# Journal Pre-proofs

Data sharing in PredRet for accurate prediction of retention time: application to plant food bioactive compounds

Dorrain Yanwen Low, Pierre Micheau, Ville Mikael Koistinen, Kati Hanhineva, László Abrankó, Ana Rodriguez-Mateos, Andreia Bento da Silva, Christof van Poucke, Conceiç ão Almeida, Cristina Andres-Lacueva, Dilip K. Rai, Esra Capanoglu, Francisco A. Tomás Barberán, Fulvio Mattivi, Gesine Schmidt, Gözde Gürdeniz, Kateřina Valentová, Letizia Bresciani, Lucie Petrásková, Lars Ove Dragsted, Mark Philo, Marynka Ulaszewska, Pedro Mena, Raúl González-Domínguez, Rocío Garcia-Villalba, Senem Kamiloglu, Sonia de Pascual-Teresa, Stéphanie Durand, Wieslaw Wiczkowski, Maria Rosário Bronze, Jan Stanstrup, Claudine Manach

Please cite this article as: Yanwen Low, D., Micheau, P., Mikael Koistinen, V., Hanhineva, K., Abrankó, L., Rodriguez-Mateos, A., Bento da Silva, A., van Poucke, C., Almeida, C., Andres-Lacueva, C., Rai, D.K., Capanoglu, E., Tomás Barberán, F.A., Mattivi, F., Schmidt, G., Gürdeniz, G., Valentová, K., Bresciani, L., Petrásková, L., Ove Dragsted, L., Philo, M., Ulaszewska, M., Mena, P., González-Domínguez, R., Garcia-Villalba, R., Kamiloglu, S., de Pascual-Teresa, S., Durand, S., Wiczkowski, W., Rosário Bronze, M., Stanstrup, J., Manach, C., Data sharing in PredRet for accurate prediction of retention time: application to plant food bioactive compounds, *Food Chemistry* (2021), doi: https://doi.org/10.1016/j.foodchem.2021.129757

1  **Data sharing in PredRet for accurate prediction of retention time: application to plant food**

2  **bioactive compounds**

3  Dorrain Yanwen Low[1,29+], Pierre Micheau[1], Ville Mikael Koistinen[2,3], Kati Hanhineva[2,3], László

4  Abrankó[4], Ana Rodriguez-Mateos[5], Andreia Bento da Silva[6,7], Christof van Poucke[8], Conceição

5  Almeida[6], Cristina Andres-Lacueva[9,10], Dilip K. Rai[11], Esra Capanoglu[12], Francisco A. Tomás

6  Barberán[13,14], Fulvio Mattivi[15,16], Gesine Schmidt[17], Gözde Gürdeniz[18], Kateřina Valentová[19],

7  Letizia Bresciani[20], Lucie Petrásková[19], Lars Ove Dragsted[18], Mark Philo[21], Marynka

8  Ulaszewska[15,30], Pedro Mena[22], Raúl González-Domínguez[9,10], Rocío Garcia-Villalba[13], Senem

9  Kamiloglu[12,23], Sonia de Pascual-Teresa[24], Stéphanie Durand[25], Wieslaw Wiczkowski[26], Maria

10  Rosário Bronze[6,27,28], Jan Stanstrup[18]*, Claudine Manach[1]*

11  [1]Université Clermont Auvergne, INRAE, UNH, F-63000 Clermont Ferrand, France.

12  [2]Institute of Public Health and Clinical Nutrition, University of Eastern Finland, FI-70211 Kuopio, Finland.

13  [3]Department of Biochemistry, University of Turku, FI-20014 Turun yliopisto, Finland.

14  [4]Faculty of Food Science, Department of Applied Chemistry, Szent István Egyetem, 29-43 Villanyi street, 1118 Budapest, Hungary.

15  [5]Department of Nutritional Sciences, School of Life Course Sciences, King's College London, SE1 9NH London, United Kingdom.

16  [6]Instituto de Tecnologia Química e Biológica, Universidade Nova de Lisboa, Av. da República, 2780-157 Oeiras, Portugal.

17  [7]Faculty of Pharmacy, University of Lisbon, Avenida Professor Gama Pinto, 1649-003 Lisbon, Portugal.

18  [8]Technology and Food Science Department, Flanders research institute for Agriculture Fisheries and Food (ILVO), Brusselsesteenweg 370, B-9090

19  Melle, Belgium.

20  [9]Nutrition, Food Science and Gastronomy Department, Pharmacy faculty, University of Barcelona, Av Joan XXIII, 08028 Barcelona, Spain.

21  [10]CIBER Fragilidad y Envejecimiento Saludable (CIBERfes), Instituto de Salud Carlos III, 28029 Madrid, Spain.

22  [11]Department of Food BioSciences, Teagasc Food Research Centre Ashtown, Dublin, D15 KN3K, Ireland.

23  [12]Faculty of Chemical and Metallurgical Engineering, Department of Food Engineering, Istanbul Technical University, 34469, Maslak, Istanbul,

24  Turkey.

25  [13]Department of Food Science and Technology, CEBAS-CSIC, Campus Universitario de Espinardo, edf 25, 30100 Murcia, Spain.

26  [14]Department of Biotechnology, College of Science, Taif University, Taif 26571, Saudi Arabia.

27  [15]Department of Food Quality and Nutrition, Metabolomics Unit, Research and Innovation Centre, Fondazione Edmund Mach, 38010 San Michele

28  all'Adige, Italy.

29  [16]Department of Cellular, Computational and Integrative Biology, CIBIO, University of Trento, 38123 Trento, Italy.

30  [17]Department of Food and Health, Nofima, Norwegian Institute of Food, Fisheries and Aquaculture Research, PB 210, NO-1433 Ås, Norway.

31  [18]Department of Nutrition, Exercise and Sports, University of Copenhagen, Frederiksberg, DK-1985, Denmark.

32  [19]Laboratory of Biotransformation, Institute of Microbiology of the CAS, Vídeňská 1083, CZ-142 20 Prague, Czechia.

33  [20]Human Nutrition Unit, Department of Veterinary Science, University of Parma, Via Volturno, 39, 43125 Parma PR, Italy.

34  [21]Quadram Institute Biosciences, Norwich Research Park, NR4 7 UQ, United Kingdom.

35  [22]Human Nutrition Unit, Department of Food and Drug, University of Pharma, Via Volturno, 39, 43125 Parma PR, Italy.

36  [23]Science and Technology Application and Research Center (BITUAM), Bursa Uludag University, 16059 Gorukle, Bursa, Turkey.

37  [24]Department of Metabolism and Nutrition, Institute of Food Science, Technology and Nutrition (ICTAN-CSIC), Jose Antonio Novais 10, 28040

38  Madrid, Spain.

39  [25]Université Clermont Auvergne, INRAE, UNH, Plateforme d'Exploration du Métabolisme, MetaboHUB Clermont, F-63000 Clermont Ferrand,

40  France.

41  [26]Institute of Animal Reproduction and Food Research of the Polish Academy of Sciences, Tuwima 10, 10-748 Olsztyn, Poland.

42  [27]Instituto de Biologia Experimental Tecnológica, Av. da República, Quinta do Marquês, Edificio iBET/ITQB, 2780-157 Oeiras, Portugal.

43  [28]Research Institute for Medicines, Faculty of Pharmacy, University of Lisbon, Avenida Professor Gama Pinto, 1649-003 Lisbon, Portugal.

44  [29]Present address: Lee Kong Chian School of Medicine, Nanyang Techological University, 59 Nanyang Drive, Singapore 636921.

45  [30]Present address: Proteomics and Metabolomics Facility, Center for Omics Sciences, IRCCS San Raffaele Scientific Institute, Via Olgettina n.60

46  20132 Milan, Italy.

47  *Claudine Manach and Jan Stanstrup share last co-authorship.

48  +Corresponding author: Dorrain Low; Lee Kong Chian School of Medicine, Nanyang Technological University, 59 Nanyang Drive, Singapore 636921;

49  email: dorrain.low@ntu.edu.sg

50

51  **Abstract**

52  Prediction of retention times (RTs) is increasingly considered in untargeted metabolomics to

53  complement MS/MS matching for annotation of unidentified peaks. We tested the performance

54  of PredRet (http://predret.org/) to predict RTs for plant food bioactive metabolites in a data

55  sharing initiative containing entry sets of 29–103 compounds (totalling 467 compounds, >30

56  families) across 24 chromatographic systems (CSs). Between 27 and 667 predictions were

57  obtained with a median prediction error of 0.03–0.76 min and interval width of 0.33–8.78 min.

58  An external validation test of eight CSs showed high prediction accuracy. RT prediction was

59  dependent on shape and type of LC gradient, and number of commonly measured compounds.

60  Our study highlights PredRet's accuracy and ability to transpose RT data acquired from one CS to

61  another CS. We recommend extensive RT data sharing in PredRet by the community interested

62  in plant food bioactive metabolites to achieve a powerful community-driven open-access tool for

63  metabolomics annotation.

64

65  **Keywords**

66  Predicted retention time, metabolomics, plant food bioactive compounds, metabolites, data

67  sharing, PredRet

68

69  **1. Introduction**

70  The dark matter in metabolomics refers to the large fraction of molecular signals that are

71  detected with untargeted analyses but remain unidentified. Part of this dark matter corresponds

72  to the food metabolome. Currently, >26,000 compounds have been described in foods

73  (https://foodb.ca), and upon ingestion and digestion, these food components are further

74  transformed into various metabolites (Scalbert et al., 2014), many of which are not identified or

75  inventoried yet in databases (Barabási, Menichetti, & Loscalzo, 2020). Plant food bioactive

76  compounds (also referred as dietary phytochemicals, e.g., (poly)phenols, carotenoids,

77  glucosinolates, alkaloids) and their phase I, -II and gut microbial metabolites represent an

78  important class of the food metabolome that receive widespread interest for their protective

79  health effects and more recently, for their usefulness as food intake biomarkers. They cover a

80    large chemical space ranging from highly polar to lipophilic compounds, and their identification

81    in untargeted methods remains a challenging feat.

82

83    Identification of unknowns in untargeted metabolomics combines multiple types of information

84    and tools, such as matching of exact mass in compound databases, comparison of experimental

85    to reference $MS^n$ spectral data and chromatographic retention time (RT) of authentic standards

86    for Metabolomics Standards Initiative (MSI) level I identification, or to publicly available spectral

87    databases for MSI level II (Sumner et al., 2007). But searches in databases often return an

88    excessive number of structurally similar hypotheses (Hall, Hill, Cawley, Hall, Chen, & Grant, 2018),

89    and purchasing all corresponding standards is not feasible due to limited availability and high cost.

90    In the case of plant food bioactive compounds and their metabolites, identification is further

91    challenged by the lack of commercial standards and the high structural similarity between many

92    isomeric compounds, which makes their MS/MS spectra indistinguishable.

93

94    Leveraging orthogonal data such as RT becomes valuable for assisting the certainty of

95    identification to MSI levels I and II, by narrowing the number of plausible hypotheses within an

96    observed RT window. Recent years have seen several approaches to adopt RT prediction models

97    for integration into untargeted analysis workflows with varying degrees of success (McEachran,

98    Mansouri, Newton, Beverly, Sobus, & Wiliiams, 2018; Witting & Bocker, 2020). Existing types of

99    RT prediction models include i) simple algorithms based on log P or gradient back-calculation

100   (Abate-Pella et al., 2015; Boswell, Schellenberg, Carr, Cohen, & Hegeman, 2011), ii) monotonically

101   constrained generalised additive model (GAM) (Stanstrup, Neumann, & Vrhovšek, 2015) of

102   retention times and iii) complex *in silico* quantitative structure-retention relationship (QSSR)

103   models based on combinations of molecular descriptors. QSSRs can be built using different

104   machine learning approaches, such as artificial neural network, random forest and support vector

105   regression models (Aalizadeh, Nika, & Thomaidis, 2019; Bade, Bijlsma, Miller, Barron, Sancho, &

106   Hernández, 2015; Bouwmeester, Martens, & Degroeve, 2020; Domingo-Almenara et al., 2019;

107   Hall et al., 2018; McEachran et al., 2018; Naylor, Catrow, Maschek, & Cox, 2019; Tada et al., 2019;

108   Wolfer, Lozano, Umbdenstock, Croixmarie, Arrault, & Vayer, 2015). However, these prediction

109   models are limited in their application, as RT data are specific to one chromatographic system

110   (CS) and the models do not provide accurate predictions outside the trained conditions.

111

112   As analytical methods are not harmonised and most laboratories tend to have their own routine

113   semi-targeted or untargeted LC methods for covering plant food bioactive compounds in various

114   types of matrices (serum, plasma, urine, digestive fluids, food materials), it is ideal that RT

115   prediction models be customisable across CSs. PredRet (Stanstrup et al., 2015) represents an

116   original approach that enables users of the scientific community to benefit from RT data sharing

117   through its open access RT database, and obtain predictions in their own CS if the RT of a

118   compound has been experimentally determined by another user or laboratory. In this aspect,

119   PredRet is relevantly applicable for transposing RTs between CSs differing in mobile phase

120   composition, gradient, flow rate and column dimensions. In the framework of the COST Action

121   POSITIVe (https://www6.inra.fr/cost-positive, FA1403), we evaluated the performance of

122   PredRet to predict the RTs of plant food bioactive compounds and their metabolites in a multi-

123    laboratory test involving 19 laboratories across Europe, using 24 gradient-based reversed-phase

124    CSs. We also expanded PredRet database with experimental RTs of 467 plant food compounds.

125

126    **2. Experimental Section**

127    **2.1. Chemical compounds**

128    All participating laboratories purchased their own chemicals, differing from one laboratory to

129    another, except that 10 laboratories previously involved in a multiplatform coverage test

130    organised by the COST Action POSITIVe, received two common standard mixtures comprising of

131    56 plant food bioactive compounds (Koistinen et al., 2018). Synthesised standards ($n$ = 49) were

132    accepted in addition to commercial standards, provided that the structure was unambiguously

133    elucidated by NMR and MS/MS spectra and that the compounds are entered in the online

134    platform for food compound exchange, FoodComEx (https://foodcomex.org/). Depending on

135    laboratories, chemicals were analysed in solvent or spiked in biological matrices (urine or plasma).

136    A full list of the 467 analysed compounds is provided in Table S1, with their common name, InChI,

137    IDs in HMDB, FooDB and PhytoHub, taxonomy, chemical structure, formula, monoisotopic mass,

138    predicted logP and the number of CSs where they were analysed.

139

140    Experimental RT datasets containing compound name, InChI and/or chemical structure were

141    provided by the involved laboratories. InChIs were used as unambiguous identifiers for

142    recognition of identical compounds between CSs and compound names were harmonised across

143    laboratories. For polyphenol metabolites, we applied the new reference KCC nomenclature (Kay

144    et al., 2020). InChIs were either extracted from databases such as PhytoHub

145    (http://phytohub.eu), PubChem (https://pubchem.ncbi.nlm.nih.gov) (Kim et al., 2019), HMDB

146    v4.0 (www.hmdb.ca) (Wishart et al., 2018) or computed from chemical structures using Marvin

147    v19.7, 2019, ChemAxon (https://www.chemaxon.com). LogP values were computed using

148    ALOGPS v2.1 (http://www.vcclab.org/lab/alogps/) (Tetko et al., 2005; VCCLAB, 2005) after

149    conversion of InChIs to SMILES via InChIToSMILES (http://www.chemspider.com/inchi.asmx)

150    (Pence & Williams, 2010). In PredRet database, the main InChI layer containing chemical formula,

151    atom connections and hydrogen atom sublayers is considered when matching compounds, and

152    information after the main layer (e.g., charge, stereochemical and isotopic layers) is ignored.

153

154    **2.2 Chromatographic systems**

155    Experimental RT data were collected from 24 CSs across 19 laboratories. These CSs were not

156    intentionally optimised for the RT prediction test but rather represent the routine semi-targeted

157    or untargeted metabolomic methods of the various laboratories. A full description of instrument,

158    column and analytical conditions used in the 24 CSs is provided in Table 1. Overall, 15 C18 reverse-

159    phase (RP) columns from various manufacturers were used with dimensions ranging from 0.5 to

160    4.6 mm (internal diameter), 50 to 250 mm (length) and 1.6 to 5 μm (particle size). HPLC or UHPLC

161    methods were used in acidic conditions. Water and acetonitrile acidified with formic acid (0.1-

162    0.9%) or trifluoroacetic acid (0.1%) were most commonly used as mobile phases A and B, while

163    three CSs used methanol or acetone as mobile phase B. The gradients utilised in 13 UHPLC

164    methods consisted of linear and multiphasic slopes with flow rates of 0.4 to 0.6 mL/min and total

165    run times ranging from 6 to 26 min. There were four HPLC methods with multiphasic slopes with

166    flow rates of 0.015 to 1.5 mL/min and longer run times of 20 to 135 min. Figure S1 shows the

167    diversity of gradient slopes in the 24 CS.

168

169    **2.3 Prediction of retention times**

170    Experimentally measured RTs (Table S3) were entered in PredRet for the 467 compounds listed

171    in Table S1. The number of measured RTs by CS varied from 29 in CS9 to 103 in CS14. For each CS,

172    the compound names, InChI and experimentally measured RTs were entered into PredRet web

173    interface (http://predret.org) along with a description of the respective CS method. PredRet is

174    then able to predict RTs for compounds that have not been previously experimentally measured

175    in one CS but have been determined in some other CS. The prediction is achieved by constructing

176    GAMs between all pairs of CSs in the PredRet database using the compounds that were measured

177    in both CSs. Empirical prediction intervals (PI) were established via bootstrapping of GAMs, as

178    described in more details by Stanstrup et al (2015). The model providing the prediction with

179    narrowest PI was then used. Predictions were flagged as suspicious by the program if the RT is

180    considered potentially incorrect, when the difference between experimental and predicted RTs

181    was ≥ twice the distance from the predicted RT to outer limits of the PI. Predictions were

182    automatically discarded if their PI widths were ≥ 2 min or ≥ 20% of the predicted RT. The total

183    number of RT predictions between CSs, as well as accuracy and coverage of PI relative to the total

184    chromatographic run time, were compared.

185

186    **2.4. Validation of predicted retention times**

187  A validation test was conducted on CSs 1, 2, 4, 5, 14, 18, 19 and 22, which had the highest number

188  of experimental RT values. These eight CSs comprise of UHPLC and UPLC methods varying in LC

189  instrument and gradient, column, mobile phases, flow rate and run time. The experimental RT

190  datasets of these CSs were split into training sets (80% data, $n$ = 79, 71, 73, 67, 82, 63, 63 and 78

191  compounds respectively) and test sets (20% data, $n$ = 20, 18, 18, 17, 21, 16, 16 and 19 compounds

192  respectively). For selection of compounds in the test sets, the datasets were split into three equal

193  sections covering the beginning, middle and end of the chromatographic run, and then 20% of

194  the compounds were randomly selected from the three sections to ensure a uniform distribution

195  of RT along the entire chromatographic run. Another criterion was to select, in the test set, the

196  same proportion of unique compounds as in the whole dataset of the selected CSs. Validation of

197  RT predictions for each of the eight selected CSs was performed in conditions where the complete

198  datasets of the remaining 23 CSs were entered into the PredRet database.

199

200  **3. Results and Discussion**

201  **3.1. Large diversity of plant food metabolites analysed**

202  A total of 1583 experimental RT values were collected for 467 plant food compounds or related

203  human metabolites in one or several of the 24 CSs used by the 19 participating platforms. The

204  467 compounds belong to >30 families including flavonoids (anthocyanins, flavonols, flavones,

205  flavanols, flavanones, isoflavones), phenolic acids, lignans, ellagitannins, coumarins and

206  furanocoumarins, nitrogen-containing compounds (i.e., alkaloids, amines, indoles),

207  glucosinolates, alkylresorcinols, thiosulfinates, tocopherols, phytosterols, carotenoids and mono,

208  di-, sesqui- and triterpenoids, and their human metabolites, e.g. glucuronidated and sulfated

209 conjugates, as well as gut microbial metabolites. They cover a large chemical space from highly

210 polar to lipophilic with predicted logP values from –3.48 to 10.40 and with monoisotopic masses

211 from 95.0371 to 934.0712 daltons (Figure 1). The PredRet database is growing continuously with

212 addition of new compounds and associated RT data by registered users. At the time of our

213 experiment, a limited number of plant food compounds was present in PredRet, and our datasets

214 represented a major update for this category of compounds.

215

216 The number of CSs in which each compound was analysed is provided in Table S1. Of the 467

217 entered compounds, 212 were analysed in one CS only, while 4'-hydroxy-3'-methoxycinnamic

218 (ferulic), 4-hydroxy-3-methoxybenzoic (vanillic), 3,4-dihydroxybenzoic (protocatechuic), 5-*O*-

219 caffeoylquinic and 4'-hydroxycinnamic (*p*-coumaric) acids were most commonly measured in 20

220 of the 24 CSs (Figure S2). The size of the datasets varied from 29 to 103 experimental RTs. CSs 1,

221 2, 4, 5, 7, 14, 17, 18, 19, 22 and 23 contained ≥ 75 RTs, as illustrated by their large node size in

222 Figure 2, in contrast to CS9 and CS16, which contained the least RT data (29 and 35 RTs,

223 respectively). Across the platforms, CSs 2, 6, 11, 13 and 15 shared the highest compound overlap

224 as evidenced by their highly connected nodes (Figure 2) while still showing relatively good overlap

225 with CSs 1, 3, 7, 14, 16, 22 and 23. Pairwise clusters of CSs 18-19 and 4-5 were observed as they

226 shared > 90% compounds similarity, corresponding to two analytical methods from the same

227 platform.

228

229 **3.2. Retention time prediction coverage and rate**

230  A total of 6382 new RT predictions were obtained for the 24 CSs, with up to 667 predictions for

231  one CS (Table 2 and Figure S3). Compounds that were entered in PredRet prior to this study (1783

232  unique compounds, ~10% were plant food bioactive compounds) contributed to prediction of

233  additional compounds beyond the 467 compounds entered in this study. We observed a general

234  trend that as more experimental RTs are entered in PredRet, more RT predictions are generated

235  for compounds not previously analysed. This is demonstrated in CSs 1, 2, 22 and 23 where 559,

236  539, 667 and 572 new RT predictions were generated from 98, 89, 97 and 75 compounds entered

237  into PredRet respectively (Table 2). However, RT prediction was also dependent on shape (Figure

238  S1) and type (i.e., UHPLC or HPLC) of the LC gradient as well as number of common compounds

239  shared with other CSs. For example, infrequently used mobile phases may limit the predictability

240  of a CS. The entry of 29 compounds for CS9 was not sufficient to obtain RT predictions. However,

241  despite relatively small RT datasets (35 to 46 compounds) were entered for CSs 11, 15 and 16,

242  they had a high prediction rate, explained by a versatile CS and/or good combination of

243  compounds.

244

245  **3.3. Retention time prediction accuracy**

246  PredRet provided RT predictions for compounds never analysed in the CSs but also for compounds

247  in the entry dataset. We used the latter to compare prediction accuracy between CSs. RT

248  predictions were highly accurate across the 24 CSs, with median prediction errors between 0.03

249  and 0.76 min (Table 2). As run times vary greatly across CSs (5 to 135 min), median prediction

250  errors were also expressed in percentage relative to the total runtime, ranging from 0.3% to 1.8%

251  (CS9 excluded).

252

253    A graph comparing experimental and predicted RTs for compounds of CS1 entry dataset is given

254    in Figure 3 as an example. Equivalent graphs for all other CSs are provided in Figure S4. In CS1,

255    accurate predictions with narrow PI were obtained for most compounds with RT ranging between

256    6.6 and 14.2 min. Predictions for eight compounds (myo-inositol, proline betaine, dopamine,

257    3,4,5-trihydroxybenzoic acid (gallic acid), 1,3-dimethyluric acid, $\alpha$-tocopherol, ursolic acid and

258    alkylresorcinol C17:0) were discarded by PredRet algorithm as their PI widths were $\geq$ 2 min or >

259    20% of the predicted RT. PI width is an important indicator of prediction accuracy as it represents

260    how accurate the projection models are, based on the number of experimentally known RTs in

261    the RT range of compounds that are being projected in the pairwise CS models (Stanstrup et al.,

262    2015). We observed that predictions were usually missing at the beginning and end of the runs,

263    where there tends to be a low density of known RTs, and conditions are approaching the

264    analytical limits of the CSs (Figure S5). In CS1, predictions were not generated before the first 1.5

265    min and after 14 min. For 15 compounds (1-methylpiperidine, arbutin, 1-methylxanthine, 1*H*-

266    pyrrole-2-carboxaldehyde, cyclo(Leu-Pro), 5-(3',4'-dihydroxyphenyl)valeric acid,

267    homoeriodictyol, tomatidine, formononetin, bergapten, nobiletin, isosakuranetin, kaempferide,

268    biochanin A and bergamottin), RT prediction was not expected, as they were not present in any

269    other CS. Globally, PredRet performed well for CS1 with a median prediction error of 0.07 min

270    (0.27% of runtime) and median PI width of 0.83 min. For 77 non-unique compounds entered into

271    PredRet, 559 new predictions for compounds never analysed in this system were obtained, in the

272    range of 1.01 to 14.18 min.

273

274    Amongst CS7, CS8, CS9, CS14, CS17, CS20 and CS24, a common trait is the high proportion of rare

275    plant compounds unique to their CSs, which indirectly resulted in a low number of common

276    compounds shared with other CSs. For example, CS7 contained sesquiterpenoids not represented

277    in other CSs, likewise for anthocyanin glycosides in CS8, urolithins and conjugated isoflavone

278    metabolites in CS14, glucosinolates and rare flavonoids in CS17, flavonolignans (e.g.,

279    dehydrosilydianin), rare flavonoids and sulfated conjugates in CS20, and urolithins and

280    conjugated flavonoids in CS24. Adding the RTs of these rare plant food compounds and

281    metabolites contributes to the richness of PredRet database; however, a caveat is that these CSs

282    themselves may not receive the benefit of good prediction coverage. In such circumstances, the

283    user is encouraged to include common plant compounds that are also frequently represented in

284    the PredRet database. As an example, we propose a list of 14 compounds frequently analysed in

285    our study ($\geq$ 67% of 24 CSs), which covers a wide RT range: 4'-hydroxy-3'-methoxycinnamic acid

286    (ferulic acid), 4-hydroxy-3-methoxybenzoic acid (vanillic acid), 5-$O$-caffeoylquinic acid, 4'-

287    hydroxycinnamic acid ($p$-coumaric acid), 3,4-dihydroxybenzoic acid (protocatechuic acid), 3',4'-

288    dihydroxycinnamic acid (caffeic acid), 3,4,5-trihydroxybenzoic acid (gallic acid), 3',5'-dimethoxy-

289    4'-hydroxycinnamic acid (sinapic acid), (–)-epicatechin, kaempferol, hippuric acid, luteolin,

290    phloretin and hesperetin (Table S4).

291

292    To further validate the predictive performance of GAM in PredRet, we performed an external

293    validation test on a subset of eight CSs, splitting the experimental datasets into 80% for training

294    sets and 20% for test sets. The training sets were used to build GAMs between CSs in PredRet

295    database to obtain predictions with PIs for the compounds in the test sets. Predictions were

296 compared to experimental data to obtain the prediction error for each compound (Table S5) and

297 the prediction statistics for each CS are provided in Table S6. Accurate predictions were achieved,

298 with the median prediction error in the test sets ranging from 0.04 to 0.41 min across the eight

299 CSs. The maximum absolute prediction error was 3.55 min for $\alpha$-tocopherol (CS2), followed by

300 catechol (2.45 min, CS5). It is difficult to compare the performance of PredRet with other RT

301 prediction tools as those only allow predictions within the same CS, while PredRet predicts RTs

302 from one CS to another CS differing in mobile phase composition, gradient and flow rate.

303

304 Despite accurate models being built for the CSs, we observed that early- and late-eluting

305 compounds were generally omitted from predictions, likely due to their extreme polarity.

306 Compounds unique to respective CSs (e.g., nobiletin in CS1, 2',5'-dihydroxyphenylacetic acid in

307 CS2 and 9-hydroxy-urolithin-3-glucuronide (isourolithin A glucuronide) in CS14) did not obtain RT

308 predictions as well as compounds that did not have sufficient RT data density in the RT area (e.g.,

309 *N*-(3-hydroxybenzoyl)glycine (2-furoylglycine) in CS19 and pinoresinol in CS22). Between 14 and

310 39% of the compounds in the CSs of the validation (test) set had experimental RTs that fall outside

311 the estimated PI, showing that the PIs should be interpreted with caution, as previously noted in

312 the original paper describing PredRet. The practical implication is that a proposed annotation for

313 an experimental RT cannot be completely discarded even if the RT falls outside the proposed

314 annotation's PI. A few limitations of PredRet were identified in our study that may be corrected

315 in the future. Firstly, users have no information about the standards that have been considered

316 for providing predictions in their CS: e.g., commercial or synthesised standard, analysed in solvent

317 or spiked in a biological matrix. Secondly, PredRet algorithm recognises the entered compounds

318    based on the main InChI layer only and therefore stereochemical information is ignored during

319    RT prediction.

320

321    **3.4. Application of PredRet predictions for identification of plant food compounds in**

322    **metabolomic studies**

323    The effectiveness of RT prediction using PredRet allowed the distinction of isomeric compounds.

324    In Figure 4A, 3-(3',4'-dihydroxyphenyl)propanoic (dihydrocaffeic) acid with a predicted RT of 8.3

325    min (PI: 8.1 to 8.5 min) could be distinguished from its isomers, 4'-hydroxy-3'-

326    methoxyphenylacetic acid (homovanillic) acid (PI: 8.5 to 9.0 min) and 3,4-dimethoxybenzoic acid

327    (veratric) acid (PI: 9.3 to 9.6 min). In Figure 4B, the predicted RTs of fisetin (PI: 9.8 to 10.7 min),

328    kaempferol (PI: 11.4 to 12 min) and luteolin (PI: 10.6 to 11.3 min) were also clearly distinguished,

329    except for the narrow overlap in the PIs of luteolin and fisetin (10.6 to 10.7 min). This is

330    particularly useful as an orthogonal parameter to eliminate hypotheses when identifying

331    unknown features with the same $m/z$ in untargeted metabolomics studies. In addition, as RTs of

332    flavonoid conjugates (glycosides, glucuronides) differ from that of their aglycones, prediction of

333    RT may help to distinguish between aglycones truly present in the samples and detected

334    aglycones that are generated during the analysis as in-source fragments of glycosides or

335    glucuronides.

336

337    Another useful application of PredRet is aiding in annotation of rare plant food compounds in

338    untargeted metabolomics studies, when the standards are not commercially available or difficult

339    to synthesise. As soon as a user enters experimental data for a rare plant food compound in a

340   CSs, PredRet provides RT prediction with PI for this compound in CSs where it has not been

341   experimentally measured. For example, the contribution of tomatidine's experimental RT (11.8

342   min) from CS1 enabled the prediction of RTs in 15 other CSs, while formononetin (CS1), 8-

343   hydroxy-urolithin-3-sulfate (CS14) and 8-deoxylactucin (CS7) enabled the prediction of RTs in 13

344   other CSs. To optimise this process, it is crucial that users who entered experimental RTs for rare

345   compounds also enter experimental RTs for common compounds such as those suggested above.

346

347   **4. Conclusion**

348   PredRet, based on pairwise GAMs, was demonstrated to be a useful tool for obtaining a good

349   number and highly accurate RT predictions for plant food bioactive compounds and their

350   metabolites. Its use in untargeted metabolomics studies can definitely help for tentative

351   identification, by eliminating hypotheses that do not fall within the predicted RT range, or when

352   commercial standards are not readily available. PredRet predictions are precise enough to

353   distinguish structural isomers. Our data sharing initiative and multi-laboratory study contributed

354   to the expansion of the PredRet database with > 1500 experimental RTs in 24 CSs for > 467 plant

355   food bioactive compounds and their metabolites (> 30 families). Importantly, as more

356   experimentally known RTs are entered, more RT predictions are generated and accuracy of the

357   predictions increases. The PredRet database has grown considerably since its introduction and

358   now contains 15,000 RT entries across 68 CSs. Overall, the database covers 4,000 unique

359   compounds, beyond plant food bioactive compounds. In comparison, spectral libraries such as

360   the MassBank of North America (MoNa), contain mass spectra for >200,000 compounds, so there

361   remains a large potential for RT sharing. If sufficiently developed to allow accurate RT prediction

362 in any CSs, PredRet would facilitate comparisons between-studies and minimise the need to

363 develop a consensus LC–MS method for plant food compounds. We thus invite the scientific

364 community to contribute to the community-driven open access PredRet database as part of the

365 global effort for annotation of the dark matter of metabolomes. We suggest that sharing of RT as

366 well as collisional cross section data should be as commonplace in the future as sharing of MS/MS

367 data to provide enough orthogonal data for unambiguous identification in metabolomics.

368

399

400    **Author contributions:** Experimental data were provided by all authors, and data analysis was

401    conducted by D.Y.L, P.M, J.S and C.M. Manuscript was drafted by D.Y.L and C.M, and reviewed by

402    all authors.

403

404    **Conflict of interest:** Kati Hanhineva and Ville Koistinen are affiliated with Afekta Technologies,

405    Ltd. All other authors declare no conflict of interest.

406

407 **Supporting information:** Further methods and analysis details; Table S1: 467 analysed

408 compounds; Table S2: 24CS methods; Table S3: 24CS experimental RTs; Table S4: frequently

409 measured compounds; Table S5: validation test results; Table S6: validation test summary; Figure

410 S1: %B gradient; Figure S2: RT data density in 24CS; Figure S3: predictions per CS; Figure S4: CS1-

411 24 RT graphs; Figure S5: pairwise GAM graphs.

412  **References**

413  Aalizadeh, R., Nika, M.-C., & Thomaidis, N. S. (2019). Development and application of retention
414       time prediction models in the suspect and non-target screening of emerging
415       contaminants. *Journal of Hazardous Materials, 363*, 277-285.
416       https://doi.org/10.1016/j.jhazmat.2018.09.047.
417  Abate-Pella, D., Freund, D. M., Ma, Y., Simón-Manso, Y., Hollender, J., Broeckling, C. D., . . .
418       Boswell, P. G. (2015). Retention projection enables accurate calculation of liquid
419       chromatographic retention times across labs and methods. *Journal of Chromatography A,*
420       *1412*, 43-51. https://doi.org/10.1016/j.chroma.2015.07.108.
421  Bade, R., Bijlsma, L., Miller, T. H., Barron, L. P., Sancho, J. V., & Hernández, F. (2015). Suspect
422       screening of large numbers of emerging contaminants in environmental waters using
423       artificial neural networks for chromatographic retention time prediction and high
424       resolution mass spectrometry data analysis. *Science of the Total Environment, 538*, 934-
425       041.
426  Barabási, A.-L., Menichetti, G., & Loscalzo, J. (2020). The unmapped chemical complexity of our
427       diet. *Nature Food, 1*, 33-37. https://doi.org/10.1038/s43016-019-0005-1.
428  Boswell, P. G., Schellenberg, J. R., Carr, P. W., Cohen, J. D., & Hegeman, A. D. (2011). Easy and
429       accurate high-performance liquid chromatography retention prediction with different
430       gradients, flow rates, and instruments by back-calculation of gradient and flow rate
431       profiles. *Journal of Chromatography A, 1218*(38), 6742-6749.
432       https://doi.org/10.1016/j.chroma.2011.07.070.
433  Bouwmeester, R., Martens, L., & Degroeve, S. (2020). Generalized calibration across liquid
434       chromatography setups for generic prediction of small-molecule retention times.
435       *Analytical Chemistry, 92*, 6571-6578. https://doi.org/10.1021/acs.analchem.0c00233.
436  Domingo-Almenara, X., Guijas, C., Billings, E., Montenegro-Burke, R., Uritboonthai, W., Aisporna,
437       A., . . . Siuzdak, G. (2019). The METLIN small molecule dataset for machine learning-based
438       retention time prediction. *Nature Communications, 10*(5811), 1-9.
439       https://doi.org/10.1038/s41467-019-13680-7.
440  Hall, M. L., Hill, D. W., Cawley, S., Hall, L. H., Chen, M. H., & Grant, D. F. (2018). Development of a
441       reverse phase HPLC retention index model for nontargeted metabolomics using synthetic
442       compounds. *Journal of Chemical Information and Modeling, 58*, 591-604.
443  Kay, C. D., Clifford, M. N., Mena, P., McDougall, G. J., Andres-Lacueva, C., Cassidy, A., . . . Crozier,
444       A. (2020). Recommendations for standardizing nomenclature for dietary (poly)phenol
445       catabolites. *The American Journal of Clinical Nutrition, 112*(4), 1051-1068.
446       https://doi.org/10.1093/ajcn/nqaa204.
447  Kim, S., Chen, J., Cheng, T., Gindulyte, A., He, J., He, S., . . . Bolton, E. E. (2019). PubChem 2019
448       update: improved access to chemical data. *Nucleic Acids Research, 47*, 1102-1109.
449       https://doi.org/10.1093/nar/gky1033.
450  Koistinen, V. M., da Silva, A. B., Abrankó, L., Low, D., Villalba, R. G., Tomás-Barberán, F., . . . Bronze,
451       M. R. (2018). Interlaboratory coverage test on plant food bioactive compounds and their
452       metabolites by mass spectrometry-based untargeted metabolomics. *Metabolites, 8*(46),
453       1-17. https://doi.org/10.3390/metabo8030046.

454    McEachran, A. D., Mansouri, K., Newton, S. R., Beverly, B. E. J., Sobus, J. R., & Wiliiams, A. J. (2018).
455        A comparison of three liquid chromatography (LC) retention time prediction models.
456        *Talanta, 182*, 371-379. https://doi.org/10.1016/j.talanta.2018.01.022.

457    Naylor, B., Catrow, L., Maschek, A., & Cox, J. (2019). QSRR automator: A tool for automating
458        retention time prediction in lipidomics and metabolomics. *Metabolites, 10*(237), 1-15.
459        https://doi.org/10.3390/metabo10060237.

460    Pence, H. E., & Williams, A. (2010). ChemSpider: An online chemical information resource. *Journal*
461        *of Chemical Education, 87*(11), 1123-1124. https://doi.org/doi.org/10.1021/ed100697w.

462    Scalbert, A., Brennan, L., Manach, C., Andres-Lacueva, C., Dragsted, L. O., Draper, J., . . . Wishart,
463        D. S. (2014). The food metabolome: a window over dietary exposure. *American Journal of*
464        *Clinical Nutrition, 99*(6), 1286-1308.

465    Stanstrup, J., Neumann, S., & Vrhovšek, U. k. (2015). PredRet: Prediction of retention time by
466        direct mapping between multiple chromatographic systems. *Analytical Chemistry, 87*(18),
467        9421-9428. https://doi.org/10.1021/acs.analchem.5b02287.

468    Sumner, L., Amberg, A., Barrett, D., Beale, M. H., beger, R., Daykein, C. A., . . . Viant, M. R. (2007).
469        Proposed minimum reporting standards for chemical analysis Chemical Analysis Working
470        Group (CAWG) Metabolomics Standards Initiative (MSI). *Metabolomics, 3*(3), 211-221.
471        https://doi.org/10.1007/s11306-007-0082-2.

472    Tada, I., Tsugawa, H., Meister, I., Zhang, P., Shu, R., Katsumi, R., . . . Chaleckis, R. (2019). Creating
473        a reliable mass spectral-retention time library for all ion fragmentation-based
474        metabolomics. *Metabolites, 9*(251), 1-15. https://doi.org/10.3390/metabo9110251.

475    Tetko, I. V., Gasteiger, J., Todeschini, R., Mauri, A., Livingstone, D., Ertl, P., . . . Prokopenko, V. V.
476        (2005). Virtual computational chemistry laboratory - design and description. *Journal of*
477        *Computer-Aided Molecular Design, 19*, 453-363. https://doi.org/10.1007/s10822-005-
478        8694-y.

479    VCCLAB.    (2005).    Virtual    Computational    Chemistry    Laboratory.    Retrieved    from:
480        http://www.vcclab.org Accessed.

481    Wishart, D. S., Feunang, Y. D., Marcu, A., Guo, A. C., Liang, K., Vázquez-Fresno, R., . . . Scalbert, A.
482        (2018). HMDB 4.0: the human metabolome database for 2018. *Nucleic Acids Research, 46*,
483        608-617. https://doi.org/10.1093/nar/gkx1089.

484    Witting, M., & Bocker, S. (2020). Current status of retention time prediction in metabolite
485        identification.    *Journal    of    Separation    Science,    43*,    1746-1754.
486        https://doi.org/10.1002/jssc.202000060.

487    Wolfer, A. M., Lozano, S., Umbdenstock, T., Croixmarie, V., Arrault, A., & Vayer, P. (2015). UPLC-
488        MS retention time prediction: a machine learning approach to metabolite identification in
489        untargeted profiling. *Metabolites, 12*(8), 1-13.

490

491 **Figure captions**

492

493 Figure 1. Chemical space covered by the 467 plant food metabolites entered in PredRet.

494

495 Figure 2. Network map illustrating compound coverage overlap. Size of node represents the

496 number of compounds present in the dataset while the thickness and colour of edges represent

497 the number of common compounds between the paired datasets. The thicker the edge, the larger

498 number of common compounds. Edge colours    and    denote low (< 10) and high (> 60) similarity

499 of compounds, respectively. E: number of RT data entered into PredRet; P: number of new RT

500 predictions made.

501

502 Figure 3. Retention time (RT) prediction accuracy and coverage of 98 compounds with

503 experimentally known RTs in Chromatographic System (CS) 1. Refer to Supporting Information

504 Tables 3 and 4 for more details of individual compounds.

505

506 Figure 4. Retention time (RT) prediction for isomers A) 3,4-dimethoxybenzoic (veratric acid), 4'-

507 hydroxy-3'-methoxyphenylacetic (homovanillic) acid and 3-(3',4'-dihydroxyphenyl)propanoic

508 (dihydrocaffeic) acid, and B) kaempferol, luteolin and fisetin in CS1. Coloured areas represent the

509 prediction interval width.

510 **Table captions**

511

512 Table 1. Instrument and conditions of chromatographic systems used by participating platforms.

513

514 Table 2. Statistics of PredRet retention time predictions for 24 liquid chromatographic systems

515 (CSs) with an entry dataset of 467 plant compounds.

516 Highlights
517 • Identifying food bioactive compounds in untargeted metabolomics is challenging.
518 • Predicted retention time is valuable towards effort in metabolite identification.
519 • 24 Chromatographic systems obtained predicted retention times from PredRet database.
520 • High accuracy and coverage of retention time predictions for new compounds obtained.
521 • We recommend extensive retention time data sharing in open access PredRet database.
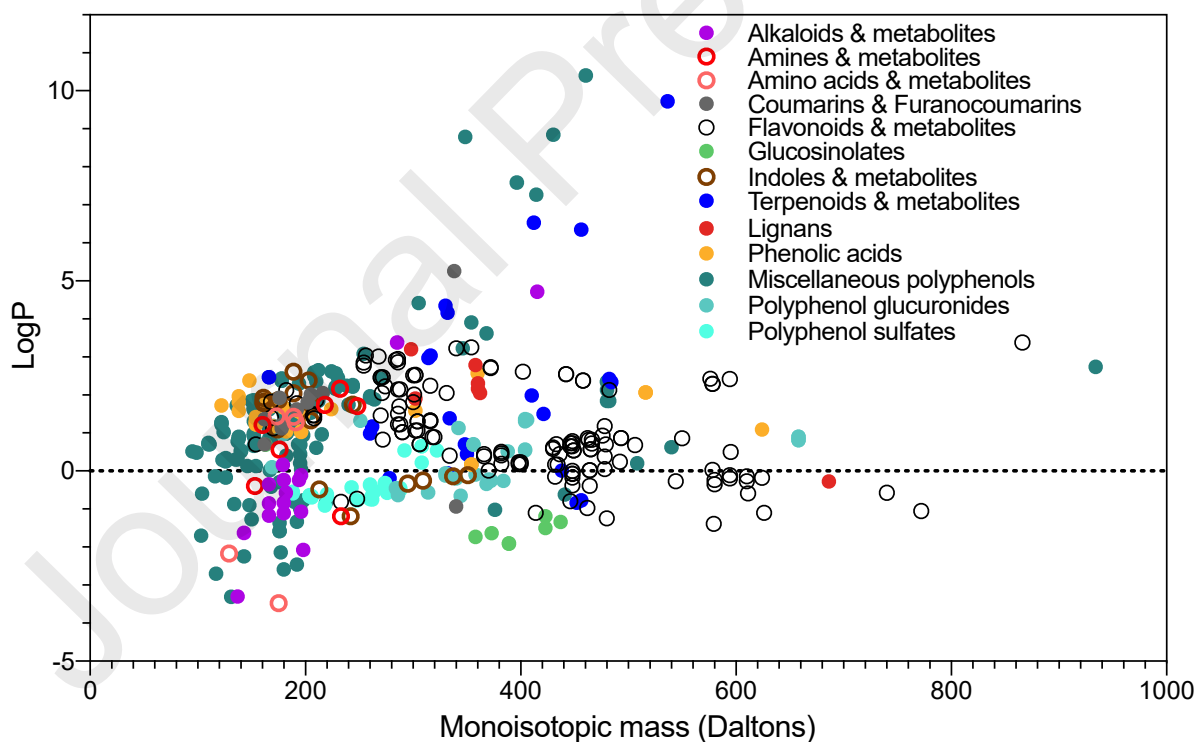522



Figure 1. Chemical space covered by the 467 plant food metabolites entered in PredRet.

523
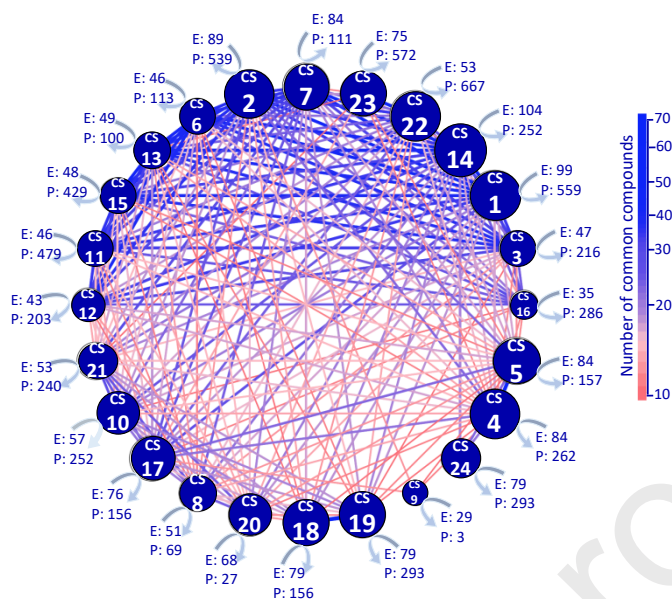524 (Color should be used for figure in print)
525
526

Figure 2. Network map illustrating compound coverage overlap. Size of node represents the number of compounds present in the dataset while the thickness and colour of edges represent the number of common compounds between the paired datasets. The thicker the edge, the larger number of common compounds. Edge colours    and    denote low (<10) and high (>60) similarity of compounds respectively. E: number of RT data entered into PredRet; P: number of new RT predictions made.

527
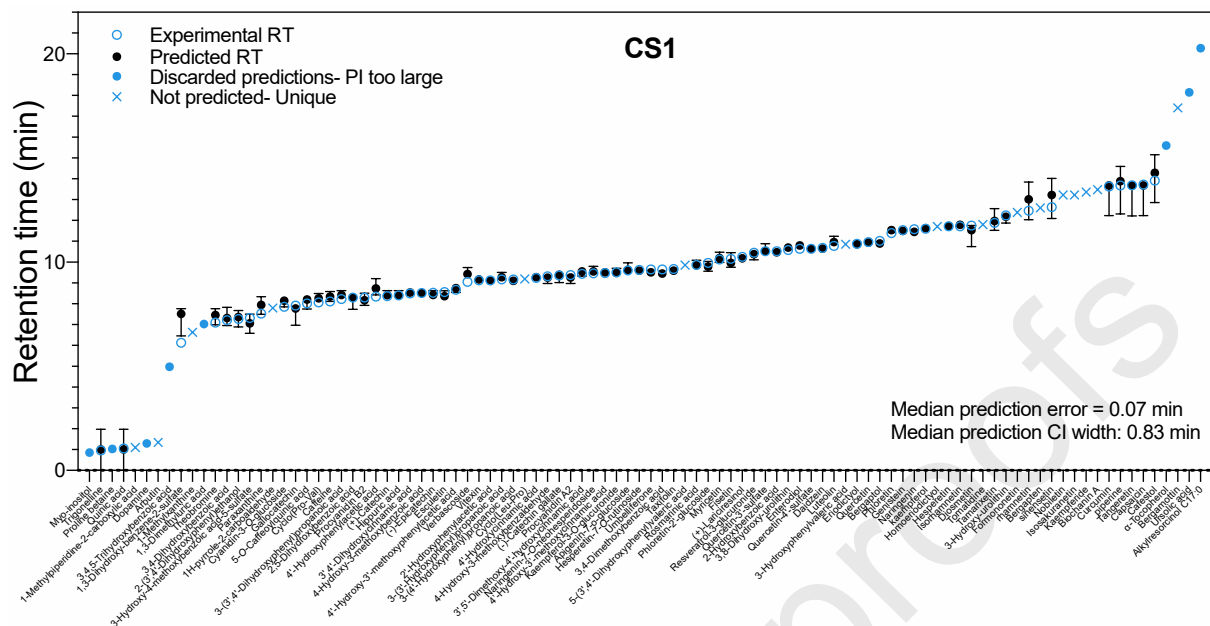528    (Color should be used for figure in print)
529
530

Figure 3. Retention time (RT) prediction accuracy and coverage of 98 compounds with experimentally known RTs in Chromatographic System (CS) 1. Refer to supporting information Tables 3 and 4 for more details of individual compounds.

531
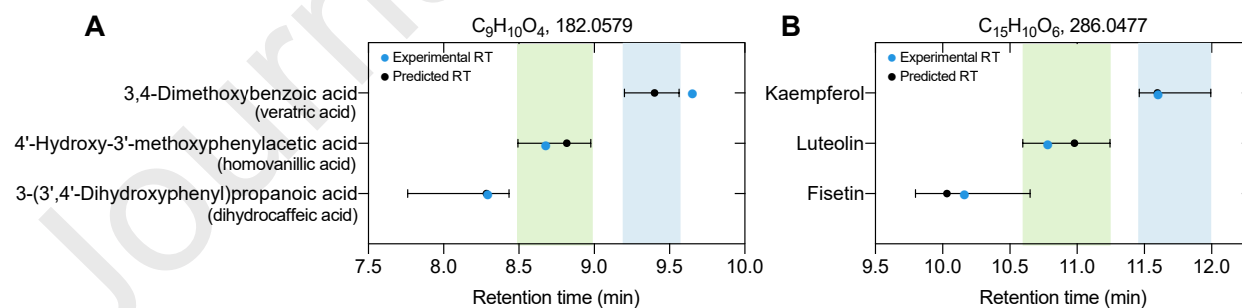532    (Color should be used for figure in print)
533
534



Figure 4. Retention time (RT) prediction for isomers A) 3,4-dimethoxybenzoic (veratric acid), 4'-hydroxy-3'-methoxyphenylacetic (homovanillic) acid and 3-(3',4'-dihydroxyphenyl)propanoic

(dihydrocaffeic) acid, and B) kaempferol, luteolin and fisetin in CS1. Coloured areas represent

the prediction interval width.

535
536    (Color should be used for figure in print)
537
538