoluted Chromatogram (IPDC) Algorithm and HaloSeeker. *Int J Environ Sci Technol* **2022**, *27*, 1-12.

(16) Li, Z.; Lu, Y.; Guo, Y.; Cao, H.; Wang, Q.; Shui, W. Comprehensive evaluation of untargeted metabolomics data processing software in feature detection, quantification and discriminating marker selection. *Anal Chim Acta* **2018**, *1029*, 50-57. DOI: 10.1016/j.aca.2018.05.001

(17) Tautenhahn, R.; Bottcher, C.; Neumann, S. Highly sensitive feature detection for high resolution LC/MS. *BMC Bioinformatics* **2008**, *9*, 504. DOI: 10.1186/1471-2105-9-504

(18) DeFelice, B. C.; Mehta, S. S.; Samra, S.; Cajka, T.; Wancewicz, B.; Fahrmann, J. F.; Fiehn, O. Mass Spectral Feature List Optimizer (MS-FLO): A Tool To Minimize False Positive Peak Reports in Untargeted

441 Liquid Chromatography-Mass Spectroscopy (LC-MS) Data Processing. *Anal Chem* **2017**, *89* (6), 3250-
442 3255. DOI: 10.1021/acs.analchem.6b04372
443 (19) Neumann, S.; Pervukhin, A.; Böcker, S. Mass decomposition with the Rdisop package. **2022**, 1.
444 (20) Duhrkop, K.; Fleischauer, M.; Ludwig, M.; Aksenov, A. A.; Melnik, A. V.; Meusel, M.; Dorrestein, P.
445 C.; Rousu, J.; Bocker, S. SIRIUS 4: a rapid tool for turning tandem mass spectra into metabolite structure
446 information. *Nat Methods* **2019**, *16* (4), 299-302. DOI: 10.1038/s41592-019-0344-8
447 (21) Ludwig, M.; Nothias, L. F.; Duhrkop, K.; Koester, I.; Fleischauer, M.; Hoffmann, M. A.; Petras, D.;
448 Vargas, F.; Morsy, M.; Aluwihare, L.; et al. Database-independent molecular formula annotation using
449 Gibbs sampling through ZODIAC. *Nat Mach Intell* **2020**, *2* (10), 629-+. DOI: 10.1038/s42256-020-00234-6
450 (22) Chen, L.; Lu, W.; Wang, L.; Xing, X.; Chen, Z.; Teng, X.; Zeng, X.; Muscarella, A. D.; Shen, Y.;
451 Cowan, A.; et al. Metabolite discovery through global annotation of untargeted metabolomics data. *Nat*
452 *Methods* **2021**, *18* (11), 1377-1385. DOI: 10.1038/s41592-021-01303-3
453 (23) Alfano, R.; Chadeau-Hyam, M.; Ghantous, A.; Keski-Rahkonen, P.; Chatzi, L.; Perez, A. E.; Herceg,
454 Z.; Kogevinas, M.; de Kok, T. M.; Nawrot, T. S.; et al. A multi-omic analysis of birthweight in newborn cord
455 blood reveals new underlying mechanisms related to cholesterol metabolism. *Metabolism* **2020**, *110*,
456 154292. DOI: 10.1016/j.metabol.2020.154292
457 (24) Wu, P.; Chen, D.; Ding, W.; Wu, P.; Hou, H.; Bai, Y.; Zhou, Y.; Li, K.; Xiang, S.; Liu, P.; et al. The
458 trans-omics landscape of COVID-19. *Nat Commun* **2021**, *12* (1), 4543. DOI: 10.1038/s41467-021-24482-
459 1
460 (25) Han, S.; Van Treuren, W.; Fischer, C. R.; Merrill, B. D.; DeFelice, B. C.; Sanchez, J. M.;
461 Higginbottom, S. K.; Guthrie, L.; Fall, L. A.; Dodd, D.; et al. A metabolomics pipeline for the mechanistic
462 interrogation of the gut microbiome. *Nature* **2021**, *595* (7867), 415-420. DOI: 10.1038/s41586-021-03707-
463 9
464 (26) Liang, L.; Rasmussen, M. H.; Piening, B.; Shen, X.; Chen, S.; Rost, H.; Snyder, J. K.; Tibshirani, R.;
465 Skotte, L.; Lee, N. C.; et al. Metabolic Dynamics and Prediction of Gestational Age and Time to Delivery
466 in Pregnant Women. *Cell* **2020**, *181* (7), 1680-1692 e1615. DOI: 10.1016/j.cell.2020.05.002
467 (27) Barupal, D. K.; Zhang, Y.; Shen, T.; Fan, S.; Roberts, B. S.; Fitzgerald, P.; Wancewicz, B.; Valdiviez,
468 L.; Wohlgemuth, G.; Byram, G.; et al. A Comprehensive Plasma Metabolomics Dataset for a Cohort of
469 Mouse Knockouts within the International Mouse Phenotyping Consortium. *Metabolites* **2019**, *9* (5). DOI:
470 10.3390/metabo9050101
471 (28) Chambers, M. C.; Maclean, B.; Burke, R.; Amodei, D.; Ruderman, D. L.; Neumann, S.; Gatto, L.;
472 Fischer, B.; Pratt, B.; Egertson, J.; et al. A cross-platform toolkit for mass spectrometry and proteomics.
473 *Nat Biotechnol* **2012**, *30* (10), 918-920. DOI: 10.1038/nbt.2377
474 (29) Kim, S.; Gindulyte, A.; Zhang, J.; Thiessen, P. A.; Bolton, E. E. PubChem Periodic Table and
475 Element Pages: Improving Access to Information on Chemical Elements from Authoritative Sources.
476 *Chem Teach Int* **2021**, *3* (1), 57-65. DOI: 10.1515/cti-2020-0006
477 (30) Currie, L. A. Nomenclature in Evaluation of Analytical Methods Including Detection and
478 Quantification Capabilities (Iupac Recommendations 1995). *Pure Appl Chem* **1995**, *67* (10), 1699-1723.
479 DOI: DOI 10.1351/pac199567101699
480 (31) Fahy, E.; Subramaniam, S. RefMet: a reference nomenclature for metabolomics. *Nat Methods* **2020**,
481 *17* (12), 1173-1174. DOI: 10.1038/s41592-020-01009-y
482 (32) Fahy, E.; Sud, M.; Cotter, D.; Subramaniam, S. LIPID MAPS online tools for lipid research. *Nucleic*
483 *Acids Res* **2007**, *35* (Web Server issue), W606-612. DOI: 10.1093/nar/gkm324
484 (33) *FDA structured product labelling*. https://www.fda.gov/industry/fda-resources-data-
485 standards/structured-product-labeling-resources (accessed.
486 (34) Williams, A. J.; Grulke, C. M.; Edwards, J.; McEachran, A. D.; Mansouri, K.; Baker, N. C.; Patlewicz,
487 G.; Shah, I.; Wambaugh, J. F.; Judson, R. S.; et al. The CompTox Chemistry Dashboard: a community
488 data resource for environmental chemistry. *J Cheminform* **2017**, *9* (1), 61. DOI: 10.1186/s13321-017-
489 0247-6
490 (35) Kind, T.; Fiehn, O. Seven Golden Rules for heuristic filtering of molecular formulas obtained by
491 accurate mass spectrometry. *BMC Bioinformatics* **2007**, *8*, 105. DOI: 10.1186/1471-2105-8-105
492 (36) Fakouri Baygi, S.; Kumar, Y.; Barupal, D. K. IDSL.IPA Characterizes the Organic Chemical Space in
493 Untargeted LC/HRMS Data Sets. *J Proteome Res* **2022**, *21* (6), 1485-1494. DOI:
494 10.1021/acs.jproteome.2c00120
495 (37) Yu, N.; Wen, H.; Wang, X.; Yamazaki, E.; Taniyasu, S.; Yamashita, N.; Yu, H.; Wei, S. Nontarget
496 Discovery of Per- and Polyfluoroalkyl Substances in Atmospheric Particulate Matter and Gaseous Phase

497  Using Cryogenic Air Sampler. *Environ Sci Technol* **2020**, *54* (6), 3103-3113. DOI:
498  10.1021/acs.est.9b05457
499  (38) Charbonnet, J. A.; McDonough, C. A.; Xiao, F.; Schwichtenberg, T.; Cao, D.; Kaserzon, S.; Thomas,
500  K. V.; Dewapriya, P.; Place, B. J.; Schymanski, E. L.; et al. Communicating Confidence of Per- and
501  Polyfluoroalkyl Substance Identification via High-Resolution Mass Spectrometry. *Environmental Science*
502  *& Technology Letters* **2022**. DOI: 10.1021/acs.estlett.2c00206
503  (39) Barupal, D. K.; Fiehn, O. Generating the Blood Exposome Database Using a Comprehensive Text
504  Mining and Database Fusion Approach. *Environ Health Perspect* **2019**, *127* (9), 97008. DOI:
505  10.1289/EHP4713
506  (40) Loos, M.; Gerber, C.; Corona, F.; Hollender, J.; Singer, H. Accelerated isotope fine structure
507  calculation using pruned transition trees. *Anal Chem* **2015**, *87* (11), 5738-5744. DOI:
508  10.1021/acs.analchem.5b00941
509  (41) Schum, S. K.; Brown, L. E.; Mazzoleni, L. R. MFAssignR: Molecular formula assignment software for
510  ultrahigh resolution mass spectrometry analysis of environmental complex mixtures. *Environ Res* **2020**,
511  *191*, 110114. DOI: 10.1016/j.envres.2020.110114
512  (42) Hughey, C. A.; Hendrickson, C. L.; Rodgers, R. P.; Marshall, A. G.; Qian, K. Kendrick mass defect
513  spectrum: a compact visual analysis for ultrahigh-resolution broadband mass spectra. *Anal Chem* **2001**,
514  *73* (19), 4676-4681. DOI: 10.1021/ac010560w
515  (43) Jacob, P.; Barzen-Hanson, K. A.; Helbling, D. E. Target and Nontarget Analysis of Per- and
516  Polyfluoralkyl Substances in Wastewater from Electronics Fabrication Facilities. *Environ Sci Technol*
517  **2021**, *55* (4), 2346-2356. DOI: 10.1021/acs.est.0c06690
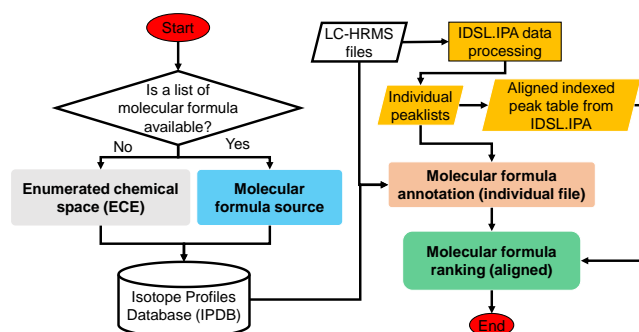
518

519

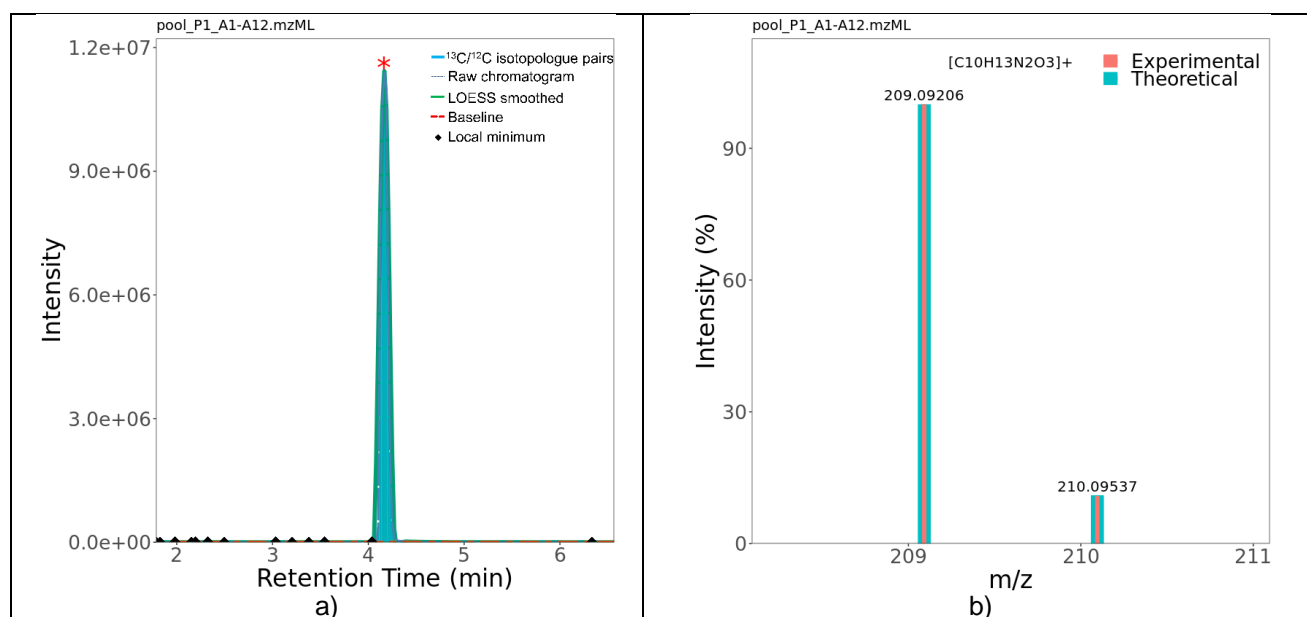520    **Figure 1.** A simplified flowchart of the IDSL.UFA software.

**Figure 2.** a) A chromatographic peak generated by the IDSL.IPA pipeline for Kynurenine ion ($[C_{10}H_{13}N_2O_3]^+$ = $[M+H]^+$) to detect peak boundaries. b) Comparison between the theoretical isotopic profile and integrated spectra across the chromatographic peak after molecular formula annotation using IDSL.UFA.
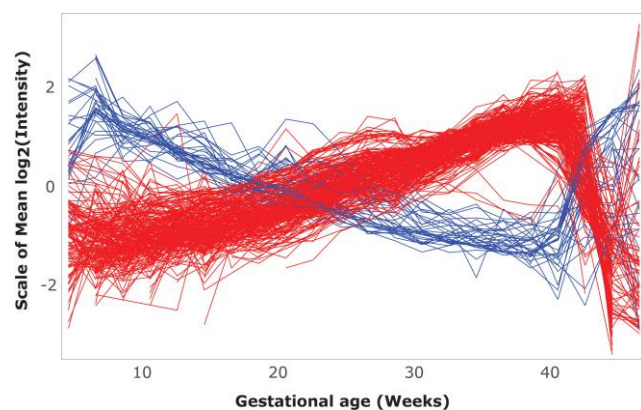
521

522

**Figure 3.** Trends of 274 peaks associated with pregnancy dynamics. Molecular formula annotated interactive plots are available at https://ufa.idsl.me/st001430 for an enumerated chemical space and IDSL.ExposomeDB IPDBs.
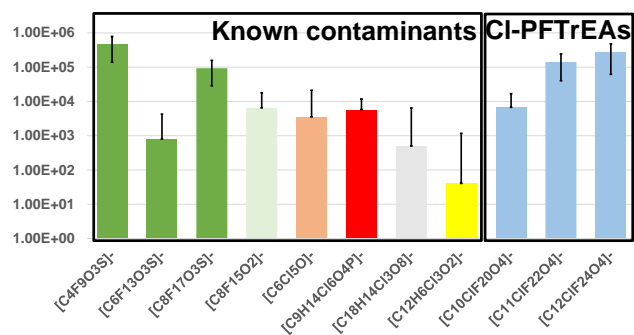
**Figure 4.** Peak area of halogenated contaminants in human blood ($[C_nF_{2n+1}O_3S]^-$ (n = 4, 6, 8), $[C_8F_{15}O_2]^-$, $[C_6Cl_5O]^-$, $[C_9H_{14}Cl_6O_4P]^-$, $[C_{18}H_{14}Cl_3O_8]^-$, $[C_{12}H_6Cl_3O_2]^-$) and Cl-PFTrEAs $[C_nClF_{2n}O_4]^-$ (n = 10-12) across 781 negative samples in the ST001430 study.