

Manuscript version: Author's Accepted Manuscript

The version presented in WRAP is the author's accepted manuscript and may differ from the published version or Version of Record.

Persistent WRAP URL:

<http://wrap.warwick.ac.uk/172971>

How to cite:

Please refer to published version for the most recent bibliographic citation information. If a published version is known of, the repository item page linked to above, will contain details on accessing it.

Copyright and reuse:

The Warwick Research Archive Portal (WRAP) makes this work by researchers of the University of Warwick available open access under the following conditions.

© 2023 Elsevier. Licensed under the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International <http://creativecommons.org/licenses/by-nc-nd/4.0/>.



Publisher's statement:

Please refer to the repository item page, publisher's statement section, for further information.

For more information, please contact the WRAP Team at: wrap@warwick.ac.uk.

1 **Analytical Chemistry Solutions to Hazard Evaluation of**
2 **Petroleum Refining Products**

3
4 Alina T. Roman-Hubers,¹ Alexandra C. Cordova,¹ Mark P. Barrow,² Ivan Rusyn^{1,*}

5
6 ¹Interdisciplinary Faculty of Toxicology and Department of Veterinary Physiology and
7 Pharmacology, Texas A&M University, College Station, Texas, USA

8 ²Department of Chemistry, University of Warwick, Coventry, CV4 7AL, United Kingdom

9
10 *Corresponding author: Ivan Rusyn, MD, PhD, Veterinary Physiology and Pharmacology, 4458
11 TAMU, Texas A&M University, College Station, TX 77843, USA; Tel: +1-979-458-9866; email:
12 irusyn@tamu.edu

13
14
15 **Abbreviations:**

16 API, American Petroleum Institute; ASTM, American Society for Testing and Methods; CAS
17 Chemical Abstract Service; CCS, collision cross section; DBE, double bond equivalents; ECHA,
18 European Chemical Agency; EINECS, European Inventory of Existing Commercial Chemical
19 Substances; US EPA, United States Environmental Protection Agency; FID, Flame ionization
20 detection; FT-ICR MS, Fourier transform ion cyclotron resonance mass spectrometry; FWHM, full
21 width at half maximum; GC, gas chromatography; GC×GC, two-dimensional gas chromatography;
22 GC-MS, gas chromatography coupled to mass spectrometry; HPLC, high-performance liquid
23 chromatography; IMMS, ion mobility mass spectrometry; IR, infrared; KMD, Kendrick mass
24 defect; LC, liquid chromatography; MS, mass spectrometry; NMR, nuclear magnetic resonance;
25 REACH, Registration, Evaluation, Authorisation and Restriction of Chemicals; TIMS, trapped ion
26 mobility spectrometry; TOF MS, time-of-flight mass spectrometry; UV, ultraviolet; UVCBs,
27 substances of unknown, variable composition, complex reaction products, or biological materials.

28 **Abstract**

29 Products of petroleum refining are substances that are both complex and variable. These
30 substances are produced and distributed in high volumes; therefore, they are heavily scrutinized in
31 terms of their potential hazards and risks. Because of inherent compositional complexity and
32 variability, unique challenges exist in terms of their registration and evaluation. Continued
33 dialogue between the industry and the decision-makers has revolved around the most appropriate
34 approach to fill data gaps and ensure safe use of these substances. One of the challenging topics
35 has been the extent of chemical compositional characterization of products of petroleum refining
36 that may be necessary for substance identification and hazard evaluation. There are several novel
37 analytical methods that can be used for comprehensive characterization of petroleum substances
38 and identification of most abundant constituents. However, translation of the advances in
39 analytical chemistry to regulatory decision-making has not been as evident. Therefore, this
40 review's goal is to bridge the divide between the science of chemical characterization of petroleum
41 and the needs and expectations of the decision-makers. Collectively, mutual appreciation of the
42 regulatory guidance and the realities of what information these new methods can deliver should
43 facilitate the path forward in ensuring safety of the products of petroleum refining.

44 **Introduction**

45 Crude oils are naturally occurring and highly complex substances which vary considerably
46 in molecular composition according to their origins; they comprise a myriad of constituents,
47 primarily hydrocarbons, but also other organic and inorganic molecules (Smith et al., 1959). Close
48 to one hundred billion barrels of crude oils are annually extracted, distributed, and processed into a
49 wide variety of refined petroleum products (Kaiser, 2017; Salvito et al., 2020). The chemical
50 composition of petroleum refining products therefore depends on both the type (*i.e.*, origins) of
51 crude oil from which it was derived, and the refining process (*i.e.*, fractional distillation and/or
52 cracking followed by additional processing through solvent extraction, hydro-desulfurization, or
53 hydrogenation) used to meet performance characteristics of the end-products (McKee et al., 2015).
54 Products of petroleum refining are high production volume substances and thus are heavily
55 scrutinized in terms of their potential human and environmental health hazards and risks. Because
56 of inherent compositional complexity and variability, petroleum substances are prototypical
57 representatives of a diverse class known as substances of unknown, variable composition, complex
58 reaction products, or biological materials (UVCBs); these substances present unique challenges to
59 regulatory agencies, especially in terms of characterization of their chemical composition (Clark et
60 al., 2013; ECHA, 2017c; Lai et al., 2022). It is worth noting that the different fields within the
61 academic, industrial, and regulatory science communities can use differing terminology. Analytical
62 researchers often refer to petroleum refining products as “hydrocarbon mixtures” which are, in turn,
63 part of the broader “complex mixtures” family of samples. The term “mixture” is avoided by the
64 industry and decision-makers as they reason that most substances in commerce that are made from
65 oil are products of refining, rather than mixing, and thus to differentiate from the mixtures found in
66 the environment, such substances are called petroleum UVCBs.

67 Studies of molecular composition of crude oils and petroleum refining products have a long
68 history spanning over 80 years (**Figure 1**). The analytical characterization of petroleum substances
69 historically tracked the physico-chemical properties that pertained to the functionality of the
70 product, such as flash points and vapor pressure. With the advent of spectroscopy and mass
71 spectrometry techniques, there came the possibility to gradually gain more detailed understanding
72 of composition; however, the granularity of information on the constituents in registered petroleum
73 products is still lacking. Recent improvements have been made in the resolution of mass
74 spectrometers, ionization methods in order to access a wider range of components, separation

75 methods to offer structural insights, and developments in data visualization through standardized
76 diagrams (Palacio Lozano et al., 2020). As a result of these advances, new opportunities emerged
77 to provide comprehensive characterization of these complex substances and satisfy regulatory needs
78 on the composition, quantity of potentially hazardous constituents, and the extent of variability
79 among manufacturing batches of these products. The range and types of mass spectrometry
80 techniques that can be used for the analysis of petroleum-related samples is quite extensive;
81 collectively, the methods for study of petroleum are now being referred to as “petroleomics,” a sub-
82 field of analytical chemistry aiming to identify the totality of constituents of crude oil and petroleum
83 refining products using high resolution mass spectrometry methods (Hsu et al., 2011; Marshall and
84 Rodgers, 2004; Palacio Lozano et al., 2020).

85 Despite major advances in the ever-improving analytical resolution of individual molecules
86 and their classes in oils and complex petroleum UVCBs (Wise et al., 2022), there has been relatively
87 little use of the data from these new methods and instruments in regulatory submissions, with the
88 exception of GC×GC-FID-derived data (Redman et al., 2014; Ventura et al., 2011), or even their
89 mention in the reviews or original research publications (**Figure 2A**). The naming conventions and
90 approaches to identification of complex petroleum UVCBs remain rather imprecise (Rasmussen et
91 al., 1999); only general compositional characteristics are used to define broad manufacturing
92 categories (Salvito et al., 2020). While such information is generally sufficient for naming and
93 identification of petroleum UVCBs (ECHA, 2017a), it is often not sufficient for evaluation of human
94 health and environmental hazards, a prerequisite to registration and authorization for their use
95 (Salvito et al., 2020).

96 The regulatory science and analytical chemistry fields run in parallel and both are highly
97 specialized, requiring significant expertise. Consequently, intricacies of the legislative mandates
98 governing regulatory decision-making are often unfamiliar to the researchers who develop and
99 refine advanced methods for petroleomics. Similarly, decision-makers may not be aware of the latest
100 opportunities that analytical chemistry has to offer. Collectively, there is a considerable gap in the
101 translation of knowledge from petroleomics-focused analytical laboratories to applied decision-
102 making. This review aims to first summarize the regulatory guidance for characterization of the
103 chemical composition of petroleum UVCBs and then demonstrate how existing petroleomics
104 techniques could be applied to address these needs. Recent additions to the European Union (EU)
105 Registration, Evaluation, Authorisation and Restriction of Chemicals (REACH) regulatory guidance

106 for demonstrating the composition of complex UVCB (ECHA, 2022), together with recent
107 advancements in petroleomics applications that are poised to address these needs (Palacio Lozano
108 et al., 2020), create a unique opportunity to bridge the divide. Here, we highlight the opportunities
109 that are already within reach for using modern analytical and data analysis/visualization techniques
110 for impactful decisions on petroleum UVCBs. We reason that it is imperative for decision-makers
111 to be aware of the possibilities and limitations of current analytical approaches so that regulations
112 can be sufficiently strict, yet realistic in terms of their attainment using best available science.

113

114 **What Information do the Decision-Makers Seek on Petroleum UVCBs?**

115 Most impactful guidance to the industry on petroleum UVCBs has been produced in the
116 United States and the European Union. These include the United States Environmental Protection
117 Agency (US EPA) High Production Volume (HPV) Challenge Program, a voluntary industry-
118 government information sharing effort that was launched in 1998 (Petroleum HPV Testing Group,
119 2017), and the REACH legislation and associated guidance documents in the European Union
120 (ECHA, 2017c) which is now being adopted in other countries around the globe. Because REACH
121 is the most recent and stringent legislative regime, it has become a *de facto* global driver for the
122 regulatory scrutiny of both new and existing chemicals, including petroleum UVCBs. Therefore,
123 most decision-making contexts discussed herein pertain to EU REACH regulation and its
124 implementation by the European Chemical Agency (ECHA) through guidance documents.

125 Simply put, REACH-based guidance for UVCBs states that data provided by the registrants
126 shall enable (i) identification of the substances that are submitted for registration, and (ii) evaluation
127 of the potential hazards to human health and the environment (ECHA, 2017c). In both instances,
128 information on chemical composition of a substance is required; however, in a slightly different
129 context (**Table 1**). For the former, data requirements are typically less stringent because complex
130 substances may be registered based on the manufacturing process, intended use, and/or physical-
131 chemical properties. For the latter, the molecular identification (both elemental composition and
132 structure) data requirements are far greater because the individual constituents that may be present
133 in a “representative” sample, their amounts, and presence of known or suspected hazardous
134 substances must be reasonably ascertained to enable grouping and read-across among petroleum
135 UVCB substances. Based on these two initial components of REACH, an EU regulation that came
136 into force in June of 2007 (European Council, 2007), decisions on authorization of use or

137 restriction(s) are made. An additional challenge for petroleum-derived and other UVCBs is that their
138 composition is inherently variable from batch to batch of the nominally the “*same*” product; this is
139 also true for product-to-product variability within broader categories (FuelsEurope, 2015). This
140 variability may be due to the source crude oil used in manufacturing; petroleomics-related research
141 has revealed significant differences in compositions of petroleum according to its origins, slight
142 variations in manufacturing processes among refineries, and naturally occurring degradation and
143 weathering processes. Because safety testing is typically conducted with a “*representative*” sample
144 of a UVCB product or category, the confidence with which these data can be extrapolated to other
145 samples of the same product or to other products that are “*similar*,” especially when sample
146 composition data may be commercially sensitive, depends on how much and what type of data is
147 available to ascertain such similarity.

148 For the purpose of substance identification, the regulatory frameworks in the United States
149 and the European Union have historically (**Figure 1**) named, grouped and categorized petroleum
150 UVCBs based on the manufacturing processes used in oil refining, as well as physical-chemical
151 properties and other broad chemical fingerprinting data (Dimitrov et al., 2015; ECHA, 2008;
152 Rasmussen et al., 1999). Manufacturing process-centric naming conventions for these complex
153 substances were originally developed by the American Petroleum Institute (API) and the US EPA
154 for the purpose of creating an inventory of petroleum products under the Toxic Substances Control
155 Act inventory (API, 1983; EPA, 1995; U.S. EPA, 1978). Both Chemical Abstract Service (CAS)
156 and European Inventory of Existing Commercial Chemical Substances (EINECS) identifications
157 have been assigned to a large number of petroleum UVCBs, even though the descriptions of each
158 of these substances under either one of these “unique identifiers” are rather imprecise and far from
159 being unique (Rasmussen et al., 1999). The broad substance categories were somewhat refined
160 under the US EPA’s HPV Challenge Program through addition of more detailed information on
161 physical-chemical properties, as well as some human and environmental hazards data (Petroleum
162 HPV Testing Group, 2017).

163 Far greater information requirements, with respect to both chemical composition and
164 potential hazards, were imposed by REACH (European Council, 2007). Petroleum UVCBs are high
165 production volume substances that were subject to the earliest deadline for registration and the most
166 stringent requirements for hazard evaluation. From 2007 to 2010, about 8,000 registrations were
167 submitted in the EU for petroleum substances that were produced or imported at >100 tons/year

168 (CONCAWE, 2022). Subsequently, the number of registrations of petroleum products was reduced
169 to 191 substances through consolidation of redundant submissions and further grouping of
170 substances deemed to be “similar” based on a variety of considerations. To put the scale of the
171 challenge in context, tens of thousands or even hundreds of thousands of unique molecular formulae
172 can be observed in a single fraction of a petroleum sample using ultrahigh resolution mass
173 spectrometry, indicating that potentially millions of different structures are present in individual
174 crude oils when allowing for isomers (Palacio Lozano et al., 2019b; Palacio Lozano et al., 2020).
175 While the registration submissions for petroleum UVCBs were completed more than 10 years ago
176 and are regularly updated (CONCAWE, 2021), discussions between trade associations and
177 regulatory agencies are ongoing to determine the most sensible ways to improve the dossier quality
178 and ensure the information is compliant with the REACH regulation, primarily by generating more
179 testing information and reinforcing read-across and category approaches. Still, industry’s attempts
180 to waive animal testing requirements through read-across have been rejected by ECHA because of
181 considerable data gaps in hazard assessment and compositional characterization (ECHA, 2020a;
182 ECHA, 2020b; ECHA, 2021).

183 For registration under REACH (Annex VI, Section 2), the data should be sufficient to enable
184 substance identification (**Table 1**). For petroleum substances, ECHA guidance is that the following
185 data should be provided: (1) accepted nomenclature; (2) appropriate identifiers such as
186 source/feedstock, refining history, boiling and carbon number ranges, physio-chemical
187 characteristics, chromatographic or spectral information, flash point, and viscosity; and (3)
188 compositional information including identification and concentration of the individual constituents
189 present at >10% and that are known to be hazardous, persistent and/or bioaccumulative,
190 identification of any additives, and generic description of unknown constituents (CONCAWE,
191 2012; ECHA, 2017a). Additionally, Articles 7(2) and 33 of REACH have defined a concentration
192 threshold of 0.1% w/w for constituents classified as “*substance of very high concern*” (ECHA,
193 2017b). Based on these guidance documents, the registrants (companies or trade associations) have
194 traditionally relied on a wide range of analytical techniques (**Figure 2A**) to furnish the information
195 on petroleum substance identification (Clark et al., 2013; CONCAWE, 2012; CONCAWE, 2014;
196 CONCAWE, 2020).

197 For the evaluation step, REACH specifies (Annex XI, Section 1.5) that substances may be
198 grouped based on “*structural similarity between substances which results in a likelihood that the*

199 *substances have similar physicochemical, toxicological and ecotoxicological properties so that the*
200 *substances may be considered as a group or category.”* Next, prediction of possible hazards of data-
201 poor substances are made through the application of “read-across” from an analogous substance that
202 has been tested in a requisite assay and is deemed “similar.” For this step, REACH regulation states
203 that “*it is required that the relevant properties of a substance within the group may be predicted*
204 *from data for reference substance(s) within the group (read-across approach).*” The registrant shall
205 establish a read-across hypothesis which explains what structural similarities or differences exist
206 between the source and target substance(s) to which read-across is applied and why a prediction for
207 a toxicological or ecotoxicological property can be made with confidence.

208 Recently, REACH regulation (Annex XI, Section 1.5) has been amended to state that for the
209 application in grouping, “*structural similarity for UVCB substances shall be established on the basis*
210 *of similarities in the structures of the constituents, together with the concentration of these*
211 *constituents and variability in the concentration of these constituents”* (European Commission,
212 2021). Once such information becomes available, and in cases where structural differences are
213 present between the source and target substances, the read-across hypothesis should explain why
214 the differences in the chemical structures within a group will not influence the toxicological or
215 ecotoxicological properties or may do so in a regular pattern.

216 Even though the regulatory language above may seem rather straightforward, in practice the
217 bar on establishing “*structural similarity*” is very elusive in the case of petroleum UVCBs. It is
218 widely acknowledged that the chemical complexity of petroleum substances far exceeds the
219 capabilities of any one method, even the highest resolution mass spectrometers; therefore, a
220 combination of techniques and approaches is often employed. However, despite the considerable
221 scientific advances achieved in the past decade (**Figure 1**) in both molecular separation
222 (encompassing chromatography, spectrometry and ionization) and detection (various modalities of
223 mass spectrometry, flame ionization detection, and spectroscopy), many challenges with precision
224 and confidence in comprehensive molecular characterization of petroleum samples persist.

225 For example, a recent study using Fourier transform ion cyclotron resonance mass
226 spectrometry (FT-ICR MS) revealed nearly a quarter of a million molecular formulae in a fraction
227 from one petroleum sample (Palacio Lozano et al., 2019b). This study is an example of the reality
228 that millions of structures (as opposed to formulae, as structure also determines toxicity) are present
229 in petroleum UVCBs, an intractable challenge both for the chemical analysts to resolve and identify,

230 as well as for the decision-makers to evaluate. Indeed, even the REACH regulation itself (Annex
231 XI, Section 1.5) acknowledges that “*if it can be demonstrated that the identification of all individual*
232 *constituents is not technically possible or impractical, the structural similarity may be demonstrated*
233 *by other means, to enable a quantitative and qualitative comparison of the actual composition*
234 *between substances*” (European Commission, 2021). Clarification regarding the definition of “*other*
235 *means*” was recently released May 2022 (ECHA, 2022), where specific justification for lack of
236 knowledge about constituents, lack of published methods to identify constituents, and explanation
237 about technical hindrance to resolution and identification of constituents comprising >20% of the
238 substance is necessary to apply such means for characterization. Where “other” means are justified,
239 quantitative comparison of constituents in common between source and target substances should be
240 included, as well as a qualitative comparison of structures that vary between the substances.
241 Fingerprinting, for example, can be used if the following are addressed: information on >95% of all
242 constituents, information on the constituents of high concern, and high analytical resolution to
243 enable accurate alignment, quantitation, variability, and structural data (beyond molecular formulas)
244 of constituents between substances (ECHA, 2022). Therefore, the subsequent sections of this review
245 are framed around three overarching critical needs/questions that REACH regulation challenges the
246 registrants of petroleum UVCBs to address to gain regulatory acceptance of the grouping and read-
247 across hypotheses that are addressable by means of standardized analytical chemistry methods:

- 248 • Critical Need 1: Providing detailed information on the structure of the constituents;
- 249 • Critical Need 2: Providing information on the concentration of the individual constituents; and
- 250 • Critical Need 3: Demonstrating compositional similarity of complex petroleum UVCBs through
251 other means when the identification of all individual constituents is not technically possible or
252 impractical.

253

254 **Conventional Methods for Characterization of Petroleum Substance Identity and** 255 **Composition**

256 There are many methods for characterization of physical-chemical properties and chemical
257 composition of oil, that have also been applied to petroleum UVCBs (**Figure 2A**). Multiple
258 techniques are used due to the broad range of substances with widely different composition,
259 volatility, and polarity (CONCAWE, 2014; Stout and Wang, 2007; Wang et al., 2011). Most
260 publications concerning characterization of petroleum UVCBs have explored approaches to define

261 elemental composition, physical properties, and gross structural information using nuclear magnetic
262 resonance (NMR) or infrared spectra of these substances. The elemental analyses for
263 characterization of petroleum UVCBs evaluate the concentrations of major elements ranging from
264 carbon and heteroatoms to metals (CONCAWE, 2019). There are specific methods for assessing
265 physical properties, which typically relate to the quality of the product and are thus more regularly
266 measured during manufacturing and more available. For example, measurements of specific gravity
267 using American Society for Testing and Methods (ASTM) methods such as *ASTM D287 Standard*
268 *Method for API Gravity of Crude Petroleum and Petroleum Products* are in wide use (Giles, 2016).
269 Other physical-chemical information, such as boiling and carbon number ranges, is typically
270 deduced through physical (*e.g.*, ASTM methods D86, D1160 and D2892), or simulated
271 (CONCAWE, 2019; CONCAWE, 2020) distillation methods.

272 Spectroscopic techniques have been widely employed to obtain broad compositional
273 information for regulatory characterization and identification of UVCB substances but their utility
274 for the analysis of petroleum UVCBs has been questioned (CONCAWE, 2020). NMR methods
275 (IP392, ASTM D5292) measure the percent of carbon or hydrogen atoms in an aromatic ring
276 (CONCAWE, 2020). Infra-red spectroscopy measures the presence of functional groups to define
277 the degree of saturation in the constituents (CONCAWE, 2012). Ultra-violet spectroscopic analysis
278 quantifies compounds by detecting unsaturated bonds such as those in olefins and aromatics, as well
279 as ketonic and heteroatom groups, but is limited in resolution for other constituents.

280 More detailed compositional information, which gives greater insight into the chemical
281 classes and carbon ranges of the substance, is obtained using chromatographic techniques that
282 enable separation of constituent groups in complex petroleum UVCBs, these include gas- and liquid-
283 based approaches. Gas chromatography (GC)-based analyses predominate; gas-based analysis of
284 hydrocarbons were first published in the 1960s leading to the development of a standardized method
285 (ASTM 2887-84) for determination of *Boiling Range Distribution of Petroleum Fractions by Gas*
286 *Chromatography* (Giles, 2016). GC is a powerful tool used for the separation and semi-quantitative
287 assessment of non-polar constituents such as hydrocarbons and polycyclic aromatic hydrocarbons
288 (PAH) (CONCAWE, 2012). It offers better separation than liquid chromatography (LC) but is
289 limited by the boiling point of compounds (affecting the accessible mass range), which is a greater
290 hindrance for the characterization of the heavy petroleum that is increasingly relied upon, and
291 sometimes compounds must be derivatized to ensure GC compatibility. LC is dependent on the

292 polarities of the constituents present, predominantly used to characterize less volatile polar
293 compounds. High-performance liquid chromatography (HPLC; e.g., ASTM D6379 and IP391
294 methods) is used to quantify mono-, di- and tri-aromatic hydrocarbons (CONCAWE, 2012;
295 CONCAWE, 2019). Meanwhile, thin layer or liquid column chromatography (ASTM D2007)
296 separation generates information on basic chemical properties (CONCAWE, 2012; CONCAWE,
297 2019). Gas chromatography coupled to mass spectrometry detection (GC-MS) is widely used in
298 forensic fingerprinting (US EPA 8270 and 8051B) and to characterize the composition of petroleum
299 UVCBs (US EPA, 1996; US EPA, 2014). Flame ionization detection (FID) has been coupled with
300 both LC and GC for the detection and quantification of hydrocarbons (CONCAWE, 2012).

301 However, conventional standardized methods detailed above are insufficient to establish
302 truly comprehensive compositional characterization of UVCBs as needed by REACH to accept
303 grouping and read-across hypotheses from registrants. Recent decisions by ECHA on testing
304 proposals provided several reasons as to why that is the case (ECHA, 2020a; ECHA, 2020b; ECHA,
305 2021). In these decisions, ECHA noted that (i) physical-chemical characterization of whole complex
306 substances does not demonstrate similarity of chemical constituents of these substances; (ii)
307 elemental and other traditional analysis methods do not provide information on the identity and
308 concentration of individual chemical constituents, but rather provide physical-chemical
309 characterization of the substance as a whole; and (iii) standard methods used in the submissions
310 provided insufficient information to estimate the variability of constituents both within and among
311 substances and groups. ECHA deems these criteria necessary to establish the applicability domain
312 of a given category to confirm membership of the source substance(s) and enable subsequent read-
313 across to the target substance(s) (ECHA, 2020a; ECHA, 2020b; ECHA, 2021).

314 Collectively, despite the use of a battery of analytical assays, expending large sample
315 quantities on some of these analyses, using specialized sample preparation techniques, and incurring
316 considerable costs to acquire these data, the registrants did not establish substance characterization
317 that would be acceptable by ECHA. Indeed, the industry itself acknowledges that spectroscopic
318 techniques are limited to bulk characterization and “*most substances [i.e., petroleum UVCBs]
319 cannot be effectively differentiated from each other by UV, IR, ¹H-NMR or ¹³C-NMR
320 spectroscopies*” (CONCAWE, 2020). Further, even though chromatography-based methods provide
321 considerable amount of information for characterization of nonpolar and relatively volatile
322 compounds (e.g., ASTM D2134, D6729, D6730, among others (CONCAWE, 2012), as well as

323 aliphatic and aromatic fractions (Reddy and Quinn, 1999; Wang and Fingas, 2003), their limited
324 resolution leaves much of the complex substance uncharacterized (Wang et al., 2011; Weng et al.,
325 2015). When assessing the needs for analysis of complex substances, it is important to consider that
326 different analytical approaches are complementary and bring their own advantages and
327 disadvantages. With the limitations of many of the more routine methods, there has been increasing
328 use of ultrahigh resolution mass spectrometry and other advanced methods.

329 More recently, two-dimensional GC (GC×GC) technique has been applied to petroleum
330 substances because it allows for an even greater separation of the multitude of constituents.
331 Specifically, the coupling of GC×GC with FID allows qualitative constituent information at the
332 level of carbon number and chemical class (ASTM International, 2011); this has been useful for
333 simplifying the composition of petroleum substances by binning molecules by “hydrocarbon block”
334 (Redman et al., 2012). Coupling GC×GC with mass spectrometry provides more structural
335 information on the specific constituents (Jennerwein et al., 2014; Mao et al., 2009).

336

337 **High and Ultrahigh Resolution Mass Spectrometry Techniques**

338 Advancements in mass spectrometry over the past 80 years (**Figure 1**) have spawned the
339 application of high-resolution approaches for the study of petroleum substances at a molecular level
340 (Palacio Lozano et al., 2019a; Wise et al., 2022). Resolving power is one of the key performance
341 metrics of any mass spectrometry and is typically defined as $\frac{m}{\Delta m}$, where m is the m/z of the ion of
342 interest and Δm is the width of the peak at half its height, using the full width at half maximum
343 (FWHM) definition. In essence, the higher the resolving power, the more peaks can be observed for
344 complex samples, as this reduces overlap of peaks (Phillips et al., 2022). “High resolution” has
345 typically been accepted to mean a resolving power of >10,000 (Xian et al., 2012), but it has become
346 well-established over the past two decades that “ultrahigh resolution,” which often refers to a
347 resolving power of >100,000, is essential for characterization of the most complex samples such as
348 petroleum.

349 The second performance metric that should always be considered is that of mass accuracy.
350 An instrument that offers high mass accuracy indicates that it typically provides data with low mass
351 errors. When an elemental composition (molecular formula) has been assigned to an observed peak,
352 the observed m/z and the theoretical m/z for the assignment can be used to calculate the mass error,
353 measured in parts per million (ppm), using the following equation:

354
$$\frac{m_{observed} - m_{theoretical}}{m_{theoretical}} \times 1,000,000$$

355 Note that a negative mass error indicates the peak appears at a lower m/z than the theoretical value
356 and a positive mass error indicates the peak appears at a higher m/z than the theoretical value; for
357 each peak, the closer the mass error is to zero, the greater the confidence in the given assignment of
358 the molecular formula. It is also important to note that, while researchers desire resolving power to
359 be as high as possible, they strive to keep mass errors as low as possible. Thus, the highest resolving
360 powers afford researchers the ability to observe more components within complex samples, while
361 the highest mass accuracies (lowest mass errors) afford greater confidence in the assignments of
362 elemental composition and structures, together establishing detailed compositional “profiles,”
363 “fingerprints,” or “signatures” for complex samples. Within the field, there has been discussion
364 about the most appropriate terminology to use, drawing parallels with how fingerprints do not
365 normally change but that signatures do, and how this understanding may be related to compositions
366 of complex substances changing when subject to anthropogenic or environmental processes.

367 The comprehensive characterization of molecular composition of petroleum through high
368 resolution mass spectrometry is an active area of investigation and includes several approaches to
369 precise detection, naming, and structural characterization of the individual constituents (Hsu et al.,
370 2011; Marshall and Rodgers, 2004; Niyonsaba et al., 2019; Palacio Lozano et al., 2019a; Palacio
371 Lozano et al., 2019b; Palacio Lozano et al., 2020; Roman-Hubers et al., 2022; Xian et al., 2012).
372 Modern time-of-flight (TOF) mass spectrometers, which are widespread and considered high
373 resolution, offer resolving powers typically in the range of 10,000-60,000. By contrast, FT-ICR MS
374 is the highest performance variety of mass spectrometer and is considered to offer researchers
375 ultrahigh resolution, at one or two orders of magnitude higher performance. FT-ICR MS is based
376 upon ions orbiting inside a cell, which is in turn housed within the bore of a superconducting magnet,
377 and the technique offers ultrahigh resolving power ($\sim 10^6$ FWHM) and mass accuracy (sub-ppm).
378 Orbitrap mass spectrometers are a newer variety of mass spectrometer based upon a Kingdon trap
379 design, rather than using magnetic fields; these instruments typically offer a resolving power of $\sim 10^5$
380 FWHM. FT-ICR MS and Orbitrap MS offer differing degrees of ultrahigh resolution, with FT-ICR
381 MS offering the highest performance; however, the advantage of TOF MS is in rapid acquisition
382 time which allows for coupling with additional separation techniques such as two-dimensional gas
383 chromatography and ion mobility spectrometry (Palacio Lozano et al., 2019a). These three
384 analytical techniques (**Figure 2B**) are actively used in petroleomics analyses because they offer

385 somewhat different approaches to determining molecular formulae present in complex substances
386 (Palacio Lozano et al., 2019b; Rodgers and McKenna, 2011), but have not yet been used for detailed
387 petroleum substance characterization for registration or evaluation purposes.

388

389 **Ultrahigh Resolution MS Data Processing and Visualizations**

390 Petroleum substances contain highly homologous series of hydrocarbon molecules; thus,
391 complex substance analysis can be facilitated by exploiting the patterns of various chemical groups.
392 The method of Kendrick mass defect (KMD) analysis facilitates sorting molecules into homologous
393 series (Kendrick, 1963). The composition of complex petroleum substances can be visualized using
394 the KMD approach because most molecules belong to homologous series comprised of (CH₂) alkyl
395 groups and other functional groups (**Figure 3A**), this method has been widely used in petroleomics
396 (Hughey et al., 2001; Marshall and Rodgers, 2004; Marshall and Rodgers, 2008; Palacio Lozano et
397 al., 2020). Due to the high resolution and mass accuracy, the molecular composition assigned to the
398 ions that fall in or out of the homologous series can be used to predict their elemental content (HC_#,
399 O_#, N_#, O_#, S_#), carbon number, and double bond equivalents ($DBE = C_{\#} - H_{\#}/2 + N_{\#}/2 + 1$). A
400 number of other visualizations have been proposed to express the molecular composition of various
401 substances based on the rings and double bonds in the carbon framework (*i.e.*, DBE) of the
402 constituents that can be plotted against their carbon number (**Figure 3B**). Van Krevelen diagrams
403 (**Figure 3C**) are used to display the degree of aromaticity and oxidation of constituents by plotting
404 the H/C versus O/C ratio of the organic compounds in the complex substance (Kim et al., 2003; Van
405 Krevelen, 1950; Van Krevelen, 1984). Relative abundance of various classes of compounds in the
406 samples is typically plotted as a chemical class bar chart (**Figure 3D**), or as two-dimensional
407 “hydrocarbon blocks” (**Figure 3E**).

408 Molecular-level analyses using ultrahigh resolution MS typically generate voluminous
409 datasets even after processing by the software provided by instrument manufacturers. A number of
410 third-party software packages have been developed recently to facilitate data analysis and
411 visualizations. *Peak-by-Peak* fee-based software (Spectroswiss, 2019) is tailored for processing FT-
412 MS raw data on transients and mass spectra and generates output files for follow-up statistical
413 analyses. *PetroOrg* (Riches et al., 2015) and *Composer* (Sierra Analytics, 2022), launched in 2014
414 and 2008, respectively, are two stand-alone fee-based packages that allow processing and
415 visualization of the molecular composition of petroleum substances data acquired through high

416 resolution MS. *UltraMassExplorer* is an open source web-based package that uses *R Studio*
417 (Leefmann et al., 2019). It applies static formula libraries for molecular formula assignment of the
418 molecular formulas based on neutral masses coupled to PubChem searches for putative structural
419 assignment. Data visualization is enabled by van Krevelen, KMD and DBE plots. *DropMS* is
420 another web-based tool that facilitates high resolution MS data processing and molecular
421 assignments with its corresponding DBE, error, signal intensity, as well as a number of
422 visualizations (Rosa et al., 2020). *KairosMS* utilizes an R Shiny interface to process complex data
423 sets produced through hyphenated MS experiments that, when combined with software for formula
424 assignments, can be used for visualization, comparison, and statistical analyses for both direct
425 infusion and hyphenated data set, using a wide variety of approaches (Gavard et al., 2020).
426 *Interactive van Krevelen* (Kew et al., 2017) and *Open van Krevelen* (Brockman et al., 2018) are
427 two packages that offer interactive diagrams for molecular-level exploration of the data from high
428 resolution MS.

429

430 **The Regulatory Needs for Characterizing Chemical Composition of Petroleum UVCBs**

431 General considerations and examples of the application of high resolution MS for the
432 analysis of petroleum substances have been reviewed elsewhere (Hsu and Shi, 2013; Niyonsaba et
433 al., 2019; Palacio Lozano et al., 2019a; Palacio Lozano et al., 2019b; Palacio Lozano et al., 2020;
434 Rodgers and McKenna, 2011; Xian et al., 2012). However, these previous reviews did not
435 specifically place the application of these techniques in the context of the regulatory needs for
436 registration and evaluation of petroleum UVCBs. Comprehensive characterization, or at least more
437 detailed information on some specific types of constituents, of complex substances based on high
438 resolving power and mass accuracy of FT-ICR MS, Orbitrap MS and TOF MS affords the
439 opportunity to attain valuable information on both molecular composition and relative abundance
440 of the constituents. Such data can also be used to quantitatively evaluate both variability within a
441 substance or group of substances, and similarity between substances. Then, it should be possible to
442 conduct read-across and complete hazard characterization. To provide specific examples of how
443 each technique can be used to address specific regulatory needs, we performed a systematic
444 literature search (with end date of December 2021) focused on the application of each technique to
445 petroleum substances (**Figure 4**). We found that even though there are many dozens to hundreds of
446 publications on each technique in general, and on their use to analyze petroleum samples in

447 particular, few studies presented data or drew conclusions in a manner relevant to address each
448 regulatory need. The Web of Science search terms and results of the literature search are included
449 in **Supplemental Table 1**.

450

451 ***Critical Need 1: Providing detailed information on the structure of the constituents.***

452 High mass measurement accuracy and resolving power of modern mass spectrometry
453 analyzers are key steps to address this regulatory need for petroleum UVCBs. Regulatory agencies
454 have repeatedly stressed that “*broad*” compositional information is insufficient to justify groupings
455 and support read-across hypotheses (ECHA, 2014; ECHA, 2020a). While many publications tout
456 ever increasing resolution as an advance in the science of analytical chemistry, we highlight recent
457 examples of studies that focused on comprehensive structural information of the individual
458 constituents in petroleum samples.

459 With the advances in the performance of mass spectrometers and development of new
460 methodologies, significant strides have been made with respect to characterization of mixtures and
461 petroleum samples. For the analysis of whole crude oils and refined products, FT-ICR MS could
462 resolve hundreds of molecular formulae in the late 1990s (Rodgers et al., 1998) and several thousand
463 in the early 2000s (Hughey et al., 2002; Qian et al., 2001). By 2019, newer FT-ICR MS
464 instrumentation and experiment design led to the identification of ~245,000 molecular formulae in
465 a non-distillable fraction of the maltenes from a heavy petroleum sample, using a resolving power
466 of more than 3,000,000 across the entire m/z range (Palacio Lozano et al., 2019b). Thus, in less than
467 two decades, there was a 60-fold increase in resolving power and 80-fold increase in the number of
468 molecular formulae being assigned within a single sample.

469 While ultrahigh resolution mass spectrometry has offered unprecedented levels of insight
470 into highly complex samples, such instrumentation typically affords molecular formulae but does
471 not lead to definitive identification of structures. Experimental parameters, such as polarity and
472 ionization behavior, can yield information about probable functionalities of the components, but
473 true structural identification would require additional data. This could be obtained by fragmentation
474 patterns associated with individual peaks (and hence molecular formulae), acquired during “tandem
475 mass spectrometry” (also known as “MS/MS”) experiments. As complex samples may comprise
476 tens of thousands or even hundreds of thousands of molecular formulae, performing hundreds of

477 thousands of tandem mass spectrometry experiments to target all of the peaks is not viable, due to
478 the time required and workload due to data analysis.

479 Combination with an orthogonal separation method, such as a form of chromatography or
480 ion mobility, is an alternative approach for accessing information about functional groups and
481 isomeric contributions. Optimization of these experiments (such as choice of columns, temperature
482 program, etc.) is similarly non-trivial, however, and resolution of components in complex
483 substances can be challenging; the concept of “unresolved complex mixtures”, also referred to as an
484 “UCM hump” on the chromatograms, is well-known (Gough and Rowland, 1990). Furthermore,
485 while chromatography is a well-known method for distinguishing isomers, this is typically
486 combined with a need for authentic standards (e.g., to determine expected retention times on the
487 chromatography column) and databases, although certainly not all available compounds are found
488 within databases. Where there may be dozens of isomers per molecular formula, this would in
489 practice mean a requirement for millions of authentic standards to be run in order to address the
490 most complex petroleum samples that have been characterized by ultrahigh mass spectrometry. This
491 is, again, impractical due to the amount of time required. It is also worth drawing attention to a
492 subtle distinction between molecular formulae (e.g., $C_{18}H_{36}O_2$) and the concept of “peaks;” a peak
493 may represent a single molecular formula in a direct infusion experiment by mass spectrometry (and
494 data is represented by m/z and intensity only), but in combination with mass spectrometry with
495 chromatography, there will be an additional dimension to the data (m/z and intensity, but now also
496 time) and multiple peaks per molecular formula may occur due to the presence of isomers.
497 Combining mass spectrometry with chromatography can yield more peaks by counting of isomers,
498 but not necessarily more molecular formulae.

499 The isobaric information afforded by ultrahigh resolution direct injection FT-ICR MS has
500 played a significant role in comprehensive understanding of petroleum substance composition;
501 however, the presence of a large number of isomeric hydrocarbons in petroleum substances
502 challenges the utility of ultrahigh resolution methods that do not provide other means for separation
503 of isomeric species. Some studies addressed a growing need to define the isomeric composition of
504 petroleum substances by coupling FT-ICR MS with separation techniques. Gas chromatography
505 prior to FT-ICR MS was used for separation of isomers in petroleum by resolving individual
506 constituents (Barrow et al., 2014; Palacio Lozano et al., 2022). Information of the structural
507 composition and the isomeric diversity can be also attained by coupling with trapped wave ion

508 mobility spectrometry (TWIMS), a variety of ion mobility spectrometry, to FT-ICR MS (Maillard
509 et al., 2021). Multidimensional separation with ultrahigh resolution allows for characterization of
510 individual constituents with accurate mass measurements and structural features. Additionally, high
511 mass accuracy and ultrahigh resolution measurements have elucidated structural information of the
512 individual constituents through post-instrumental data analysis (Cho et al., 2011; Hu et al., 2018).

513 A number of studies have coupled chromatography to Orbitrap MS for analysis of complex
514 environmental samples (MacLennan et al., 2018; Pereira et al., 2013; Sorensen et al., 2019; Yang et
515 al., 2019). Some studies applied this technique for the analysis of oils or fractions thereof. A
516 combination of GC separation with different ionization methods for Orbitrap MS was used to study
517 volatile components of a petroleum refining product, gas condensate. The authors showed that
518 separation of different isomeric compounds could be achieved using this hyphenated method thus
519 aiding in deeper characterization of a complex substance (Kondyli and Schrader, 2019). More
520 recently, a reverse-phase liquid chromatography method was applied to the analysis of petroleum
521 refining-derived UVCBs and compared to direct injection Orbitrap and FT-ICR MS (Xia et al.,
522 2021). The authors showed that not only could they obtain elemental formulas for a large number
523 of hydrocarbon and heteroatom species, but chromatography-informed retention patterns could also
524 be used to distinguish among isomeric species, hence increasing confidence in structural
525 identification of the individual constituents in petroleum substances.

526 TOF MS detection, despite its lower mass resolution as compared to FT-ICR and Orbitrap
527 MS, affords an advantage of rapid data acquisition; this technique has been explored extensively by
528 coupling with chromatography or ion mobility spectrometry separations (Palacio Lozano et al.,
529 2019a). The most common type of front-end chromatography in petroleum substance analysis is
530 GC×GC. The additional separation afforded by this technique offers substantial improvement for
531 resolving isomeric constituents and compounds that would otherwise coelute (Ball and Aluwihare,
532 2014; Luna et al., 2014a; Ngo et al., 2012; Rowland et al., 2011). For example, this technique was
533 used to identify a large number of isomeric molecules, both aliphatic and aromatic, that could not
534 be resolved by a typical GC-MS technique in a complex milieu of hydrocarbons containing a wide
535 (C₁₂-C₃₆) range of carbon numbers (Alam et al., 2016). Similarly, the ability GC×GC-TOF MS to
536 resolve isomers of classical and sulfur-containing naphthenic acids in oil-contaminated
537 environmental samples enabled contamination-source fingerprinting (Bowman et al., 2019) and
538 identification of potentially hazardous substances (Bowman et al., 2020). TOF MS has also been

539 coupled with ultrahigh pressure liquid chromatography to isolate isomeric species in petroleum
540 derivatives (Lv et al., 2013; Mahmoud and Dabek-Zlotorzynska, 2018).

541 The compositional characterization of petroleum samples has typically relied upon either the
542 use of chromatography or, more recently, the use of ultrahigh resolution mass spectrometry. The
543 separation on the basis of retention time using gas chromatography can be sufficient for monitoring
544 targeted components (*e.g.*, steranes for differentiation of sample origins) and revealing the most
545 significant contributions. As samples become more complex, however, particularly with the
546 increasing use of heavy petroleum sources, there is the increasing probability of coeluting
547 components, as mentioned above, where compounds have the same retention time and therefore
548 cannot be distinguished by one-dimensional chromatography alone. The use of a second column of
549 a different type to result in a two-dimensional approach, GC×GC, can significantly improve the
550 separation and therefore the number of compounds observed.

551 Due to the temperature ranges of the ovens, gas chromatography-based methods have limited
552 retention times that are accessible, as these are, in turn, linked to the boiling points of the individual
553 compounds. Direct infusion methods coupled with ultrahigh resolution mass spectrometry (*i.e.*,
554 where samples are injected directly into an ion source, with no preceding chromatography) do not
555 have the limitations associated with boiling point, and ultrahigh resolution approaches with direct
556 injection delivery have become increasingly significant, especially for heavy petroleum samples.

557 There has also been limited coupling of chromatography with ultrahigh resolution mass
558 spectrometry, which affords researchers the ability to separate compounds on the basis of two
559 different dimensions: retention time and m/z . In this way, co-eluting compounds (*i.e.*, the same
560 retention time associated with the GC column) can be separated due to the additional use of the m/z
561 dimension (Barrow et al., 2014). The data sets acquired can be large (*e.g.*, ~25-50 GB) and therefore
562 present data processing and data analysis challenges. The combination of orthogonal approaches,
563 yielding so-called “hyphenated” techniques, represents means by which to access the advantages of
564 the individual methods. Ultrahigh resolution mass spectrometry affords the ability to resolve a
565 greater number of peaks and assign many thousands of unique molecular formulae with confidence,
566 while chromatography or ion mobility affords the ability to separate isomers (same molecular
567 formulae and therefore same m/z , but differing structures due to different arrangements of the
568 atoms). While GC and low resolution GC-MS methods have seen widespread usage for many years,

569 it can be expected that multidimensional GC and the coupling of orthogonal methods with ultrahigh
570 resolution mass spectrometry will both become increasingly used.

571 Obtaining structural isomeric information is also possible using ion mobility mass
572 spectrometry (IMMS), a post-ionization separation technique, often coupled to TOF MS and used
573 for the analysis of petroleum samples (Santos et al., 2015). IMMS provides structural information
574 that is complementary to the observed m/z of a compound by characterizing the spatial conformation
575 of individual constituents via their drift time through an inert gas and the subsequent derivation of
576 collision cross section (CCS) values (Dodds and Baker, 2019). CCS values can also be compared
577 with those determined through computational means to determine structures for observed
578 compounds. Currently, this is laborious for complex samples and so not viable for samples with, for
579 example, tens of thousands of molecular formulae and their associated isomers. The combination of
580 ion mobility, mass spectrometry, and computational chemistry does, however, hold potential for
581 providing greater structural insights during characterization of samples. IMMS allows for separation
582 of isomeric compounds by their structural composition (Hoskins et al., 2011; Lalli et al., 2015; Lalli
583 et al., 2017; Mahmoud and Dabek-Zlotorzynska, 2018). The high mass accuracy and resolution of
584 TOF MS, coupled with structural characterization using ion mobility, has allowed comprehensive
585 elucidation of the composition of complex petroleum substances and byproducts (Lalli et al., 2015;
586 Lalli et al., 2017). Confident molecular formula assignment to the IMMS-derived features in a
587 gasoline standard and a crude oil sample was demonstrated through the use of KMD analyses based
588 on CH_2 and H functional units (Roman-Hubers et al., 2021).

589 While instrumentation has increased in performance, the data analysis methods have often
590 struggled to keep pace. It is evident that highly complex samples are being characterized better than
591 ever before. The field is acutely aware that handling the increasingly complex data is already
592 challenging, but that also there is a need to obtain greater structural insights, going beyond molecular
593 formulae alone and/or ensuring all of the many thousands of peaks are associated with a definitive
594 structure, not only those that are mostly easily targeted. While the progress with complex substance
595 analysis has been remarkable in recent years, there remain mountains to climb.

596
597 ***Critical Need 2: Providing information on the concentration of the individual constituents.***

598 Quantification of individual constituents in complex and multi-constituent substances is a
599 required step for hazard evaluation to ensure no underestimation of the potential human health and

600 environmental hazards (CONCAWE, 2012; ECHA, 2017a). Thus, the ability of high resolution
601 mass spectrometers to determine the thousands of individual constituents is not sufficient without
602 determination of their abundance. Traditional quantitative approaches relying on mass spectrometry
603 detection require the use of various extraction and detection standards, preferably isotopically-
604 labelled ones (Urban, 2016). However, because of the complexity of petroleum substances, the use
605 of standards for absolute quantitation of numerous, rather than a small number of targeted
606 constituents, would potentially mean the need for millions of authentic standards, which is both
607 impractical and impossible. Instead of quantifying absolute concentrations of multiple individual
608 constituents, traditional analysis techniques such as GC-MS and GC×GC-FID derive relative (*e.g.*,
609 fraction of total) amounts for groups of compounds based on a limited number of standards for key
610 classes of constituents. For example, quantitation of individual n-alkanes, selected isoprenoids,
611 polycyclic aromatic and alkyl polycyclic aromatic hydrocarbons, and biomarker compounds is
612 possible by GC-MS (Wang et al., 1994); however, the “UCM hump” of high molecular weight
613 hydrocarbons limits the utility of this technique for hazard evaluation. GC×GC-FID technique is
614 also commonly used for petroleum analyses to derive “hydrocarbon blocks” (ASTM International,
615 2011). This technique is more amenable for hazard evaluation as it can separate polycyclic
616 compounds with known or suspected hazardous properties (Bierkens and Geerts, 2014); however,
617 this technique is not considered sufficiently informative by some decision-makers, especially for
618 the substances that contain C30 or greater hydrocarbons (ECHA, 2020a). Therefore, there is a great
619 need to determine the ability of high resolution MS techniques to provide quantitation of the
620 individual constituents in petroleum substances.

621 There are significant challenges with respect to quantifying the abundance of constituents in
622 petroleum UVCBs. The signal observed for a given compound will be influenced by a number of
623 experimental parameters. One of these is, of course, concentration. Other factors include solubility
624 in the chosen solvents, pH of the sample solution, the suitability of ionization method chosen, the
625 polarity of the ion source (*e.g.*, when using electrospray ionization, acidic species will be observed
626 in negative-ion mode while basic species will be observed in positive-ion mode), and instrument
627 tuning, amongst other variables. As one example, alkanes can represent as much as half or more of
628 the composition of some petroleum samples and yet if performing analysis using electrospray
629 ionization, a widespread ionization method which is suitable for observation of polar and ionic
630 species, then the alkanes would essentially go unobserved.

631 The matter of quantification is also closely linked to structure, as previously discussed.
632 Different isomers of a molecular formula may have differing functional groups which, in turn, may
633 influence the ionization response of the compounds. For instance, an organic molecule containing
634 two oxygen atoms may incorporate two hydroxyl groups or a carboxylic acid, where the carboxylic
635 acid would ionize much more readily than the former, and therefore give a much stronger signal.
636 For this reason, quantification typically involves the coupling of chromatography to mass
637 spectrometry, in order that isomers are separated the signals associated with each isomer can be
638 measured and considered separately. Such experiments are labor intensive, however. Authentic
639 standards must be used for each isomer of each molecular formula, as mentioned previously, but the
640 standards must also be prepared as a series of sample solutions spanning a range of concentrations
641 to measure the signal intensity as a function of concentration, leading to calibration curves. The
642 needs for authentic standards and for calculation of the calibration curves leads to increases in
643 researcher workload akin to orders of magnitude. “Untargeted” analyses could necessitate the need
644 for potentially millions of experiments for quantification purposes, depending on the objectives. It
645 is much more common that “targeted” analyses are instead used, where a short list of compounds of
646 concern are searched for, such as benzene, toluene, ethylbenzene, and xylene compounds. Many
647 advanced mass spectrometry approaches which are commonly used are accepted to be semi-
648 quantitative, balancing sample complexity, experimental design, data processing, and time.

649 Most studies that use FT-ICR MS for the analysis of petroleum samples report tens to
650 hundreds of thousands of detectable constituents; however, this technique is semi-quantitative where
651 abundance of each molecule would depend on multiple factors (*e.g.*, solubility, concentration,
652 ionization response, tuning, etc.). FT-ICR MS studies typically report relative abundances of various
653 hydrocarbon or heteroatom classes rather than that of individual constituents (Bae et al., 2010; Chen
654 et al., 2012; Jennerwein et al., 2014; Kim et al., 2015; Oldenburg et al., 2014; Walters et al., 2015).
655 Some publications focused on detection and quantitation of specific constituents, such as organic
656 sulfur compounds (Lu et al., 2013), or metalloporphyrin complexes (Cho et al., 2014); however,
657 these constituents have uncertain relevance for the purpose of hazard evaluation. The ultrahigh
658 resolution Orbitrap MS technique has not been used for the characterization of the individual
659 constituents; instead, the abundances of the detected ions with inferred elemental composition are
660 used for semi-quantitative evaluation of various fractions and broad chemical classes (Castiblanco
661 et al., 2020; Liu et al., 2020; Rodrigues Covas et al., 2020; Vanini et al., 2020).

662 Similarly, quantification of the individual constituents using TOF MS is challenging and
663 most studies focused on quantifying abundances of broad classes of compounds (Scarlett et al.,
664 2008). A combination of GC×GC-TOF MS and GC×GC-FID was used to conduct qualitative and
665 quantitative analysis of polycyclic aromatic hydrocarbons for several petroleum UVCBs; however,
666 quantitation was performed only for the constituents for which chemical standards were used (Ristic
667 et al., 2018). Even though absolute quantitation of the individual constituents in complex petroleum
668 samples may be unattainable, the combination of confident molecular formula assignments aided
669 by structural information provided by IMMS and data on relative abundance does enable
670 quantitative evaluation of the most abundant constituents (Roman-Hubers et al., 2021). For example,
671 a study of variability in chemical composition of petroleum UVCBs used the average relative
672 abundance of each constituent in a product across production cycles to determine what molecules
673 may be present at relatively high (*e.g.*, REACH threshold of concern at 0.1% (ECHA, 2017b))
674 amounts and whether those molecules vary significantly between production cycles (Roman-Hubers
675 et al., 2022). This study is an example of how high resolution MS can be used to not only
676 characterize individual constituents, but also to determine their abundance for consideration as
677 potential substances of concern in hazard evaluation. Indeed, the ability to quantify, even in relative
678 terms, the abundance of the identifiable constituents that comprise most of the petroleum sample is
679 of utmost interest under the REACH framework (**Table 1**). The most recent advice from ECHA
680 states that “*all constituents present in a concentration at or above 1% must be identified*” when
681 grouping or read-across is proposed (ECHA, 2022). The bar is even greater for specific constituents
682 that may possess hazardous properties, *i.e.*, “*0.1% for constituents that are classified as*
683 *carcinogenic or mutagenic and 0.3% for substances that are toxic to reproduction or development.*”
684 Collectively, “*identified constituents above the thresholds given above must account for a minimum*
685 *of 80% of the mass of a UVCB substance*” (ECHA, 2022).

686
687 ***Critical Need 3: Demonstrating compositional similarity of complex petroleum UVCBs through***
688 ***other means when the identification of all individual constituents is not technically possible or***
689 ***impractical.***

690 Even though challenges remain in the confident identification of structures for individual
691 constituents in complex substances, the multi-dimensional data from high resolution MS analyses
692 (FT-ICR, Orbitrap and TOF MS) is highly valuable for evaluating broad similarities among complex

693 petroleum UVCB substances and for identifying the degree of variability within a class of
694 substances or between production batches of the same substance (**Figure 4**). Indeed, a number of
695 studies have used these data to perform statistical analyses and visualize the relationships between
696 samples in large datasets to demonstrate that even broad molecular compositional data can be
697 effective as a means for screening, evaluating similarity and determining what constituents may be
698 variable to prioritize selected samples and constituents that may warrant further targeted quantitative
699 analyses.

700 The variation in composition of hydrocarbons and heteroatoms resolved using FT-ICR MS
701 has been explored for the analysis of crude oils (Hosseini et al., 2021; Rocha et al., 2018; Silva et
702 al., 2020) and refining products (Benassi et al., 2013; Mennito and Qian, 2013; Oldenburg et al.,
703 2014; Oldenburg et al., 2017; Orrego- Ruíz, 2018; Silva et al., 2020; Walters et al., 2015; Wang et
704 al., 2020), including studies that used such data for categorization of new products (Abib et al.,
705 2020; Hourani et al., 2013; Liu et al., 2014). Through detailed chemical profiles, compositional
706 characterization of different environmental samples allowed detection of the asphaltenes (Neumann
707 et al., 2021; Ruger et al., 2015), fulvic acids (Stenson et al., 2003), oil sands process-affected water
708 (Barrow et al., 2016), bitumen (Lacroix-Andrivet et al., 2021), biochar-derived organic matter (Li
709 et al., 2022a), as well as soil and sedimentary organic matter (Zhong et al., 2011). Studies of oil
710 weathering (Wozniak et al., 2019) and transformation in the environment (Jaggi et al., 2019; Li et
711 al., 2022b; Wozniak et al., 2019) enabled to not only determine the trends in compositional changes,
712 but also to group samples, allow for forensic identification of the related samples, predict the
713 potential environmental impact of oil spills, and designing mitigation strategies.

714 While there will be limitations to the lowest m/z that can be detected by FT-ICR MS, based
715 upon the magnetic field strength and highest frequency that can be detected (*e.g.*, $\sim m/z$ 37 for a
716 modern 12 T instrument), with appropriate instrument tuning, both Orbitrap and FT-ICR mass
717 spectrometers can offer comprehensive insights into the lower molecular weight compounds of the
718 chemical profile of complex substances, not afforded by other high resolution MS instruments (Chen
719 et al., 2018; Cheng and Hous, 2021; Cho et al., 2017; Headley et al., 2011). Through the application
720 of Orbitrap MS the variations of the mass spectrum profiles can be traced at a molecular level with
721 low mass error to determine the overall chemical composition (Castiblanco et al., 2020; Dong et al.,
722 2019; Porto et al., 2019; Silva et al., 2019; Vanini et al., 2020). Liquid chromatography separation
723 techniques coupled to Orbitrap MS have provided qualitative and quantitative monitoring of the

724 organic species through the chemical profile complex substances. (Folkerts et al., 2019; Miles et al.,
725 2020; Sorensen et al., 2019; Xia et al., 2021)

726 When designing future strategies for screening petroleum-related compounds, it is important
727 to consider what is being measured by different approaches. For example, with electrospray
728 ionization, acidic species would typically deprotonate and be observed using negative-ion mode,
729 while basic species would protonate and be observed using positive-ion mode; acidic and basic
730 species would not be observed in the same experiment, therefore. While this might initially be
731 considered a disadvantage because of the need for two experiments instead of one, on the other hand
732 this can be used to the researcher's advantage when needing to differentiate by functionality. An
733 organic molecule with a heteroatom content of only one nitrogen atom (*e.g.*, $C_cH_hN_1$) could be
734 pyrrolic (weakly acidic) if observed in negative-ion mode, but would more likely be pyridinic (basic)
735 if observed in positive-ion mode. Where electrospray ionization is suitable for polar and ionic
736 compounds, it does not readily ionize non-polar species and so, for these, an alternative ionization
737 method, such as atmospheric pressure photoionization, would be used. To cover acidic, basic, and
738 non-polar species in an untargeted manner, three experiments may be used, but the number of
739 methods employed could be reduced if instead adopting a targeted approach. Through an awareness
740 of the advantages and disadvantages of the different ionization methods, combined with clear
741 objectives of the screening process (which could become of a targeted nature), it is possible to tailor
742 a system which balances required information and workload.

743 Multidimensional GC separation coupled with high resolution detection in TOF MS has been
744 widely employed to map out the composition of complex substances. The high-throughput
745 acquisition has been readily employed to screen for saturate and aromatic hydrocarbons to
746 characterize the broad compositional profile (Haitao et al., 2013; Hao et al., 2017; Kulkarni and
747 Thies, 2012; Luo et al., 2016; Qian et al., 2004; Rui et al., 2012). Nevertheless, to achieve a better
748 compositional information and distribution patterns of complex volatile constituents, TOF MS has
749 been coupled with gas chromatography for qualitative and quantitative assessment (Haitao et al.,
750 2013; Hao et al., 2017). When coupled with an HPLC, non-volatile components can be readily
751 separated based on isomeric and isobaric information for comprehensive high resolution
752 characterization (Cao et al., 2020). The drawbacks observed from the above addressed separation
753 techniques prompted the coupling of capillary electrophoresis with TOF MS for qualitative
754 assessment of high molecular weight constituents in a heavy gas oil (Nolte et al., 2013). Compared

755 to other TOF MS hyphenations, GC×GC offers multidimensional high resolution characterization
756 to map out the composition of complex substances (Alam et al., 2018; Frenzel et al., 2010; Gabetti
757 et al., 2021; Luna et al., 2014b; Muller et al., 2020; O'Reilly et al., 2019; Qian and Wang, 2019;
758 Ristic et al., 2018; Zhu et al., 2020). The plots generated from GC×GC-TOF MS analysis plotting
759 constituents in a two-dimensional space (1st retention time versus 2nd retention time) help elucidate
760 the composition of complex substances and facilitates the molecular classification of compounds
761 (Damasceno et al., 2014). Notably, physical characteristics (*i.e.*, density) can be directly correlated
762 based on the detailed chemical composition defined through GC×GC-TOF MS (Vozka et al., 2019).

763

764 **Petroleum UVCB composition: A regulatory challenge, but what are the solutions?**

765 As summarized above, novel high resolution mass spectrometers offer comprehensive
766 characterization of petroleum substances to qualitatively assess and measure the broad chemical
767 composition of a substance and its individual constituents. The information available through high
768 resolution characterization of the chemical composition of UVCBs can address the shortcomings
769 with regards to the prediction of toxicological properties when practicing read-across assessment.
770 Comprehensive characterization can provide sufficient information on the molecular composition
771 and their relative abundance within a substance to define the commonality between substances and
772 their similarity at a molecular level. With this information, it is then possible to determine
773 constituents which may be used to infer human health hazard properties of the whole substance.

774 Still, complex UVCBs, especially those produced by refining of oil, remain to be an evolving
775 competency in regulatory science. In the recent past, both decision-makers and industry relied on
776 rather imprecise substance categorization and read-across hypotheses to predict toxicological
777 properties (Clark et al., 2013; McKee et al., 2015; Salvito et al., 2020). Recognizing the lack of
778 clarity in the original regulations and guidelines (ECHA, 2017c), and persistent issues identified in
779 regulatory submissions of petroleum UVCBs (ECHA, 2020a; ECHA, 2020b; ECHA, 2021),
780 additional advice was recently provided with respect to the information on chemical composition
781 for UVCB substances (ECHA, 2022). Indeed, there appears to be a major gap in what information
782 is perceived as sufficient in terms of chemical characterization of petroleum UVCBs by either
783 decision-makers or industry. At the same time, the science of analytical chemistry delivered a
784 number of major advances in terms of novel techniques and visualizations to aid in deep
785 characterization of complex petroleum UVCBs. Scores of research studies have been published,

786 reviews written, and lectures delivered. A “meeting in the middle” is required between those
787 working in regulatory and analytical disciplines, so that the translation of petroleomics science into
788 the practice of decision making on petroleum UVCBs can begin to be realized.

789 One barrier to such a translation could be the difference of opinions on what “science” is
790 “evidence” and whether the scientific data are “sufficient” to satisfy the regulatory requirements for
791 making decisions (National Research Council, 2009). What is sufficient for a peer-reviewed
792 scholarly publication in a specialized scientific journal may not be sufficient for a regulatory
793 decision. For example, in the EU, it is generally noted that “*available scientific and technical data*”
794 shall be taken into account “*in preparing its policy on the environment*” (Allio et al., 2006). It was
795 noted that this legislative mandate does not call for the available data to be of “best available”
796 quality, even though some of the agencies are required to “*provide the EU institutions and the*
797 *Member States with the best possible scientific opinions*” (Allio et al., 2006). In the US, the concept
798 of “best available science” has been more defined as a statutory requirement for risk evaluation of
799 chemical substances (US EPA, 2017). Specifically, the statute defines *best available science* as
800 “*science that is reliable and unbiased. Use of best available science involves the use of supporting*
801 *studies conducted in accordance with sound and objective science practices, including, when*
802 *available, peer reviewed science and supporting studies and data collected by accepted methods or*
803 *best available methods (if the reliability of the method and the nature of the decision justifies use of*
804 *the data). Additionally, EPA will consider as applicable: (1) The extent to which the scientific*
805 *information, technical procedures, measures, methods, protocols, methodologies, or models*
806 *employed to generate the information are reasonable for and consistent with the intended use of the*
807 *information; (2) The extent to which the information is relevant for the Administrator's use in*
808 *making a decision about a chemical substance or mixture; (3) The degree of clarity and*
809 *completeness with which the data, assumptions, methods, quality assurance, and analyses employed*
810 *to generate the information are documented; (4) The extent to which the variability and uncertainty*
811 *in the information, or in the procedures, measures, methods, protocols, methodologies, or models,*
812 *are evaluated and characterized; and (5) The extent of independent verification or peer review of*
813 *the information or of the procedures, measures, methods, protocols, methodologies or models.”* In
814 principle, all of these can apply to characterization of the chemical composition of petroleum
815 UVCBs when data are submitted to the authorities for a decision.

816 Related to the concept of *best available science* is the concept of “*reasonably available*
817 *information*” which means information that the agency “*possesses or can reasonably generate,*
818 *obtain, and synthesize for use in risk evaluations, considering the deadlines [...] for completing*
819 *such evaluation*” (US EPA, 2017). To relate the latter concept to the data on chemical composition
820 of petroleum UVCBs, one may argue that a range of technically valid analytical methods applicable
821 to the challenge of characterizing their composition and variability are already available and the data
822 could be generated in a reasonable period of time. There is a balance between cost, availability,
823 time, and the level of information required; the most advanced techniques which ultimately yield
824 the fullest possible understanding may not be suitable for routine screening, whilst more common
825 techniques may afford incomplete information or increased risk of misinterpretation. In this regard,
826 novel analytical advances, such as those described herein, have greatly increased our ability to
827 access valuable information on the composition of UVCBs at a molecular level. At the same time,
828 the advances have revealed in much greater detail the compositional complexity of petroleum and
829 areas where development is still required in order to fully access the range of components and, most
830 significantly, establish structures. Still, strong arguments have been made about the impractical and
831 potentially unnecessary regulatory requirements to deconvolute every possible constituent in a
832 UVCB, even if at some arbitrary abundance cut-off level(s) (**Table 1**), and that hazard assessment
833 may be sufficiently informed by using bulk compositional data and hydrocarbon blocks (Redman et
834 al., 2012; Salvito et al., 2020). Clearly, the fields of analytical chemistry and regulatory science have
835 not always worked in tandem.

836 To bridge this chasm, all sides will need to meet in the middle to establish baseline
837 requirements with an understanding of what is analytically possible (**Figure 4**), including
838 considerations of instrumentation availability, cost and time required for the analysis and data
839 processing. It has been well-documented that different analytical methods, when applied to
840 petroleum substances, have different windows of applicability with respect to the types of molecules
841 they may ionize and/or detect, and may suffer from being semi-quantitative (Aeppli, 2022;
842 Fernandez-Lima et al., 2009; McKenna et al., 2013; Prince and Walters, 2022; Rodgers and
843 McKenna, 2011). Indeed, **Figure 4** illustrates the point that novel high resolution analytical
844 methods (FT-ICR, Orbitrap-MS and TOF-MS) have been already applied by a number of authors
845 to the research questions that are directly relevant to the specific regulatory needs in petroleum
846 UVCB space. However, the number of such publications remains very small as compared to studies

847 of the traditional analytical approaches (**Figure 2**), and seldom do the authors acknowledge the
848 potential for their work to be applied in the regulatory context(s). It is clear that the analytical
849 chemists could better appreciate the regulatory issues and strive to provide solutions that not only
850 advance the science but also address specific challenges. For example, the focus on the critical needs
851 identified herein will serve well both the researchers and the ultimate end-users if targeted
852 collaborations and regulatory-informed case studies to build confidence in their performance and
853 utility. In addition, cross-validation collaborative trials are needed to achieve standardization of the
854 new techniques and analysis methods to build confidence among stakeholders, a relevant example
855 comes from the evolution of the analytical methods used for oil spill response (Faksness et al., 2002).

856 Concomitantly, the government bodies tasked with enforcing REACH and other relevant
857 regulations, while providing additional useful guidance about the details of chemical
858 characterization of petroleum UVCBs (ECHA, 2022), should have realistic expectations as to what
859 is achievable using even the most advanced analytical methods. The most comprehensive analytical
860 information may or may not be actionable in terms of hazard evaluation in absence of anchoring to
861 the toxicological data. In this regard, existing data on a handful of known hazardous components in
862 petroleum UVCBs, be it “priority PAHs” (ATSDR, 2005) or other constituents with existing human
863 hazard evaluations, are widely regarded as insufficient for assessment of the whole substance(s).
864 Thus, additional data from *in vitro* and other “bioactivity” data streams may be needed to determine
865 the toxicologically relevant compositional features and variability within and among substances and
866 categories (Grimm et al., 2016; House et al., 2022; House et al., 2021).

867 Concomitantly, the size and complexity of data sets pertaining to sample composition, which
868 are increasingly accompanied with orthogonal high-dimensional information such as
869 (eco)toxicological data, present challenges with respect to data handling and subsequent
870 interpretation. These comprehensive high-/multi-dimensional datasets will typically be subject to
871 dimensionality reduction and used to support grouping and/or classification of the individual
872 petroleum UVCBs. A recent study demonstrated that the choices of the data analysis and
873 visualization methods can not only potentially aid in the communication of “sufficient similarity”
874 among complex substances, but also yield different outcomes in terms of grouping and classification
875 (Onel et al., 2019). The proverbial “black box” of data processing and analysis may create “*a variety*
876 *of methodological and scientific concerns which mean that it is impossible to independently assess*
877 *the methods and results*” (ECHA, 2020a) in terms of reliably predicting the properties of the

878 substances that are being evaluated. Hence, it is important to involve not only the analytical chemists
879 when examining the data, but those in other scientific disciplines, such as statistics, mathematics,
880 and computer science, and ensure that the data processing and analysis methods are transparent, and
881 that the decision-makers are sufficiently familiar with the bioinformatics aspects of the information
882 presented to them. Such involvement can expedite assessments of which components are
883 particularly relevant for regulatory needs and may warrant greater focus and analytical precision.

884 Overall, while statutory and advisory language in the government agencies-produced
885 documents may seem clear, the analytical science may not be available to fully identify and quantify
886 all, or even most, components within complex petroleum UVCBs. It is being acknowledged that “*it*
887 *is not required to provide detailed structural information on all constituents of a UVCB substance,*
888 *but there must be sufficient characterization of constituents so as to demonstrate structural*
889 *similarity, and consecutively provide a basis for predicting the properties of the substance in read-*
890 *across*” (ECHA, 2020a). Even though additional clarifications have been recently made with respect
891 to what “sufficient” may mean in the context of petroleum UVCBs (ECHA, 2022), ultimately, the
892 flexibility is needed to consider new science and determine if it is “best available” and also fit for a
893 specific decision-context purpose so that the statutes and regulations are ultimately enforceable.
894 While adherence to the most common and established methods is understandable, the industry needs
895 to be more open to the application of the modern analytical chemistry methods to the analysis of the
896 samples and inclusion of such data into regulatory submissions. The traditional approaches that are
897 used for broad characterization of petroleum UVCBs still have an important role to play with respect
898 to substance identification, greater consideration of nontargeted approaches (and how these may
899 also shed new light on emerging hazards), followed by targeted approaches in some cases, is needed.
900 It is also important, however, to recognize that the measurement of the composition of a complex
901 substance is very much influenced by the methods and techniques used; comparisons of data from
902 different laboratories must demonstrate appreciation of such factors if the comparisons are to be
903 meaningful. Therefore, detailed characterization of the technical performance of novel analytical
904 methods is required, in addition to the generation of data for a wide range of samples and
905 applications.

906

907

908

909 **Declaration of Competing Interest**

910 The authors declare that they have no known competing financial interests or personal
911 relationships that could have appeared to influence the work reported in this paper.

912

913 **Acknowledgements**

914 The authors wish to thank Dr. Delina Lyon (Shell and CONCAWE) and Stuart Forbes
915 (CONCAWE) for critical suggestions on this manuscript's text and organization. A.T. Roman-
916 Hubers and A. C. Cordova were supported, in part, by a training grant from the National Institute of
917 Environmental Health Sciences (T32 ES026568).

918 **Table 1.** Summary of the chemical characterization needs for UVCB that are registered and
 919 evaluated under REACH in the European Union.

Information Needed	Substance Identification for <u>Registration</u> (REACH Annex VI, Section 2)	Hazard Characterization for <u>Evaluation</u> (REACH Annex XI, Section 1.5)		
		Read-Across Assessment Framework: Considerations on multi-constituent substances and UVCBs (ECHA, 2017c)		Advice on Using Read-Across for UVCB substances (ECHA, 2022)
		Category Approach	Analogue Approach	
Overall Substance Identification	<ul style="list-style-type: none"> Substance name [from a trade association and/or nomenclature system] Substance identifiers [source, manufacturing process, carbon and boiling range, phys-chem properties, etc.] 	<ul style="list-style-type: none"> Substance(s) to be grouped in a category Compositions to be included Manufacturing process description 	<ul style="list-style-type: none"> Source substance(s) and target substance Compositions to be included Manufacturing process description 	<ul style="list-style-type: none"> Similarity may be established based on: (i) presence of <i>identical</i> constituents OR (ii) variation in concentration and variability in constituents Constituents present at >1% must be identified; lower thresholds if constituents of concern are present (>0.1% for carcinogenic/mutagenic, >0.3% for repro/developmental) >80% of constituents in the substance must be identified
Composition Characterization	<ul style="list-style-type: none"> Constituents present >10% Constituents <10% that may be impacting hazard classification 	Category “domain” needs to be defined (constituent-specific concentration determination not specified)	<ul style="list-style-type: none"> Source substance(s) Target substance (constituent-specific concentration determination not specified) 	<ul style="list-style-type: none"> Must be characterized up to 100% When full characterization is impractical/impossible need to provide (i) justification AND (ii) demonstration of similarity “by other means”
Variability	Not Required	<ul style="list-style-type: none"> Structural similarity for category based on worst-case scenario Determine if quantitative differences or patterns in predicted properties may be reflected in structural similarity 	<ul style="list-style-type: none"> Structural Similarity between source & target based on worst-case scenario Determine if quantitative differences or patterns in predicted properties may be reflected in structural similarity 	<ul style="list-style-type: none"> Structural similarity explained based on quantitative and qualitative comparison of composition When full characterization is impractical/impossible need to provide (i) comparison of constituents AND (ii) demonstration of similarity “by other means” (e.g. analytical information for >95% constituents, constituents of high concern, high resolution for confidence fingerprinting) Analysis of at least 5 independent (i.e., production batches) samples analyzed from ALL registrants of a substance

920

921

References

- 922
923
- 924 Abib, G. A. P., et al., 2020. Assessing raw materials as potential adsorbents to remove acidic
925 compounds from Brazilian crude oils by ESI (-) FT-ICR MS. *An Acad Bras Cienc.* 92,
926 e20200214.
- 927 Aeppli, C., 2022. Recent advance in understanding photooxidation of hydrocarbons after oil spills.
928 *Curr Opin Chem Eng.* 36, 100763.
- 929 Alam, M. S., et al., 2016. Using Variable Ionization Energy Time-of-Flight Mass Spectrometry
930 with Comprehensive GCxGC To Identify Isomeric Species. *Anal Chem.* 88, 4211-20.
- 931 Alam, M. S., et al., 2018. Mapping and quantifying isomer sets of hydrocarbons
932 ($\geq C_{12}$) in diesel exhaust, lubricating oil and diesel fuel samples
933 using GC \times GC-ToF-MS. *Atmos Meas Tech.* 11, 3047-3058.
- 934 Allio, L., et al., 2006. Enhancing the role of science in the decision-making of the European
935 Union. *Regul Toxicol Pharmacol.* 44, 4-13.
- 936 API, Petroleum Process Stream Terms Included in the Chemical Substances Inventory Under the
937 Toxic Substances Control Act (TSCA). In: Control, H. S. R. c. T. F. o. T. S., (Ed.).
938 American Petroleum Institute, 1983.
- 939 ASTM International, UOP Method 990-11: Organic Analysis of Distillate by Comprehensive
940 Two-Dimensional Gas Chromatography with Flame Ionization Detection. ASTM
941 International, West Conshohocken, PA, 2011.
- 942 ATSDR, Toxicology profile for polyaromatic hydrocarbons. CRC Press Boca Raton City, FL,
943 2005.
- 944 Bae, E., et al., 2010. Identification of about 30 000 Chemical Components in Shale Oils by
945 Electrospray Ionization (ESI) and Atmospheric Pressure Photoionization (APPI) Coupled
946 with 15 T Fourier Transform Ion Cyclotron Resonance Mass Spectrometry (FT-ICR MS)
947 and a Comparison to Conventional Oil. *Energ Fuel.* 24, 2563-2569.
- 948 Ball, G. I., Aluwihare, L. I., 2014. CuO-oxidized dissolved organic matter (DOM) investigated
949 with comprehensive two dimensional gas chromatography-time of flight-mass
950 spectrometry (GC \times GC-TOF-MS). *Org Chem.* 75, 87-98.
- 951 Barrow, M. P., et al., 2014. An added dimension: GC atmospheric pressure chemical ionization
952 FTICR MS and the Athabasca oil sands. *Anal Chem.* 86, 8281-8.
- 953 Barrow, M. P., et al., 2016. Effects of Extraction pH on the Fourier Transform Ion Cyclotron
954 Resonance Mass Spectrometry Profiles of Athabasca Oil Sands Process Water. *Energy &*
955 *Fuels.* 30, 3615-3621.
- 956 Benassi, M., et al., 2013. Petroleum crude oil analysis using low-temperature plasma mass
957 spectrometry. *Rapid Commun Mass Spectrom.* 27, 825-34.

- 958 Bierkens, J., Geerts, L., 2014. Environmental hazard and risk characterisation of petroleum
959 substances: a guided "walking tour" of petroleum hydrocarbons. *Environ Int.* 66, 182-93.
- 960 Bowman, D. T., et al., 2019. Profiling of individual naphthenic acids at a composite tailings
961 reclamation fen by comprehensive two-dimensional gas chromatography-mass
962 spectrometry. *Sci Total Environ.* 649, 1522-1531.
- 963 Bowman, D. T., et al., 2020. Isomer-specific monitoring of naphthenic acids at an oil sands pit
964 lake by comprehensive two-dimensional gas chromatography-mass spectrometry. *Sci Total*
965 *Environ.* 746, 140985.
- 966 Brockman, S. A., et al., 2018. Van Krevelen diagram visualization of high resolution-mass
967 spectrometry metabolomics data with OpenVanKrevelen. *Metabolomics.* 14, 48.
- 968 Cao, X. H., et al., 2020. Sequential thermal dissolution of two low-rank coals and characterization
969 of their structures by high-performance liquid chromatography/time-of-flight mass
970 spectrometry and gas chromatography/mass spectrometry. *Rapid Commun in Mass*
971 *Spectrom.* 34, e8887.
- 972 Castiblanco, J. E. B., et al., 2020. Molecular behavior assessment on initial stages of oil spill in
973 terrestrial environments. *Environ Sci Pollut Res.* 28, 13595-13604.
- 974 Chen, X., et al., 2018. Separation and Molecular Characterization of Ketones in a Low-
975 Temperature Coal Tar. *Energy Fuels.* 32, 4662-4670.
- 976 Chen, X. B., et al., 2012. Characterization and Comparison of Nitrogen Compounds in
977 Hydrotreated and Untreated Shale Oil by Electrospray Ionization (ESI) Fourier Transform
978 Ion Cyclotron Resonance Mass Spectrometry (FT-ICR MS). *Energy & Fuels.* 26, 1707-
979 1714.
- 980 Cheng, X., Hous, D., 2021. Characterization of Severely Biodegraded Crude Oils Using Negative-
981 Ion ESI Orbitrap MS, GC-NCD and GC-SCD: Insights into Heteroatomic Compounds
982 Biodegradation. *Energies.* 14.
- 983 Cho, Y., et al., 2017. Extension of the Analytical Window for Characterizing Aromatic
984 Compounds in Oils Using a Comprehensive Suite of High-Resolution Mass Spectrometry
985 Techniques and Double Bond Equivalence versus Carbon Number Plot. *Energy & Fuels.*
986 31, 7874-7883.
- 987 Cho, Y., et al., 2011. Planar limit-assisted structural interpretation of
988 saturates/aromatics/resins/asphaltenes fractionated crude oil compounds observed by
989 Fourier transform ion cyclotron resonance mass spectrometry. *Anal Chem.* 83, 6068-73.
- 990 Cho, Y., et al., 2014. Evaluation of Laser Desorption Ionization Coupled to Fourier Transform Ion
991 Cyclotron Resonance Mass Spectrometry To Study Metalloporphyrin Complexes. *Energy*
992 *& Fuels.* 28, 6699-6706.
- 993 Clark, C. R., et al., 2013. A GHS-consistent approach to health hazard classification of petroleum
994 substances, a class of UVCB substances. *Regul Toxicol Pharmacol.* 67, 409-20.

- 995 CONCAWE, REACH – Analytical characterisation of petroleum UVCB substances. Brussels,
996 Belgium, 2012.
- 997 CONCAWE, Guidance on Reporting Analytical Information for Petroleum Substances in REACH
998 Registration Dossiers V2. Brussels, 2014.
- 999 CONCAWE, Concawe Substance Identification Group Analytical Program Report (Abridged
1000 Version). Brussels, Belgium, 2019.
- 1001 CONCAWE, Guidance to Registrants on Methods for Characterisation of Petroleum UVCB
1002 Substances for REACH Registration Purposes. Brussels, Belgium, 2020.
- 1003 CONCAWE, Petroleum Substances and REACH. 2021.
- 1004 CONCAWE, REACH Background. 2022.
- 1005 Damasceno, F. C., et al., 2014. Characterization of naphthenic acids using mass spectroscopy and
1006 chromatographic techniques: study of technical mixtures. *Anal Methods*. 6, 807-816.
- 1007 Dimitrov, S. D., et al., 2015. UVCB substances: methodology for structural description and
1008 application to fate and hazard assessment. *Environ Toxicol Chem*. 34, 2450-62.
- 1009 Dodds, J. N., Baker, E. S., 2019. Ion Mobility Spectrometry: Fundamental Concepts,
1010 Instrumentation, Applications, and the Road Ahead. *J Am Soc Mass Spectrom*. 30, 2185-
1011 2195.
- 1012 Dong, X., et al., 2019. Evaluation of elemental composition obtained by using mass spectrometer
1013 and elemental analyzer: A case study on model compound mixtures and a coal-derived
1014 liquid. *Fuel*. 245, 392-397.
- 1015 ECHA, Guidance on information requirements and chemical safety assessment. Chapter R.6:
1016 QSARS and grouping of chemicals. 2008.
- 1017 ECHA, Decision on a testing proposal set out in a registration pursuant to Article 40(3) of
1018 regulation (EC) No 1907/2006 for Asphalt, CAS No 8052-42-4 (EC No 232-490-9).
1019 European Chemical Agency, Helsinki, Finland, 2014.
- 1020 ECHA, Guidance for identification and naming of substances under REACH and CLP. Vol. 2.1.
1021 European Chemical Agency, Helsinki, Finland, 2017a.
- 1022 ECHA, Guidance on requirements for substances in articles. Vol. Version 4.0. European
1023 Chemicals Agency, Helsinki, Finland, 2017b.
- 1024 ECHA, Read-Across Assessment Framework (RAAF) - considerations on multi-constituent
1025 substances and UVCBs. European Chemical Agency, Helsinki, Finland, 2017c.
- 1026 ECHA, Testing Proposal Decision on Substance EC 295-332-8 "Extracts (petroleum), deasphalted
1027 vacuum residue solvent". European Chemicals Agency, Helsinki, Finland, 2020a.

- 1028 ECHA, Testing Proposals Decision on Substance EC 265-110-5 "Extracts (petroleum), residual
1029 oil solvent". 2020b.
- 1030 ECHA, Testing Proposal Decision on Substance EC 265-182-8 "Gas oils (petroleum),
1031 hydrodesulfurized". European Chemicals Agency, Helsinki, Finland, 2021.
- 1032 ECHA, Advice on using read-across for UVCB substances - obligations arising from Commission
1033 Regulation 2021/979, amending REACH annexes. European Chemicals Agency, Helsinki,
1034 Finland, 2022.
- 1035 EPA, U. S., Toxic Substances Control Act Inventory Representation for Chemical Substances of
1036 Unknown or Variable Composition, Complex Reaction Products and Biological Materials:
1037 UVCB Substances. 1995.
- 1038 European Commission, COMMISSION REGULATION (EU) 2021/979 of 17 June 2021
1039 amending Annexes VII to XI to Regulation (EC) No 1907/2006 of the European
1040 Parliament and of the Council concerning the Registration, Evaluation, Authorisation and
1041 Restriction of Chemicals (REACH). Official Journal of the European Union, Brussels,
1042 Belgium, 2021.
- 1043 European Council, Corrigendum to Regulation (EC) No 1907/2006 of the European Parliament
1044 and of the Council of 18 December 2006 concerning the Registration, Evaluation,
1045 Authorisation and Restriction of Chemicals (REACH), establishing a European Chemicals
1046 Agency, amending Directive 1999/45/EC and repealing Council Regulation (EEC) No
1047 793/93 and Commission Regulation (EC) No 1488/94 as well as Council Directive
1048 76/769/EEC and Commission Directives 91/155/EEC, 93/67/EEC, 93/105/EC and
1049 2000/21/EC. Vol. L 136/3. European Commission, Brussels, Belgium, 2007.
- 1050 Faksness, L. G., et al., 2002. Round Robin Study—Oil Spill Identification. *Environ Forensics*. 3,
1051 279-291.
- 1052 Fernandez-Lima, F. A., et al., 2009. Petroleum crude oil characterization by IMS-MS and FTICR
1053 MS. *Anal Chem*. 81, 9941-7.
- 1054 Folkerts, E. J., et al., 2019. Toxicity in aquatic model species exposed to a temporal series of three
1055 different flowback and produced water samples collected from a horizontal hydraulically
1056 fractured well. *Ecotoxicol Environ Saf*. 180, 600-609.
- 1057 Frenzel, M., et al., 2010. Complications with remediation strategies involving the biodegradation
1058 and detoxification of recalcitrant contaminant aromatic hydrocarbons. *Sci Total Environ*.
1059 408, 4093-4101.
- 1060 FuelsEurope, FuelsEurope position paper on REACH and the Refining industry. Brussels,
1061 Belgium, 2015.
- 1062 Gabetti, E., et al., 2021. Chemical fingerprinting strategies based on comprehensive two-
1063 dimensional gas chromatography combined with gas chromatography-olfactometry to

- 1064 capture the unique signature of Piemonte peppermint essential oil (*Mentha x piperita* var
1065 *Italo-Mitcham*). *J Chromatogr A*. 1645, 462101.
- 1066 Gavard, R., et al., 2020. KairosMS: A New Solution for the Processing of Hyphenated Ultrahigh
1067 Resolution Mass Spectrometry Data. *Anal Chem*. 92, 3775-3786.
- 1068 Giles, H. N., 2016. Crude Oil Analysis: History and Development of Test Methods From 1854 to
1069 2016. *Mater Perform Charact*. 5, 1-169.
- 1070 Gough, M. A., Rowland, S. J., 1990. Characterization of unresolved complex mixtures of
1071 hydrocarbons in petroleum. *Nature*. 344, 648-650.
- 1072 Grimm, F. A., et al., 2016. A chemical-biological similarity-based grouping of complex
1073 substances as a prototype approach for evaluating chemical alternatives. *Green Chem*. 18,
1074 4407-4419.
- 1075 Haitao, S., et al., 2013. Hydrocarbon composition of different VGO feedstocks and the correlation
1076 with FCC product distribution. *China Petrol Proc Petrochem Techn*. 15.
- 1077 Hao, J., et al., 2017. Thermal cracking behaviors and products distribution of oil sand bitumen by
1078 TG-FTIR and Py-GC/TOF-MS. *Energy Convers Manag*. 151, 227-239.
- 1079 Headley, J. V., et al., 2011. Preliminary fingerprinting of Athabasca oil sands polar organics in
1080 environmental samples using electrospray ionization Fourier transform ion cyclotron
1081 resonance mass spectrometry. *Rapid Commun Mass Spectrom*. 25, 1899-909.
- 1082 Hoskins, J. N., et al., 2011. Architectural Differentiation of Linear and Cyclic Polymeric Isomers
1083 by Ion Mobility Spectrometry-Mass Spectrometry. *Macromolecules*. 44, 6915-6918.
- 1084 Hosseini, S. H., et al., 2021. Characterization of crude oils derived from carbonate and siliciclastic
1085 source rocks using FTICR-MS. *Org Geochem*. 159, 104286.
- 1086 Hourani, N., et al., 2013. Atmospheric pressure chemical ionization Fourier transform ion
1087 cyclotron resonance mass spectrometry for complex thiophenic mixture analysis. *Rapid
1088 Commun Mass Spectrom*. 27, 2432-8.
- 1089 House, J. S., et al., 2022. Grouping of UVCB substances with dose-response transcriptomics data
1090 from human cell-based assays. *ALTEX*. 39, 388-404.
- 1091 House, J. S., et al., 2021. Grouping of UVCB substances with new approach methodologies
1092 (NAMs) data. *ALTEX*. 38, 123-137.
- 1093 Hsu, C. S., et al., 2011. Petroleomics: advanced molecular probe for petroleum heavy ends. *J Mass
1094 Spectrom*. 46, 337-43.
- 1095 Hsu, C. S., Shi, Q., 2013. Prospects for petroleum mass spectrometry and chromatography. *Sci
1096 China Chem*. 56, 833-839.

- 1097 Hu, M., et al., 2018. Collision cross section (CCS) measurement by ion cyclotron resonance mass
1098 spectrometry with short-time Fourier transform. *Rapid Commun Mass Spectrom.* 32, 751-
1099 761.
- 1100 Hughey, C. A., et al., 2001. Kendrick mass defect spectrum: a compact visual analysis for
1101 ultrahigh-resolution broadband mass spectra. *Anal Chem.* 73, 4676-81.
- 1102 Hughey, C. A., et al., 2002. Resolution of 11,000 compositionally distinct components in a single
1103 electrospray ionization Fourier transform ion cyclotron resonance mass spectrum of crude
1104 oil. *Anal Chem.* 74, 4145-9.
- 1105 Jaggi, A., et al., 2019. Composition of the dissolved organic matter produced during in situ
1106 burning of spilled oil. *Org Geochem.* 138, 103926.
- 1107 Jennerwein, M. K., et al., 2014. Complete Group-Type Quantification of Petroleum Middle
1108 Distillates Based on Comprehensive Two-Dimensional Gas Chromatography Time-of-
1109 Flight Mass Spectrometry (GC×GC-TOFMS) and Visual Basic Scripting. *Energy Fuels.*
1110 28, 5670-5681.
- 1111 Kaiser, M. J., 2017. A review of refinery complexity applications. *Pet Sci.* 14, 167-194.
- 1112 Kendrick, E., 1963. A Mass Scale Based on $CH_2 = 14.0000$ for High Resolution Mass
1113 Spectrometry of Organic Compounds. *Anal. Chem.* 35, 2146-2154.
- 1114 Kew, W., et al., 2017. Interactive van Krevelen diagrams - Advanced visualisation of mass
1115 spectrometry data of complex mixtures. *Rapid Commun Mass Spectrom.* 31, 658-662.
- 1116 Kim, D., et al., 2015. Combination of ring type HPLC separation, ultrahigh-resolution mass
1117 spectrometry, and high field NMR for comprehensive characterization of crude oil
1118 compositions. *Fuel.* 157, 48-55.
- 1119 Kim, S., et al., 2003. Graphical method for analysis of ultrahigh-resolution broadband mass
1120 spectra of natural organic matter, the van Krevelen diagram. *Anal Chem.* 75, 5336-44.
- 1121 Kondyli, A., Schrader, W., 2019. High-resolution GC/MS studies of a light crude oil fraction. *J*
1122 *Mass Spectrom.* 54, 47-54.
- 1123 Kulkarni, S. U., Thies, M. C., 2012. Quantitative analysis of polydisperse systems via solvent-free
1124 matrix-assisted laser desorption/ionization time-of-flight mass spectrometry. *Rapid*
1125 *Commun Mass Spectrom.* 26, 392-398.
- 1126 Lacroix-Andrivet, O., et al., 2021. Molecular Characterization of Aged Bitumen with Selective
1127 and Nonselective Ionization Methods by Fourier Transform Ion Cyclotron Resonance
1128 Mass Spectrometry. 2. Statistical Approach on Multiple-Origin Samples. *Energy & Fuels.*
1129 35, 16442-16451.
- 1130 Lai, A., et al., 2022. The Next Frontier of Environmental Unknowns: Substances of Unknown or
1131 Variable Composition, Complex Reaction Products, or Biological Materials (UVCBs).
1132 *Environ Sci Technol.* 56, 7448-7466.

- 1133 Lalli, P. M., et al., 2015. Isomeric Separation and Structural Characterization of Acids in
1134 Petroleum by Ion Mobility Mass Spectrometry. *Energy Fuels*. 29, 3626-3633.
- 1135 Lalli, P. M., et al., 2017. Functional Isomers in Petroleum Emulsion Interfacial Material Revealed
1136 by Ion Mobility Mass Spectrometry and Collision-Induced Dissociation. *Energy & Fuels*.
1137 31, 311-318.
- 1138 Leefmann, T., et al., 2019. UltraMassExplorer: a browser-based application for the evaluation of
1139 high-resolution mass spectrometric data. *Rapid Commun Mass Spectrom*. 33, 193-202.
- 1140 Li, S., et al., 2022a. Molecular characteristics of biochar-derived organic matter sub-fractions
1141 extracted by ultrasonication. *Sci Total Environ*. 806, 150190.
- 1142 Li, Y., et al., 2022b. Comprehensive chemical characterization of dissolved organic matter in
1143 typical point-source refinery wastewaters. *Chemosphere*. 286, 131617.
- 1144 Liu, D., et al., 2014. Direct hydro-liquefaction of sawdust in petroleum ether and comprehensive
1145 bio-oil products analysis. *Bioresour Technol*. 155, 152-60.
- 1146 Liu, Y., et al., 2020. The acid and neutral nitrogen compounds characterized by negative ESI
1147 Orbitrap MS in a heavy oil before and after oxidation. *Fuel*. 277.
- 1148 Lu, H., et al., 2013. Geochemical Explication of Sulfur Organics Characterized by Fourier
1149 Transform Ion Cyclotron Resonance Mass Spectrometry on Sulfur-Rich Heavy Oils in
1150 Jinxian Sag, Bohai Bay Basin, Northern China. *Energy Fuels*. 27, 5861-5866.
- 1151 Luna, N., et al., 2014a. Identification And Characterization Of Sulfur Compounds In Straight-Run
1152 Diesel Using Comprehensive Two-Dimensional Gas Chromatography Coupled To Time-
1153 Of-Flight Mass Spectrometry. *China Pet Process Petrochem*. 16, 10-18.
- 1154 Luna, N., et al., 2014b. Identification and Characterization of Sulfur Compounds in Straight-Run
1155 Diesel Using Comprehensive Two-Dimensional GC Coupled with TOF MS. *China Petrol
1156 Proc Petrochem Techn*. 16, 10-18.
- 1157 Luo, T., et al., 2016. A novel characterization of furfural-extract oil from vacuum gas oil and its
1158 application in solvent extraction process. *Fuel Process Technol*. 152, 356-366.
- 1159 Lv, H., et al., 2013. Application of UPLC-Quadrupole-TOF-MS Coupled with Recycling
1160 Preparative HPLC in Isolation and Preparation of Coumarin Isomers with Similar Polarity
1161 from *Peucedanum praeruptorum*. *Chromatographia*. 76, 141-148.
- 1162 MacLennan, M. S., et al., 2018. Characterization of Athabasca lean oil sands and mixed surficial
1163 materials: Comparison of capillary electrophoresis/low-resolution mass spectrometry and
1164 high-resolution mass spectrometry. *Rapid Commun Mass Spectrom*. 32, 695-702.
- 1165 Mahmoud, M. Y., Dabek-Zlotorzynska, E., 2018. Investigation of isomeric structures in a
1166 commercial mixture of naphthenic acids using ultrahigh pressure liquid chromatography
1167 coupled to hybrid traveling wave ion mobility-time of flight mass spectrometry. *J.
1168 Chromatogr. A*. 1572, 90-99.

- 1169 Maillard, J. F., et al., 2021. Structural analysis of petroporphyrins from asphaltene by trapped ion
1170 mobility coupled with Fourier transform ion cyclotron resonance mass spectrometry.
1171 *Analyst*. 146, 4161-4171.
- 1172 Mao, D., et al., 2009. Combining HPLC-GCXGC, GCXGC/ToF-MS, and selected ecotoxicity
1173 assays for detailed monitoring of petroleum hydrocarbon degradation in soil and leaching
1174 water. *Environ Sci Technol*. 43, 7651-7.
- 1175 Marshall, A. G., Rodgers, R. P., 2004. Petroleomics: the next grand challenge for chemical
1176 analysis. *Acc Chem Res*. 37, 53-9.
- 1177 Marshall, A. G., Rodgers, R. P., 2008. Petroleomics: chemistry of the underworld. *Proc Natl Acad
1178 Sci U S A*. 105, 18090-5.
- 1179 McKee, R. H., et al., 2015. Characterization of the toxicological hazards of hydrocarbon solvents.
1180 *Crit Rev Toxicol*. 45, 273-365.
- 1181 McKenna, A. M., et al., 2013. Expansion of the analytical window for oil spill characterization by
1182 ultrahigh resolution mass spectrometry: beyond gas chromatography. *Environ Sci Technol*.
1183 47, 7530-9.
- 1184 Mennito, A. S., Qian, K., 2013. Characterization of Heavy Petroleum Saturates by Laser
1185 Desorption Silver Cationization and Fourier Transform Ion Cyclotron Resonance Mass
1186 Spectrometry. *Energy & Fuels*. 27, 7348-7353.
- 1187 Miles, S. M., et al., 2020. Oil sands process affected water sourced *Trichoderma harzianum*
1188 demonstrates capacity for mycoremediation of naphthenic acid fraction compounds.
1189 *Chemosphere*. 258, 127281.
- 1190 Muller, H., et al., 2020. Innate Sulfur Compounds as an Internal Standard for Determining
1191 Vacuum Gas Oil Compositions by APPI FT-ICR MS. *Energy Fuels*. 34, 8260-8273.
- 1192 National Research Council, 2009. *Science and Decisions: Advancing Risk Assessment*. National
1193 Academies Press, Washington, DC.
- 1194 Neumann, A., et al., 2021. Investigation of Island/Single-Core- and Archipelago/Multicore-
1195 Enriched Asphaltenes and Their Solubility Fractions by Thermal Analysis Coupled with
1196 High-Resolution Fourier Transform Ion Cyclotron Resonance Mass Spectrometry. *Energy
1197 & Fuels*. 35, 3808-3824.
- 1198 Ngo, H. I., et al., 2012. Improved synthesis and characterization of saturated branched-chain fatty
1199 acid isomers. *Eur J Lipid Sci Technol*. 114, 213-221.
- 1200 Niyonsaba, E., et al., 2019. Recent Advances in Petroleum Analysis by Mass Spectrometry. *Anal
1201 Chem*. 91, 156-177.
- 1202 Nolte, T., et al., 2013. Desulfurized Fuels from Athabasca Bitumen and Their Polycyclic Aromatic
1203 Sulfur Heterocycles. Analysis Based on Capillary Electrophoresis Coupled with TOF MS.
1204 *Energy Fuels*. 27, 97-107.

- 1205 O'Reilly, K. T., et al., 2019. Oxygen-Containing Compounds Identified in Groundwater from Fuel
1206 Release Sites Using GCxGC-TOF-MS. *Ground Water Monit R.* 39, 32-40.
- 1207 Oldenburg, T. B. P., et al., 2014. The impact of thermal maturity level on the composition of crude
1208 oils, assessed using ultra-high resolution mass spectrometry. *Org Geochem.* 75, 151-168.
- 1209 Oldenburg, T. B. P., et al., 2017. The controls on the composition of biodegraded oils in the deep
1210 subsurface – Part 4. Destruction and production of high molecular weight non-hydrocarbon
1211 species and destruction of aromatic hydrocarbons during progressive in-reservoir
1212 biodegradation. *Org Geochem.* 114, 57-80.
- 1213 Onel, M., et al., 2019. Grouping of complex substances using analytical chemistry data: A
1214 framework for quantitative evaluation and visualization. *PLoS One.* 14, e0223517.
- 1215 Orrego- Ruíz, J. A., 2018. Finding a relationship between the composition and the emulsifying
1216 character of asphaltenes through FTICR-MS. *Cienc Tecn Fut.* 8, 45-52.
- 1217 Palacio Lozano, D. C., et al., 2019a. Chapter 32 | Mass Spectrometry in the Petroleum Industry,"
1218 in *Fuels and Lubricants Handbook: Technology, Properties, Performance, and Testing.*
1219 *MNL37-2ND-EB Fuels and Lubricants Handbook: Technology, Properties, Performance,*
1220 *and Testing.* 2.
- 1221 Palacio Lozano, D. C., et al., 2019b. Pushing the analytical limits: new insights into complex
1222 mixtures using mass spectra segments of constant ultrahigh resolving power. *Chem Sci.*
1223 10, 6966-6978.
- 1224 Palacio Lozano, D. C., et al., 2022. Revealing the Reactivity of Individual Chemical Entities in
1225 Complex Mixtures: the Chemistry Behind Bio-Oil Upgrading. *Anal Chem.* 94, 7536-7544.
- 1226 Palacio Lozano, D. C., et al., 2020. *Petroleomics: Tools, Challenges, and Developments.* *Annu*
1227 *Rev Anal Chem (Palo Alto Calif).* 13, 405-430.
- 1228 Pereira, A. S., et al., 2013. Characterization of oil sands process-affected waters by liquid
1229 chromatography orbitrap mass spectrometry. *Environ Sci Technol.* 47, 5504-13.
- 1230 Petroleum HPV Testing Group, HPV Challenge Overview. API, 2017.
- 1231 Phillips, A. L., et al., 2022. A Framework for Utilizing High-Resolution Mass Spectrometry and
1232 Nontargeted Analysis in Rapid Response and Emergency Situations. *Environ Toxicol*
1233 *Chem.* 41, 1117-1130.
- 1234 Porto, C. F. C., et al., 2019. Characterization of organosulfur compounds in asphalt cement
1235 samples by ESI(+)-FT-ICR MS and ¹³C NMR spectroscopy. *Fuel.* 256, 115923.
- 1236 Prince, R. C., Walters, C. C., 2022. Modern analytical techniques are improving our ability to
1237 follow the fate of spilled oil in the environment. *Curr Opin Chem Eng.* 36, 100787.

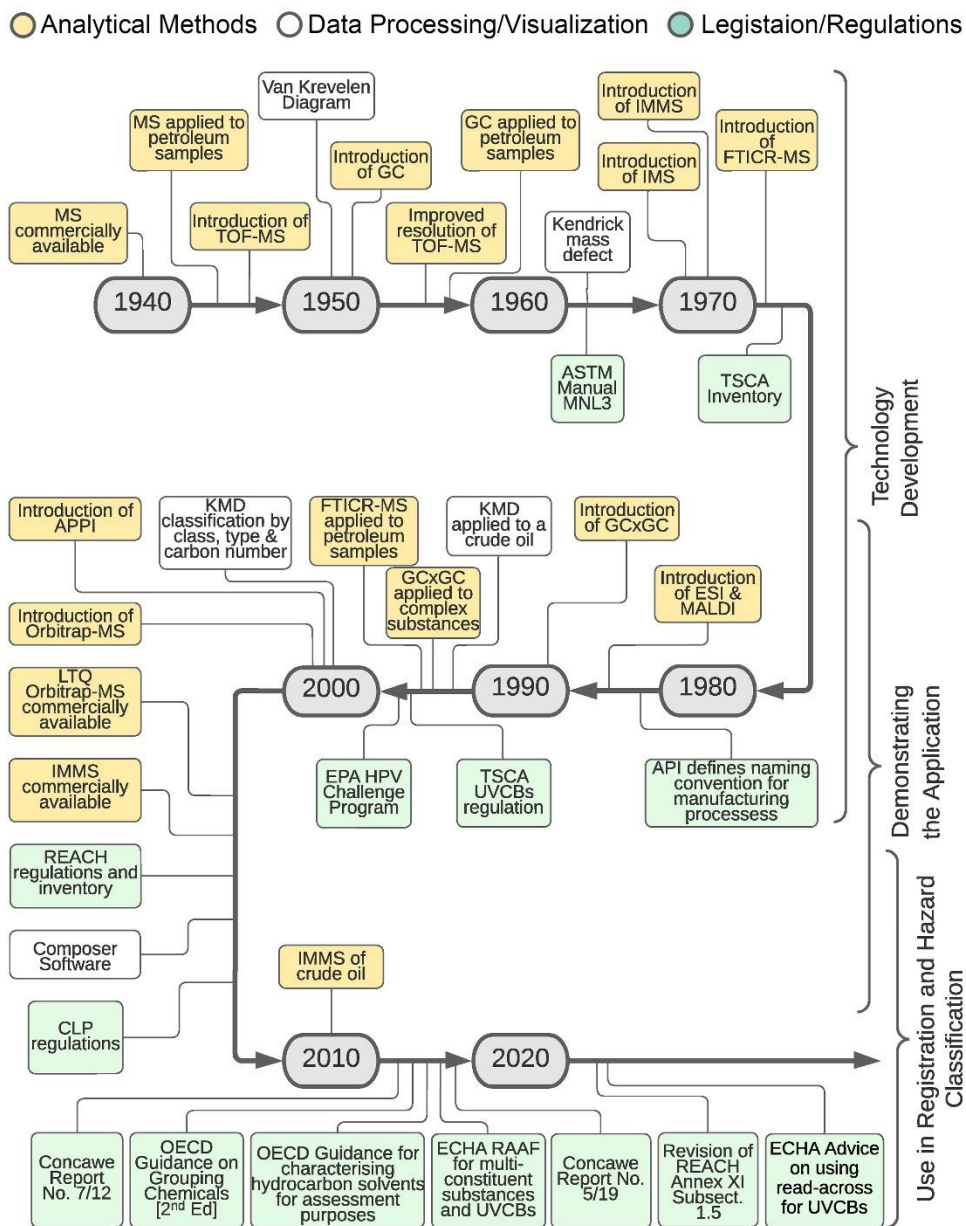
- 1238 Qian, K., et al., 2004. The Coupling of Supercritical Fluid Chromatography and Field Ionization
1239 Time-of-Flight High-Resolution Mass Spectrometry for Rapid and Quantitative Analysis
1240 of Petroleum Middle Distillates. *Eur J Mass Spectrom.* 10, 187-196.
- 1241 Qian, K., et al., 2001. Reading chemical fine print: Resolution and identification of 3000 nitrogen-
1242 containing aromatic compounds from a single electrospray ionization Fourier transform
1243 ion cyclotron resonance mass spectrum of heavy petroleum crude oil. *Energy & Fuels.* 15,
1244 492-498.
- 1245 Qian, K., Wang, F. C., 2019. Compositional Analysis of Heavy Petroleum Distillates by
1246 Comprehensive Two-dimensional Gas Chromatography, Field Ionization and High-
1247 resolution Mass Spectrometry. *J Am Soc Mass Spectrom.* 30, 2785-2794.
- 1248 Rasmussen, K., et al., 1999. Compilation of EINECS: Descriptions and definitions used for UVCB
1249 substances: Complex reaction products, plant products, (post-reacted) naturally occurring
1250 substances, micro-organisms, petroleum products, soaps and detergents, and metallic
1251 compounds. *Toxicol Environ Chem.* 69, 403-416.
- 1252 Reddy, C. M., Quinn, J. G., 1999. GC-MS analysis of total petroleum hydrocarbons and
1253 polycyclic aromatic hydrocarbons in seawater samples after the North Cape oil spill. *Mar
1254 Pollut Bull.* 38, 126-135.
- 1255 Redman, A. D., et al., 2014. PETRORISK: a risk assessment framework for petroleum substances.
1256 *Integr Environ Assess Manag.* 10, 437-48.
- 1257 Redman, A. D., et al., 2012. PETROTOX: an aquatic toxicity model for petroleum substances.
1258 *Environ Toxicol Chem.* 31, 2498-506.
- 1259 Riches, E., et al., *Ion Mobility & PetroOrg Software: Novel Techniques for Petroleomics
1260 Investigations In: Waters, (Ed.), 2015.*
- 1261 Ristic, N. D., et al., 2018. Compositional Characterization of Pyrolysis Fuel Oil from Naphtha and
1262 Vacuum Gas Oil. *Energy Fuels.* 32, 1276-1286.
- 1263 Rocha, Y. D. S., et al., 2018. Geochemical characterization of lacustrine and marine oils from off-
1264 shore Brazilian sedimentary basins using negative-ion electrospray Fourier transform ion
1265 cyclotron resonance mass spectrometry (ESI FTICR-MS). *Org Geochem.* 124, 29-45.
- 1266 Rodgers, R. P., McKenna, A. M., 2011. Petroleum analysis. *Anal Chem.* 83, 4665-87.
- 1267 Rodgers, R. P., et al., 1998. Resolution, elemental composition, and simultaneous monitoring by
1268 Fourier transform ion cyclotron resonance mass spectrometry of organosulfur species
1269 before and after diesel fuel processing. *Anal Chem.* 70, 4743-4750.
- 1270 Rodrigues Covas, T., et al., 2020. Fractionation of polar compounds from crude oils by hetero-
1271 medium pressure liquid chromatography (H-MPLC) and molecular characterization by
1272 ultrahigh-resolution mass spectrometry. *Fuel.* 267, 117289.

- 1273 Roman-Hubers, A. T., et al., 2021. Data Processing Workflow to Identify Structurally Related
1274 Compounds in Petroleum Substances Using Ion Mobility Spectrometry-Mass
1275 Spectrometry. *Energy Fuels*. 35, 10529-10539.
- 1276 Roman-Hubers, A. T., et al., 2022. Characterization of Compositional Variability in Petroleum
1277 Substances. *Fuel*. 317, 123547.
- 1278 Rosa, T. R., et al., 2020. DropMS: Petroleomics Data Treatment Based in Web Server for High-
1279 Resolution Mass Spectrometry. *J Am Soc Mass Spectrom*. 31, 1483-1490.
- 1280 Rowland, S. J., et al., 2011. Identification of individual acids in a commercial sample of
1281 naphthenic acids from petroleum by two-dimensional comprehensive gas
1282 chromatography/mass spectrometry. *Rapid Commun Mass Spectrom*. 25, 1741-1751.
- 1283 Ruger, C. P., et al., 2015. Hyphenation of Thermal Analysis to Ultrahigh-Resolution Mass
1284 Spectrometry (Fourier Transform Ion Cyclotron Resonance Mass Spectrometry) Using
1285 Atmospheric Pressure Chemical Ionization For Studying Composition and Thermal
1286 Degradation of Complex Materials. *Anal Chem*. 87, 6493-9.
- 1287 Rui, D., et al., 2012. Molecular Characterization of Hydrotreated Atmospheric Residue Derived
1288 from Arabian Heavy Crude by GC FI/FD TOF MS and APPI FT-ICR MS. *China Pet
1289 Process Petrochem*. 14, 80-88.
- 1290 Salvito, D., et al., 2020. Improving the Environmental Risk Assessment of Substances of
1291 Unknown or Variable Composition, Complex Reaction Products, or Biological Materials.
1292 *Environ Toxicol Chem*. 39, 2097-2108.
- 1293 Santos, J. M., et al., 2015. Petroleomics by ion mobility mass spectrometry: resolution and
1294 characterization of contaminants and additives in crude oils and petrofuels. *Anal Methods*.
1295 7, 4450-4463.
- 1296 Scarlett, A., et al., 2008. Chronic sublethal effects associated with branched alkylbenzenes
1297 bioaccumulated by mussels. *Environ Toxicol Chem*. 27, 561-567.
- 1298 Sierra Analytics, Composer: State of the art visualization software for petroleomics. 2022.
- 1299 Silva, R. C., et al., 2020. Mechanistic insights into sulfur rich oil formation, relevant to geological
1300 carbon storage routes. A study using (+) APPI FTICR-MS analysis. *Org Geochem*. 147.
- 1301 Silva, R. V. S., et al., 2019. Comprehensive study of the liquid products from slow pyrolysis of
1302 crambe seeds: Bio-oil and organic compounds of the aqueous phase. *Biomass Bioenerg*.
1303 123, 78-88.
- 1304 Smith, H. M., et al., 1959. Keys to the mystery of crude oil. *Proc. Amer. Petrole. Inst*. 39, 433-
1305 465.
- 1306 Sorensen, L., et al., 2019. Establishing a link between composition and toxicity of offshore
1307 produced waters using comprehensive analysis techniques - A way forward for discharge
1308 monitoring? *Sci Total Environ*. 694, 133682.

- 1309 Spectroswiss, FTMS Data Processing Tools PeakByPeak FTMS Data Analysis. In: Spectroswiss,
1310 (Ed.), 2019.
- 1311 Stenson, A. C., et al., 2003. Exact masses and chemical formulas of individual Suwannee River
1312 fulvic acids from ultrahigh resolution electrospray ionization Fourier transform ion
1313 cyclotron resonance mass spectra. *Anal Chem.* 75, 1275-84.
- 1314 Stout, S. A., Wang, Z., 2007. Chemical fingerprinting of spilled or discharged petroleum —
1315 methods and factors affecting petroleum fingerprints in the environment. in: Wang, Z.,
1316 Stout, S. A., (Eds.), *Oil Spill Environmental Forensics: Fingerprinting And Source*
1317 *Identification*. Academic Press, Cambridge, MA, pp. 1-53.
- 1318 U.S. EPA, Toxic Substance Control Act (TSCA) PL 94-469 Candidate List of Chemical
1319 Substances Addendum 1 Generic Terms Covering Petroleum Refinery Process Streams.
1320 US Environmental Protection Agency,, Washignton, D.C., 1978.
- 1321 Urban, P. L., 2016. Quantitative mass spectrometry: an overview. *Philos Trans A Math Phys Eng*
1322 *Sci.* 374, 20150382.
- 1323 US EPA, Method 8051B: Non-halogenated organics using GC-FID. US Environmental Protection
1324 Agency, Washington, DC, 1996.
- 1325 US EPA, Method 8270E (SW-846): Semivolatile Organic Compounds by Gas Chromatography/
1326 Mass Spectrometry (GC/MS). US Environmental Protection Agency, Washington, DC,
1327 2014.
- 1328 US EPA, Procedures for Chemical Risk Evaluation Under the Amended Toxic Substances Control
1329 Act. 82 FR 33726, Vol. 40 CFR 702. US Environmental Protection Agency, Washington,
1330 DC, 2017.
- 1331 Van Krevelen, D. W., 1950. Graphical-statistical method for the study of structure and reaction
1332 processes of coal. *Fuel.* 29, 269-284.
- 1333 Van Krevelen, D. W., 1984. Organic geochemistry - old and new. *Org Geochem.* 6, 1-10.
- 1334 Vanini, G., et al., 2020. Characterization of nonvolatile polar compounds from Brazilian oils by
1335 electrospray ionization with FT-ICR MS and Orbitrap-MS. *Fuel.* 282, 118790.
- 1336 Ventura, G. T., et al., 2011. Analysis of petroleum compositional similarity using multiway
1337 principal components analysis (MPCA) with comprehensive two-dimensional gas
1338 chromatographic data. *J Chromatogr A.* 1218, 2584-92.
- 1339 Vozka, P., et al., 2019. Jet fuel density via GC x GC-FID. *Fuel.* 235, 1052-1060.
- 1340 Walters, C. C., et al., 2015. Petroleum alteration by thermochemical sulfate reduction – A
1341 comprehensive molecular study of aromatic hydrocarbons and polar compounds. *Geochim*
1342 *Cosmochim Acta.* 153, 37-71.

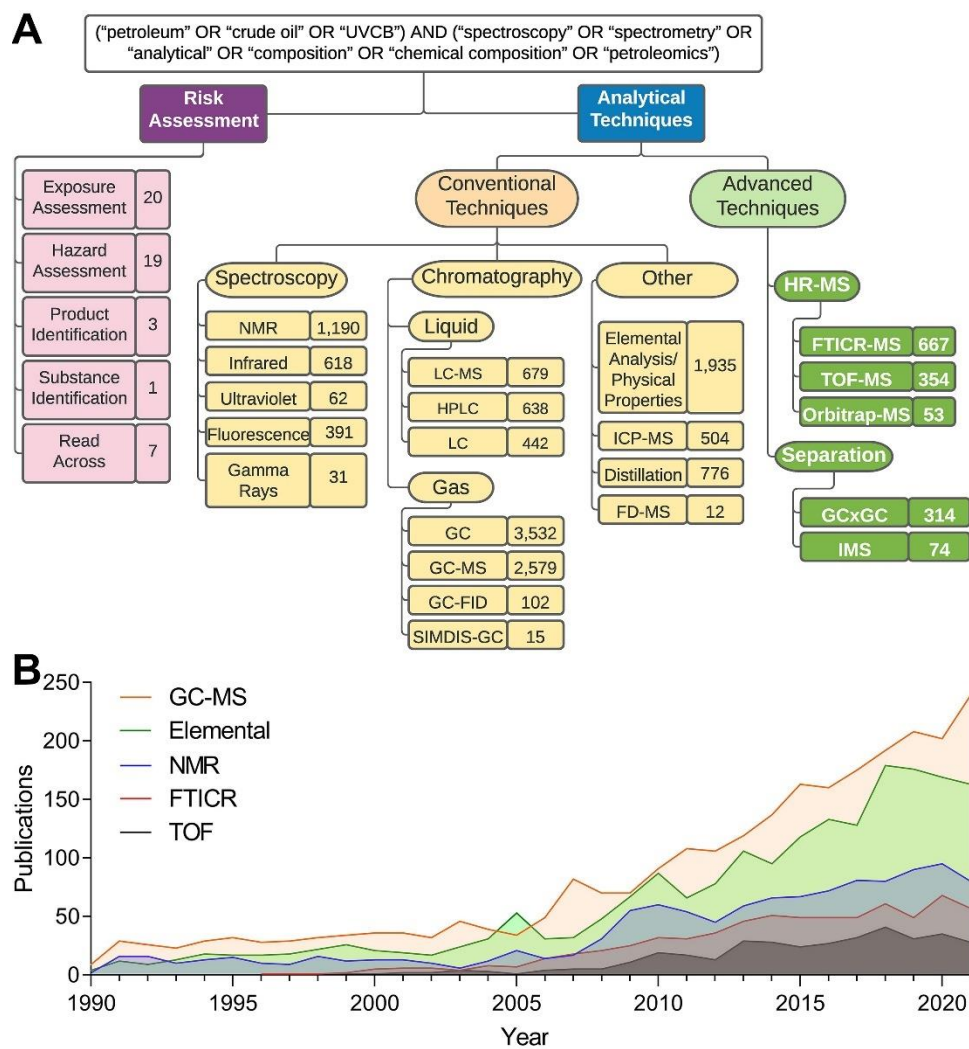
- 1343 Wang, X., et al., 2020. Characterization of wastewater effluent organic matter with different solid
1344 phase extraction sorbents. *Chemosphere*. 257, 127235.
- 1345 Wang, Z., et al., 1994. Fractionation of a Light Crude Oil and Identification and Quantitation of
1346 Aliphatic, Aromatic, and Biomarker Compounds by GC-FID and GC-MS, Part I. *J*
1347 *Chromatogr Sci*. 32, 361-366.
- 1348 Wang, Z., et al., 2011. Forensic fingerprinting and source identification of the 2009 Sarnia
1349 (Ontario) oil spill. *J Environ Monit*. 13, 3004-17.
- 1350 Wang, Z. D., Fingas, M., 2003. Fate and identification of spilled oils and petroleum products in
1351 the environment by GC-MS and GC-FID. *Energy Sources*. 25, 491-508.
- 1352 Weng, N., et al., 2015. Insight into unresolved complex mixtures of aromatic hydrocarbons in
1353 heavy oil via two-dimensional gas chromatography coupled with time-of-flight mass
1354 spectrometry analysis. *J Chromatogr A*. 1398, 94-107.
- 1355 Wise, S. A., et al., 2022. Advances in Chemical Analysis of Oil Spills Since the Deepwater
1356 Horizon Disaster. *Crit Rev Anal Chem*. 1-60.
- 1357 Wozniak, A. S., et al., 2019. Rapid Degradation of Oil in Mesocosm Simulations of Marine Oil
1358 Snow Events. *Environ Sci Technol*. 53, 3441-3450.
- 1359 Xia, Y., et al., 2021. Characterization of nitrogen-containing compounds in petroleum fractions by
1360 online reversed-phase liquid chromatography-electrospray ionization Orbitrap mass
1361 spectrometry. *Fuel*. 284, 119035.
- 1362 Xian, F., et al., 2012. High resolution mass spectrometry. *Anal Chem*. 84, 708-19.
- 1363 Yang, L., et al., 2019. Gas chromatography-Orbitrap mass spectrometry screening of organic
1364 chemicals in fly ash samples from industrial sources and implications for understanding
1365 the formation mechanisms of unintentional persistent organic pollutants. *Sci Total*
1366 *Environ*. 664, 107-115.
- 1367 Zhong, J. Y., et al., 2011. Combining advanced NMR techniques with ultrahigh resolution mass
1368 spectrometry: A new strategy for molecular scale characterization of macromolecular
1369 components of soil and sedimentary organic matter. *Org Geochem*. 42, 903-916.
- 1370 Zhu, G. Y., et al., 2020. Discovery and Molecular Characterization of Organic Caged Compounds
1371 and Polysulfanes in Zhongba81 Crude Oil, Sichuan Basin, China. *Energy Fuels*. 34, 6811-
1372 6821.
- 1373

1374 **Figure 1.** A timeline of major developments in the fields of analytical chemistry and data analysis
 1375 of petroleum UVCB, and the concomitant evolution of the regulatory frameworks for registration
 1376 and hazard classification of these substances See abbreviations in the text.
 1377



1378
 1379
 1380

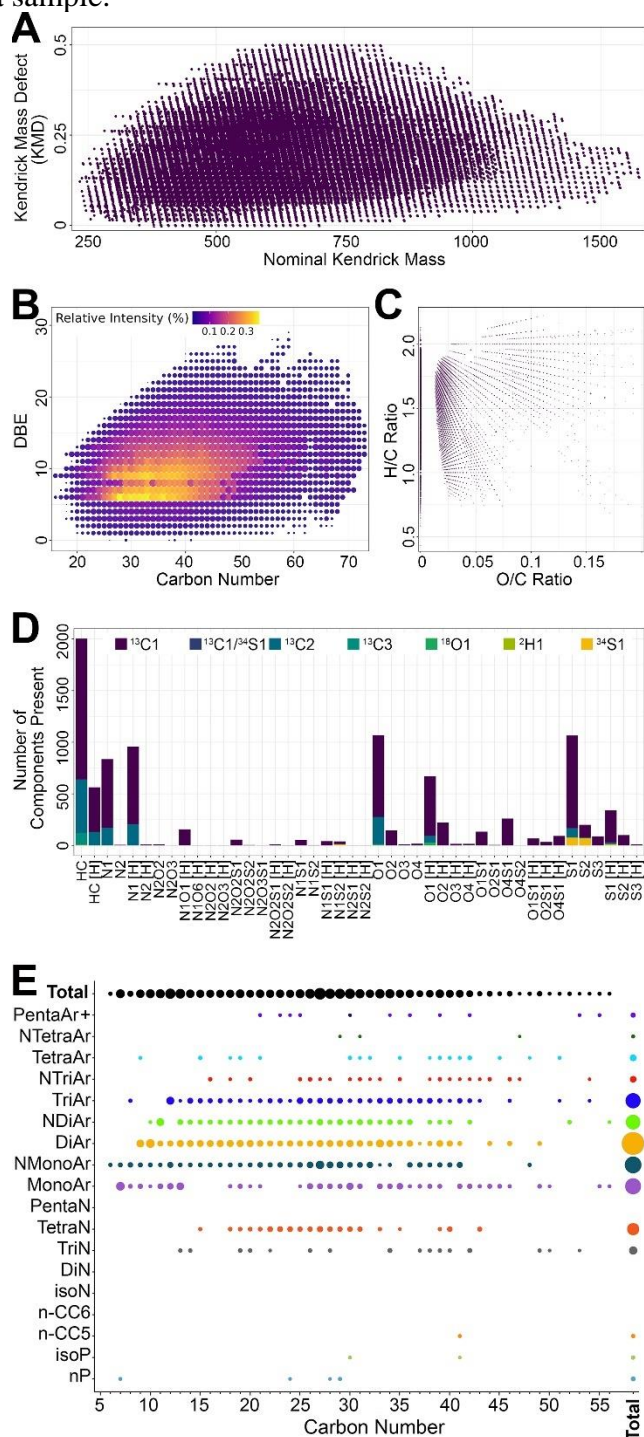
1381 **Figure 2.** Literature review of the major analytical methods and their application for the analysis of
 1382 petroleum substances. **(A)** A dendrogram of the major searches. The numbers indicate the quantity
 1383 of publications for each search. See Supplemental Table 1 for information on the exact search terms
 1384 and hyperlinks to the publications. **(B)** Cumulative histograms indicating the number of publications
 1385 across time through December 2021. Colors correspond to the methods indicated in the inset.
 1386



1387

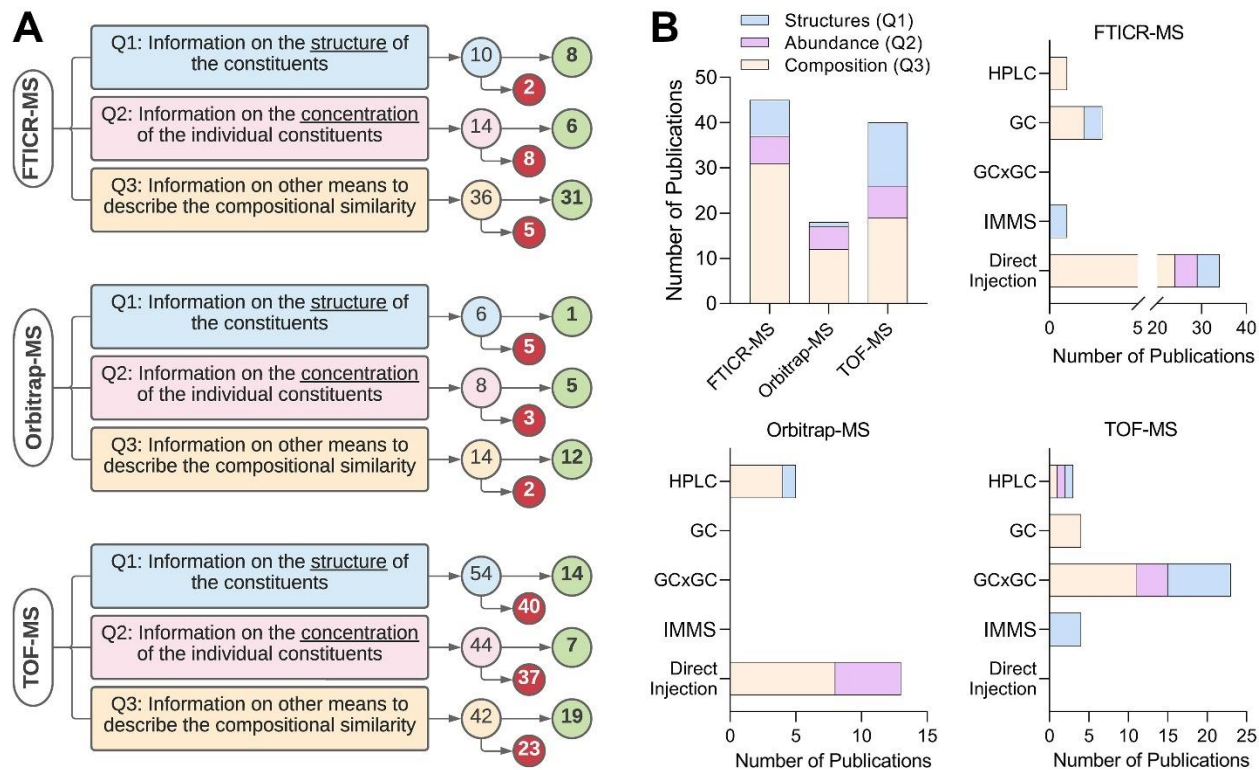
1388

1389 **Figure 3.** Example visualizations commonly used to represent hydrocarbon composition of
 1390 petroleum substances analyzed using high resolution MS techniques. **(A)** A Kendrick mass defect
 1391 (KMD) plot demonstrating repetitive patterns of CH₂-containing molecules in a petroleum sample.
 1392 **(B)** A DBE vs carbon number plot indicating the relative proportions of molecules varying by their
 1393 degree of aromaticity and carbon number in a sample. **(C)** van Krevelen plot display the degree of
 1394 oxidation by plotting the H/C versus O/C ratio in a sample. **(D)** A stacked bar plot showing relative
 1395 proportions of constituents from various chemical groups. **(E)** A plot of relative amounts of various
 1396 hydrocarbon blocks in a sample.



1397

1398 **Figure 4.** The use of various novel MS methods to address specific regulatory needs identified in
 1399 this review. **(A)** Scholarly publications that were identified as relevant to each regulatory
 1400 need/question (identified by colors). A total number of publications identified by a literature search
 1401 is listed in the first circle (see Supplemental Table 2 for details). Upon examination of each study's
 1402 content, a number of publications were deemed not relevant (red circles); the remaining studies are
 1403 shown in green circles. **(B)** The number of relevant publications as a function of the high-/ultra-high
 1404 MS technique. Top left, a stacked bar graph indicating the number of publications as they pertain to
 1405 each regulatory need/question in (A). Remaining stacked bar plots show the number of studies that
 1406 used various separation (HPLC, GC, GC×GC or IMMS) or direct injection with each MS technique.
 1407



1408