

AperTO - Archivio Istituzionale Open Access dell'Università di Torino

Features for non-targeted cross-sample analysis with comprehensive two-dimensional chromatography

This is the author's manuscript

Original Citation:

Availability:

This version is available <http://hdl.handle.net/2318/90589> since 2016-01-26T12:49:29Z

Published version:

DOI:10.1016/j.chroma.2011.07.046

Terms of use:

Open Access

Anyone can freely access the full text of works made available as "Open Access". Works made available under a Creative Commons license can be used according to the terms and conditions of said license. Use of all other works requires consent of the right holder (author or publisher) if not exempted from copyright protection by the applicable law.

(Article begins on next page)



UNIVERSITÀ DEGLI STUDI DI TORINO

This Accepted Author Manuscript (AAM) is copyrighted and published by Elsevier. It is posted here by agreement between Elsevier and the University of Turin. Changes resulting from the publishing process - such as editing, corrections, structural formatting, and other quality control mechanisms - may not be reflected in this version of the text. The definitive version of the text was subsequently published in [*Journal of Chromatography A*, Volume: 1226 Pages: 140-148 Published: FEB 24 2012, <http://dx.doi.org/10.1016/j.chroma.2011.07.046>].

You may download, copy and otherwise use the AAM for non-commercial purposes provided that your license is limited by the following restrictions:

- (1) You may use this AAM for non-commercial purposes only under the terms of the CC-BY-NC-ND license.
- (2) The integrity of the work and identification of the author, copyright owner, and publisher must be preserved in any copy.
- (3) You must attribute this AAM in the following format: Creative Commons BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/deed.en>), <http://dx.doi.org/10.1016/j.chroma.2011.07.046>

Features for non-targeted cross-sample analysis with comprehensive two-dimensional chromatography

Stephen E. Reichenbach^{a,*}, Xue Tian^a, Chiara Cordero^b, Qingping Tao^c

^a*University of Nebraska – Lincoln, Computer Science and Engineering Department,
Lincoln NE 68588-0115, USA*

^b*Dipartimento di Scienza e Tecnologia del Farmaco, Università degli Studi di Torino,
Via P. Giuria 9, I-10125 Torino, Italy*

^c*GC Image, LLC, PO Box 57403, Lincoln NE 68505-7403, USA*

Abstract

This review surveys different approaches for generating features from comprehensive two-dimensional chromatography for non-targeted cross-sample analysis. The goal of non-targeted cross-sample analysis is to discover relevant chemical characteristics (such as compositional similarities or differences) from multiple samples. In *non-targeted analysis*, the relevant characteristics are unknown, so individual features for all chemical constituents should be analyzed, not just those for targeted or selected analytes. *Cross-sample analysis* requires matching the corresponding features that characterize each constituent across multiple samples so that relevant characteristics or patterns can be recognized. Non-targeted, cross-sample analysis requires generating and matching all features across all samples. Applications of non-targeted cross-sample analysis include

*Corresponding author: +1.402.472.5007.

Email addresses: reich@cse.unl.edu (Stephen E. Reichenbach),
xtian@cse.unl.edu (Xue Tian), chiara.cordero@unito.it (Chiara Cordero),
qtao@gcimage.com (Qingping Tao)

URL: <http://www.gcimage.com> (Qingping Tao)

sample classification, chemical fingerprinting, monitoring, sample clustering, and chemical marker discovery. Comprehensive two-dimensional chromatography is a powerful technology for separating complex samples and so is well suited for non-targeted cross-sample analysis. However, two-dimensional chromatographic data is typically large and complex, so the computational tasks of extracting and matching features for pattern recognition are challenging. This review examines five general approaches that researchers have applied to these difficult problems: visual image comparisons, datapoint feature analysis, peak feature analysis, region feature analysis, and peak-region feature analysis.

Keywords: Comprehensive two-dimensional gas chromatography (GC×GC), Comprehensive two-dimensional liquid chromatography (LC×LC), Non-targeted analysis, Cross-sample analysis, Feature generation and matching, Pattern recognition

1 1. Introduction

2 The goal of non-targeted cross-sample analysis is to discover relevant
3 chemical characteristics (such as compositional similarities or differences)
4 from multiple samples. Some applications of non-targeted cross-sample
5 analysis are:

- 6 • **Classification.** Given a sample from an unknown class and
7 exemplary samples from a set of known classes, determine the class of
8 the unknown sample. For example, given samples of cancerous tumors
9 labeled by grade, determine the tumor grade for an ungraded
10 sample.[1]

- 11 • **Chemical fingerprinting.** Given a sample from an unknown source
12 and exemplary samples from multiple known sources, determine the
13 source of the unknown sample. For example, given a sample of
14 environmental pollution from an unknown source and labeled samples
15 from several possible sources of the pollution, identify the source for
16 the pollution.[2] Fingerprinting is a type of classification problem
17 except that each class is restricted to a single source, whereas the
18 general classification problem allows each class to have multiple
19 similar sources.

- 20 • **Monitoring.** Given a sequence of samples, identify samples that
21 have uncharacteristic differences with other samples, e.g., for quality
22 assurance. Monitoring also can be used to discover trends in sample
23 sequences, even recognizing subtle changes if they are progressive or
24 cyclical. For example, use a time-sequence of samples from an
25 environmental oil spill to track and understand the weathering
26 processes on oil constituents.[3]

- 27 • **Clustering.** Given a set of samples, partition subsets such that
28 samples within each subset are relatively similar and samples in
29 different subsets are relatively dissimilar. For example, given multiple
30 samples from oil reservoirs, use clustering to determine the number of
31 distinct reservoirs.[4]

- 32 • **Marker discovery.** Given a set of exemplary samples from known
33 classes, determine the chemical characteristics that are most relevant
34 for distinguishing the classes. For example, given samples of tumors

35 labeled by grade, determine which characteristics (i.e., biomarkers)
36 are most useful in distinguishing different tumor grades.[1]

37 Non-targeted cross-sample analysis should evaluate each and every
38 constituent in each and every sample. For *non-targeted analysis*, the
39 relevant chemical characteristics are not known, so the analysis should
40 generate characteristic *feature(s)* for each and every constituent. Typically,
41 detector intensities or mass spectral (total and/or selected ion) intensities
42 are used as characteristic features because they indicate the analyte
43 concentrations (or amounts) and provide information for chemical
44 identification. *Cross-sample analysis* should compare the same chemical
45 characteristics across multiple samples, so it is necessary to correctly match
46 the corresponding features that characterize the same analyte in different
47 samples. For example, peak matching would establish which peaks in
48 different samples result from the same analyte. Typically, other features,
49 such as retention times and/or mass spectral signatures, are used to match
50 the characteristic features.

51 Non-targeted cross-sample analysis requires comprehensive, selective,
52 matched, accurate features. If the features aren't comprehensive, then
53 relevant characteristics may not be analyzed. If the features aren't selective,
54 then relevant trace constituents may be obscured by more prevalent but
55 less relevant constituents. If the features aren't matched, then the analysis
56 is confounded by incorrect comparisons. If the features aren't accurate,
57 then the analysis may be unable to detect subtle differences.

58 Comprehensive two-dimensional gas chromatography (GC×GC) and
59 related techniques are well-suited for non-targeted cross-sample analysis

60 because they offer increased separation capacity, higher-dimensional
61 structure-retention relationships, and improved signal-to-noise ratio (SNR),
62 compared to traditional one-dimensional chromatography. Comprehensive
63 two-dimensional chromatography preserves separations at each stage and
64 submits the entire sample to analysis, providing for comprehensive features.
65 Increased separation capacity enables more selective features. The
66 higher-dimensional structure relationships can be exploited for better
67 matched features. And, the improved SNR increases the quantitative
68 accuracy of characteristic features.

69 Comprehensive two-dimensional chromatography offers unprecedented
70 information on compositional characteristics of complex samples, but the
71 size and complexity of the data makes data analysis to extract that
72 information a challenging problem. The most relevant features for a
73 particular cross-sample analysis may be related to trace constituents and/or
74 unidentified compounds. Relevant patterns may involve subtle relationships
75 among multiple features. So, the goal of non-targeted cross-sample analysis
76 is to extract and analyze all of the information that could be relevant. In
77 some sense, it is the ultimate information processing challenge.

78 The typical data processing sequence for non-targeted cross-sample
79 analysis is:

- 80 1. Preprocess individual chromatograms.
- 81 2. Generate features for each chromatogram.
- 82 3. Match features across chromatograms.
- 83 4. Recognize relevant patterns.

84 The purpose of this review is to examine various approaches that
85 researchers have applied to Steps 2 and 3 — feature generation and
86 matching — but Steps 1 and 4 merit a brief discussion. Preprocessing (Step
87 1) involves operations (e.g., baseline correction, [5]–[8], peak
88 detection[9]–[13], coeluted peak detection[14]–[25], and alignment
89 [24][26]–[36]) that prepare data for further analysis, but which are not
90 specific to non-target cross-sample analysis. Therefore, general
91 preprocessing methods can be used for these operations. In pattern
92 recognition (Step 4), the matched comparative features are analyzed to
93 recognize relevant characteristics or patterns among samples. Such pattern
94 recognition is not specific to chromatographic analysis and so can be
95 performed with various general-purpose methods, including statistical
96 methods such as principal component analysis (PCA), analysis of variance
97 (ANOVA), and discriminant function analysis (DFA), and machine-learning
98 methods such as support vector machines (SVM), neural networks, and
99 decision trees[1][4][31][35]–[58]. Of course, research continues to improve
100 methods for preprocessing and pattern recognition and to evaluate their
101 effectiveness for non-targeted cross-sample chromatographic analysis, but
102 that research is not the focus of this review.

103 This review describes five different types of features that have been
104 used for non-targeted cross-sample analyses with comprehensive
105 two-dimensional chromatography: visual images, datapoints, peaks, regions,
106 and peak-regions. Visual images present chromatograms using various
107 methods for two-dimensional data, including pseudo-colorization, contour
108 plots, and three-dimensional projections. Datapoint analyses treat each

109 datapoint as a feature, allowing chromatograms to be compared intensity
110 by intensity. Peak-based approaches attempt to separately integrate the
111 intensities from multiple datapoints that induced by each individual
112 analyte. Regional features aggregate datapoints in separate regions of the
113 two-dimensional chromatographic plane. Peak-region methods attempt to
114 define a region for each individual analyte.

115 Some examples of previous research illustrate each approach to
116 features for two-dimensional chromatographic analyses, with most research
117 involving GC×GC. The order of presentation roughly follows the historical
118 development. The discussion of each approach presents advantages and
119 problematic issues. Other authors have provided more general reviews of
120 GC×GC and related technologies and provide a broader context for this
121 review.[59]–[77]

122 **2. Visual Features**

123 The earliest non-targeted cross-sample analyses with comprehensive
124 two-dimensional chromatography were conducted without benefit of
125 software specifically designed for operating on two-dimensional
126 chromatographic data. Therefore, most early cross-sample comparisons
127 were primarily qualitative visual comparisons using general-purpose
128 software. In particular, two-dimensional chromatograms can be regarded as
129 digital images of the chromatographic plane. Digital images are
130 two-dimensional arrays of intensities and the datapoint intensities of
131 two-dimensional chromatographs are represented naturally in
132 two-dimensional arrays arranged so that the abscissa (X-axis, left-to-right)

133 is the elapsed time for the first-column separation and the ordinate (Y-axis,
134 bottom-to-top) is the elapsed time for the second-column separation. Then,
135 digital image visualization and processing methods can be used for
136 two-dimensional chromatograms.

137 In 1990, Bushey and Jorgenson[78] demonstrated comprehensive
138 two-dimensional liquid chromatography LC×LC and showed data from a
139 UV detector as surface plots with three-dimensional projection to two
140 dimensions. They presented side-by-side visualizations of reconstituted
141 serum from a human and from a horse, but did not make explicit
142 comparisons of the samples.

143 Blomberg et al.[79] showed side-by-side two-dimensional contour plots
144 of GC×GC data from a flame ionization detector (FID) for distillation
145 fractions of a heavy catalytic cracked cycle oil before and after
146 hydrogenation to illustrate the conversion of olefins and sulfur compounds.
147 Their results showed that “a clear distinction between different products is
148 visible immediately.”[79, p. 544] For perspective on the computers of the
149 time, they used a computer with 100MHz processor, 32 megabytes of
150 memory, and generic scientific data processing and visualization software.
151 The authors noted the need for more automated processing to characterize
152 and compare samples: “The vast amount of data generated, necessitate
153 that considerable effort has to be put in software and hardware
154 developments for automated interpretation.”[79, p. 544]

155 Gaines et al.[2] presented GC×GC-FID data from an oil spill sample
156 and from two potential sources for the spill as pseudo-colored images with
157 a cold-to-hot color scale for qualitative visual comparison. Their goal was

158 to demonstrate GC×GC for oil spill source identification, an application of
159 fingerprinting. The visual comparison allowed them to note that one of the
160 sources exhibited considerably fewer peaks in the heavy aromatic region
161 than the spill, which suggested that it was not the source for the spill.
162 They also made selected quantitative comparisons for fingerprinting, as
163 described here in subsequent sections.

164 Reddy et al.[80] used a side-by-side sequence of pseudo-colored
165 images to visualize GC×GC-FID data from progressively weathered
166 samples of a fuel oil standard for comparison to an image of data from a
167 sample of a decades-old fuel oil spill. Their goal was to understand
168 progressive changes in the oil. The visual comparisons allowed them to
169 observe that 70% evaporative weathering of the standard was required to
170 effect the same level of reduction of naphthalenic compounds observed in
171 the oil spill sample, but that level of weathering also removed other
172 components that still were present in the oil spill sample. They were able to
173 conclude that evaporative weathering could not solely account for the
174 GC×GC pattern observed in the oil-spill sample and that other factors,
175 such as water washing, preferential biodegradation, and microbial
176 degradation were required to explain the actual weathering of the oil spill.

177 Others have used visual comparisons for similar purposes. Janssen et
178 al.[81] visualized LC×GC-FID data for samples of edible oils and fats as
179 two-dimensional bubble plots with circles indicating detected peaks (with
180 dot locations determined by retention from LC and carbon number from
181 GC and dot areas determined by intensity). Perera et al.[82] showed a
182 region of GC×GC-FID data as contour plots to fingerprint headspace

183 volatiles from plant samples. Hope et al.[83] used contour plots to compare
184 total intensity counts (TICs) of data from GC×GC with time-of flight
185 (TOF) mass spectrometry (MS) for pre and post harvest lawn grass
186 extracts. Shellie et al.[39] used GC×GC-TOFMS to analyze mouse spleen
187 samples, then (a) visually compared averaged chromatograms from obese
188 mice to averaged chromatograms from control mice, (b) computed the
189 difference between the averaged chromatograms and showed images of the
190 positive and negative values, (c) compared bubble plots for averaged peaks,
191 and (d) compared bubble plots for relative weighted differences of averaged
192 peaks (dividing by the average standard deviation among sample groups).

193 Hollingsworth et al.[32] developed software methods for automatically
194 aligning chromatograms using reference peaks, normalizing intensities, and
195 visualizing the differences by various image-based methods, including
196 time-loop flicker (switching between images) and colorized differences.
197 Figure 1 illustrates a small chromatographic region with benzene, toluene,
198 ethylbenzene, and xylene (BTEX) peaks and a visualization of the
199 differences between two aligned chromatograms. Nelson et al.[84] and
200 Wardlaw et al.[85] used these methods to illustrate weathering of an oil spill
201 and oil seep. Cordero et al.[51] used these methods to compare
202 chromatograms from coffee samples. Such visualizations of pointwise
203 differences provide a segue to the next approach for non-targeted
204 multi-sample analyses — pointwise feature analysis.

205 Visual comparisons continue to be used both as a preliminary tool and
206 as an investigatory and confirmatory method for automated methods.
207 However, visual analyses are insufficient in several respects: the approach is

208 not quantitative, subtle differences and complex patterns may not be
209 visible, and the approach is not well suited for cross-sample analysis with
210 large sample sets.

211 **3. Datapoint Features**

212 Quantitative pointwise comparison is a natural progression from visual
213 image comparison. In a pointwise approach, chromatograms are compared
214 point-by-point (or in imaging terms pixel-by-pixel). With this approach,
215 each datapoint is a feature and the datapoint features at the same retention
216 times are implicitly matched.

217 In 2002, Johnson and Synovec[37] used quantitative datapoint features
218 (i.e., the chromatographic intensities at each datapoint) of GC×GC-FID
219 data to recognize patterns in different jet fuel mixtures. Their first
220 experiments involved five replicates for each of nine different mixtures of
221 two fuels for a total of 45 chromatograms each with 120K datapoints. Their
222 second experiments involved three replicates for each of thirteen different
223 classes for a total of 39 chromatograms each with 105K datapoints. The
224 potential relevance of each feature was computed by ANOVA, as the Fisher
225 f ratio — the variance between classes divided by the variance within
226 classes. Then, features were selected based on a f -ratio threshold that
227 yielded good class separation in the space defined by the first two
228 components of PCA. In this way, they reduced the number of features to a
229 few hundred, which gave good PCA separation of classes and good
230 organization in a K-means dendrogram.

231 Mohler et al.[40] and Pierce et al.[41] applied PCA to

232 GC×GC-TOFMS datapoint intensities at selected mass-to-charge (m/z)
233 channels to show class separations for yeast[40] and plant[41] samples.
234 Pierce et al.[42] analyzed organic acid metabolites in urine samples with
235 GC×GC-TOFMS by computing the f ratios at every mass-to-charge (m/z)
236 channel of each chromatographic datapoint and then summing the f ratios
237 along the m/z dimension (i.e., for each datapoint). Then, they selected
238 peaks with features (i.e., datapoints) having the largest weighted and
239 unweighted f -ratio sums. For peaks indicated by the f -ratio sums, the
240 ratios of the peak volumes between samples from non-pregnant women to
241 samples from pregnant women indicated that those components
242 significantly differentiated between the two classes.

243 Guo and Lidstrom[46] applied the same approach with
244 GC×GC-TOFMS data to investigate differences in metabolite profiles of
245 methylotrophic bacteria. Mohler et al.[43] used the same approach to
246 GC×GC-TOFMS data for yeast metabolites and then performed the
247 Student's t -test as a check on the volumes of the peaks indicated by the
248 summed f ratios. Subsequently, Mohler et al.[47] used the ratios of the
249 largest and smallest signals in GC×GC-TOFMS data to distinguish
250 datapoints and then peaks that changed in concert with the dissolved
251 oxygen cycle of yeast. Vial et al.[35, 58] used dynamic peak alignment
252 followed by PCA for GC×GC-MS data for several tobacco extracts and
253 later used correlation with class members to assess the discriminatory
254 power of each datapoint to analyze a large set of GC×GC-MS
255 chromatograms for tobacco extracts in three different classes. Gröger et
256 al.[45] used multidimensional scaling, hierarchical clustering, and PCA on

257 datapoint intensities to perform clustering and Fisher criterion to identify
258 discriminating datapoints for illicit drug samples. Gröger and
259 Zimmermann[36] used *t*-tests to select significant datapoint features from
260 selected channels of GC×GC-TOFMS data for partial least-squares (PLS)
261 discriminant analysis (DA). Ventura et al.[57] recently used multiway PCA
262 on GC×GC-FID data for maltene fractions of crude oils.

263 Hollingsworth et al.[32], Mohler et al.[40, 47], Almstetter et al.[34],
264 Gröger and Zimmermann[36], and others have noted the importance of data
265 alignment for datapoint feature analysis. Hollingsworth et al.[32],
266 Almstetter et al.[34], and others have developed alignment algorithms.
267 Gröger and Zimmermann[36] implemented alignment and other
268 preprocessing operations with parallel processing. The scope of this review
269 does not include alignment algorithms.

270 Chromatographic misalignment and peak shape variations pose serious
271 problems for pointwise cross-sample analysis. The features are individual
272 datapoints, so if there is any misalignment between any pairs of samples,
273 even as small as a fraction of a datapoint interval, then the features are
274 incorrectly matched. Misalignments, both global and local, naturally occur
275 even in well controlled conditions. Analytes normally elute over multiple
276 datapoints, so the effects of small misalignments are mitigated, but
277 misalignment is a fundamental issue that is difficult to eliminate. Like
278 differences due to alignment, peak-shape differences are erroneously seen as
279 quantitative differences in datapoint features. Another issue is that
280 pointwise analysis involves many features and many of those features are
281 highly redundant. Both the number of features and feature redundancy

282 complicate pattern recognition. In view of these issues, it can be argued
283 that datapoint features may be too selective, thereby generating numerous
284 features for slightly varying retention times within individual
285 chromatographic peaks.

286 **4. Peak Features**

287 Peak features aggregate multiple datapoints with the goal of
288 characterizing individual analytes (e.g., summing all datapoint intensities
289 that are attributed to each detected peak). Peak features characterize
290 larger, more meaningful chromatographic structures, resulting in fewer
291 features that are less redundant than datapoint features. Peak features also
292 are less sensitive to misalignment and peak-shape variations than datapoint
293 features because peaks typically span many datapoints. However, unlike
294 datapoint features, peak features are not implicitly matched. So, after
295 preprocessing and peak detection, the detected peaks in each
296 chromatogram that are induced by same analyte must be matched. Feature
297 matching is a critical challenge for peak-feature analysis.

298 Gaines et al.[2] provided an early demonstration of using quantitative
299 characterizations of individual peaks and groups of peaks (i.e., the
300 aggregation of several detected peaks) in GC×GC-FID data to fingerprint
301 samples of an oil spill and potential sources in order to identify the source
302 of the spill. Their analysis used summed intensities of four peaks and nine
303 peak groups that were selected because of their suitability for source
304 determination, so the analysis was not comprehensive, but was quite
305 advanced given the lack of software for two-dimensional chromatography at

306 the time. Also, the selections were performed by hand and so were not
307 automated. Bar charts with the intensities of the selected features showed
308 that one potential source was compositionally more similar to the spill than
309 the other was.

310 Mispelaar et al.[38, 4] used a much larger number of peaks to
311 distinguish samples from different oil reservoirs with GC×GC-FID. Their
312 peak detection found about 6000 peaks per chromatogram. They used
313 retention-time based alignment and filtering to match 3904 peaks, but the
314 results of their multi-variate analysis (MVA) were unsatisfactory. They
315 attributed the poor initial results to an inadequate number of samples with
316 many non-informative peaks and peak detection, quantification, and
317 matching errors. They then selected 292 peaks using an automated
318 criterion for the relative standard deviations (RSDs) between duplicate
319 samples to indicate peak detection and quantification errors. Most of the
320 automatically selected 292 peaks were in regions of the chromatogram with
321 lower peak density. Then, they manually selected 65 peaks for relevance
322 and absence of interference. This small fraction of the peaks (about 1% of
323 the detected peaks) was adequate for clustering the samples according to
324 reservoir, but the feature reduction is indicative of the difficulties of reliable
325 peak detection and matching. Such selective processing could exclude
326 highly informative peaks.

327 In their work with mouse spleen samples, Shellie et al.[39] matched
328 peaks in each chromatogram to reference data using tolerances on retention
329 times and mass spectral matching similarity. The TIC of each peak that
330 matched the same reference peak was placed on the same row in a matrix

331 with a column each chromatogram. They did not report how many peaks
332 were detected or how many of the detected peaks were matched. Student's
333 *t*-tests were used to indicate the eleven metabolites exhibiting the most
334 significant differences between obese and control mice.

335 Qiu et al.[44] performed GC×GC-FID on volatile oils from Qianghuo,
336 a traditional Chinese medicine, from five regions. They did not report
337 parameters for rejecting peaks with low SNR nor the number of peaks
338 detected. They developed and implemented peak alignment and matching
339 methods (using retention times relative to reference peaks) to create a
340 matrix with 1544 peaks in fifteen samples. PCA analysis produced three
341 clusters, with separate clusters for samples from two of the five regions.
342 They used variable importance in the projection (VIP)[86, p. 397] to
343 identify potential marker compounds, finding some statistically significant
344 features, then used GC×GC-TOFMS for chemical identification of those
345 compounds.

346 Wardlaw et al.[85] developed an algorithm to track peaks between
347 similar samples based on retention times. The algorithm tracked about
348 1400 of about 4500 peaks in GC×GC chromatograms from oil samples from
349 the reservoir, sea floor, and sea surface.

350 Analyzing human serum with GC×GC-TOFMS, Oh et al.[87]
351 developed a peak sorting method to recognize peaks from the same
352 metabolite in different chromatograms. Their algorithm used several search
353 criteria with retention times and mass spectra, with options to eliminate
354 non-target peaks. Peaks with low signal-to-noise ratio were discarded
355 during peak detection. The matched peaks showed high correlation for

356 retention times and mass spectra, but only 105 peaks were matched across
357 all fifteen chromatograms, even with five replicates for each of three
358 samples.

359 Gaquerel et al.[48] used GC×GC-TOFMS to analyze the effect of oral
360 secretions on volatile plant emissions. Peak detection yielded about 600
361 peaks in each of the 108 samples (subject to a threshold SNR of 10). The
362 authors noted that inconsistencies in the numbers of the detected peaks in
363 each chromatogram complicated matching. In each of three sample periods,
364 the peak set of the chromatogram with the largest number of detected
365 peaks was used as reference data for matching (with the matching
366 procedure developed by Shellie et al.[39]), reducing the number of matched
367 peaks to about 400, which then were corrected for false positives from the
368 alignment and matching procedure. ANOVA followed by another manual
369 check for false positives from the peak alignment and matching was used to
370 select about 15% the peaks for MVA with hierarchical clustering analysis
371 (HCA) and PCA.

372 Li et al.[49] analyzed blood plasma with GC×GC-TOFMS. They used
373 a mass spectral filter to extract peaks for trimethylsilylated metabolites,
374 then applied a peak alignment method and a peak matching algorithm to
375 create a matrix with 492 metabolites in 79 chromatograms. They tried
376 several modeling methods, including PLS-DA, in which some problems that
377 were attributed to missing values from peak matching were resolved by
378 additional peak filtering. Then, VIP was used to indicate potential
379 biomarkers.

380 Reichenbach et al.[88] developed Smart Templates™ for peak matching.

381 The template records a prototypical pattern of peaks with retention times
382 and associated metadata, such as chemical identities and compound-group
383 membership. Then, the template pattern is matched to the detected peaks
384 in subsequent chromatograms and the metadata are copied from the
385 template to identify the matched peaks. The matching process explores the
386 space of affine geometric transforms to maximize the number of matched
387 peaks and minimize the residual geometric error. Smart Templates employ
388 rule-based constraints (e.g., multispectral matching) to increase matching
389 accuracy. Smart templates also carry other structures, such as text and
390 chemical-structure annotations and polygonal regions (which can be used
391 for region features, described below). They demonstrated the approach and
392 associated methods on urine samples analyzed by LC×LC with a
393 ultraviolet (UV) diode array detector (DAD). Figure 2 illustrates template
394 peak matching with a template derived from the detected peaks of one
395 chromatogram matched to the detected peaks of another chromatogram.

396 Cordero et al.[89] analyzed volatile fractions of roasted hazelnuts with
397 GC×GC-MS, then performed peak matching with templates in two different
398 ways. In the first approach, they aligned and summed the chromatograms,
399 then created a feature template comprised by the 411 peaks detected in the
400 cumulative chromatogram. That template then was matched to each
401 individual chromatogram, with matching rates ranging from 68% to 79%.
402 In the second approach, they performed a sequential template matching
403 that used both retention-time patterns and mass spectral matching criteria.
404 At each step of the sequence, unmatched peaks were added to build a
405 comprehensive template. At the end of the sequence, the comprehensive

406 template was matched to each chromatogram and any peak matching with
407 at least two chromatograms were retained in a consensus template. The
408 consensus template contained 422 peaks and the matching rates ranged
409 from 52% to 78%, with 196 peaks matching for all nine chromatograms.
410 For both peak matching methods, the feature fingerprints of samples from
411 nine regions were sifted for the largest normalized intensities and many of
412 the indicated compounds have a known role in defining sensory properties.

413 Castillo et al.[55] used GC×GC-TOFMS to analyze a variety of
414 samples for metabolomic characteristics. They developed a processing
415 sequence of peak detection, matching, filtering, normalization, and
416 identification. The matching algorithm used a scoring metric to choose
417 some matches over others. For a set of 60 serum samples, almost 15,000
418 prospective compounds were filtered to 1540 on the basis of matching a
419 sufficient number of chromatograms, then to 1013 compounds by mass
420 spectral and chromatographic constraints. The resulting feature vectors
421 were analyzed by PCA, which separated samples by their storage
422 temperature.

423 Koek et al.[56] evaluated the analyst and computer time required to
424 process GC×GC-TOFMS datasets for mouse liver samples to produce a
425 table of 170 metabolites in 29 samples. The analysis required
426 approximately 50h of analyst time and 60h of computer time, with
427 substantial analyst time required for optimization and construction of the
428 reference target table and dealing with problems of missing peak values.
429 These times are indicative that reliable peak matching, even with recent
430 software for GC×GC, is not yet automated. Subsequently, they evaluated

431 the resulting metabolite profiles with PCA and PCA-DA.

432 Peak detection errors as well as the inherent ambiguity of matching
433 both contribute to make comprehensive peak matching (i.e., matching all
434 peaks) across many samples intractable. Trace peaks may be detected in
435 some samples, but not in others. Coeluting analytes may be detected as
436 separate peaks in some chromatograms but as one peak in other
437 chromatograms. The peaks of different analytes may be incorrectly
438 matched, especially if constituents differ from sample to sample. To
439 overcome these challenges, researchers filter the peaks that are used for
440 cross-sample analysis. However, such filtering is unreliable and difficult to
441 automate. And, to the extent that peaks are correctly filtered, the analysis
442 is no longer truly comprehensive. Despite extensive research, methods for
443 automated peak matching still are error-prone and/or not comprehensive.
444 Despite these problems, peak features can be effectively used in many
445 applications for non-targeted cross-sample analysis.

446 **5. Region Features**

447 Region features characterize multiple datapoints (e.g., summing the
448 intensities at all datapoints in each region). Like peak features, region
449 features can characterize larger, more meaningful chromatographic
450 structures than datapoint features, resulting in fewer features that are less
451 redundant. Like peak features, region features are less sensitive to
452 misalignment than datapoint analysis.

453 For non-targeted analysis, the feature regions should be defined to
454 cover the entire chromatographic space in which analytes are present. When

455 used for cross-sample analysis, the same regions in different chromatograms
456 are implicitly matched, thereby avoiding the matching problem that is
457 inherent with peak features. However, either the chromatograms should be
458 aligned or the regions should be adjusted geometrically so that the same
459 regions in different chromatograms encompass the same analyte(s). As
460 geometric shapes, regions are amenable to geometric transformations to fit
461 different chromatograms in cases of variable retention times.

462 Two concerns with region features are that a region may encompass
463 more than one analyte and that one analyte may be spread across more
464 than one region. In the first case, selectivity is reduced as compared with
465 peak features (although peak features also may not separate coeluted
466 peaks). In the second case, multiple features for a single analyte are more
467 susceptible to errors related to misalignment as compared with peak
468 features (although peak features also may incorrectly split analyte peaks).

469 Mispelaar[4, 38] created a hand-drawn mesh of contiguous polygons to
470 subjectively encompass different groups of interest in diesel samples and
471 demonstrated the utility of geometric transformations to better fit different
472 chromatograms. Figure 3 illustrates a similar mesh for GC \times GC-FID[90]
473 with automatically drawn vertical lines at linear retention indices based on
474 the *n*-alkanes and hand-drawn lines to separate compound groups. As
475 Mispelaar noted, some prior knowledge of the sample is required to define
476 regions related to its components and component groups. And, as can be
477 seen, there are regions with multiple analytes and analyte peaks spread
478 across multiple regions.

479 To quantify weathering of an oil spill by GC \times GC-FID, Arey et al.[3]

480 created a grid with region boundaries defined by computed contours of
481 hydrocarbon vapor pressure and aqueous solubility. With this approach, no
482 prior knowledge of the nature of the sample is required, but regions may
483 contain multiple analytes and analyte peaks may straddle multiple regions.
484 To mitigate the effect of misalignment, they used trapezoidal weighting
485 functions at the borders between regions. With contour lines that are
486 roughly orthogonal, the grid can be remapped naturally to a rectangular
487 array and colorized according to intensity for convenient visualization.
488 They applied the analysis to investigate different weathering processes on
489 oil spills, including evaporation, dissolution, biodegradation,
490 photodegradation, and other processes. Wardlaw et al.[85] used these same
491 lines to warp chromatographic images.

492 To analyze Chinese medicine volatile oils with GC×GC-TOFMS, Qiu
493 et al.[44] used integration in four regions (mostly, but not fully covering the
494 analytes) to compute averages and show differences among five geographical
495 classes. Mullins et al.[91] used seven large regions to characterize compound
496 groups in downhole fluid analysis with GC×GC-FID and GC×GC-TOFMS.
497 They plotted ratios of the summed peak intensities within each region in a
498 spider diagram to visualize similarities and differences. Betancourt et al.[92]
499 used spider diagrams to visualize features for nine large compound-based
500 regions and subdivisions of those regions split by retention indices. Ventura
501 et al.[93] extended the approach to twelve regions. Vaz-Freire et al.[50]
502 divided chromatograms from olive oil samples into twelve rectangular
503 regions, then performed ANOVA and PCA with the regional features.

504 The principal issue with region features is that selectivity is reduced to

505 the extent that peaks of multiple analytes are included in the same region.
506 For some applications, such as petroleum analysis, the goal may be
507 comprehensive group-type analysis, so loss of selectivity within groups is
508 not problematic. However, the loss of selectivity could be a problem in
509 many applications, especially if a critical trace analyte is in the same region
510 as a predominant analyte that is irrelevant to the application.

511 **6. Peak-Region Features**

512 The final type of feature surveyed in this review is the peak-region.
513 Peak-region features attempt to define one region per peak. This approach
514 seeks to achieve the one-feature-to-one-analyte selectivity of peak features
515 but with the implicit matching of region features.

516 Schmarr et al.[53, 54] and Reichenbach and co-workers[51, 52, 1]
517 described similar approaches to defining regions for individual peaks
518 detected across multiple samples. Schmarr and Bernhardt indicated that
519 this general approach is common for 2D gel electrophoresis. After
520 preprocessing, including alignment, the chromatograms are combined (e.g.,
521 simply by addition or other fusion operations [94]) to form a single
522 chromatogram that is reflective of all of the constituents in all samples.
523 Then, the boundaries that delineate each peak are recorded as a region in a
524 template. That template is then geometrically mapped back to each
525 chromatogram and each region defines a feature for each chromatogram.
526 The features are comprehensive, accounting for every analyte, and feature
527 matching is implicitly performed by the retention-time mapping.

528 Schmarr and Bernhardt [53] analyzed 32 samples of volatiles of

529 different fruits by GC×GC-MS. They performed baseline correction with
530 the rolling-ball method, then manually generated warp graphs to determine
531 warping transforms to align 31 chromatograms to a reference
532 chromatogram. Then, each of the chromatograms was aligned by the
533 warping transform and combined using a weighted-mean “union fusion” [94].
534 They manually detected more than 700 spots indicative of peaks in the
535 fused chromatogram. Then, the spot patterns were mapped back to each
536 chromatogram according to the inverse of its warping transform and the
537 intensities for each region in each chromatogram were computed. The
538 software package that they used was optimized for gel electrophoresis rather
539 than GC×GC, so much of the processing was manual, requiring about 5h of
540 an analyst’s time for the 32 samples. They used HCA and PCA with the
541 resulting peak-region features to cluster samples. The different fruits
542 (apples, pears, and quince) formed clear clusters. The two pear varieties
543 and some of the six apple varieties formed sub-clusters. The mass-spectral
544 signatures were used for compound identification of spots which were
545 statistically relevant for differentiation. Using a similar approach for
546 analyzing red wines subjected to microoxygenation (MOX), Schmarr et
547 al.[54] were able to differentiate MOX treatments and specific varietal and
548 technological effects. They were able to identify areas in the 2D
549 chromatograms that were most responsible for discrimination among
550 different MOX treatments and the loadings of individual aroma compounds
551 suggested a set of markers for the MOX-induced modifications of volatiles.

552 Cordero et al.[51] analyzed samples of coffees and junipers by
553 GC×GC-MS. After preprocessing including peak detection, they identified

554 peaks that could be matched reliably across all chromatograms. These
555 reliable peaks were the basis of a registration template with mass spectral
556 matching rules that then was used to determine a geometric transform to
557 align the chromatograms. After alignment, the chromatograms were
558 summed to create a cumulative chromatogram. In three chromatograms of
559 coffee samples, about 1700 peaks were detected, about half of which were
560 reliable. They manually drew a mesh of about 1100 regions which were
561 combined with the registration peaks to create a feature template that
562 could be matched to individual chromatograms thereby transforming the
563 regions to maintain their positions relative to the reliable peaks. They
564 sifted the features by intensity, standard deviation, and relative standard
565 deviation to select relevant features but did not perform MVA because of
566 the small number of samples. Many of the indicated compounds were
567 known botanical, technological, and/or aromatic markers for coffee. For the
568 analysis of five chromatograms of juniper samples, there were about 100
569 reliable peaks and 727 peak-regions were drawn. Reichenbach et al.[52]
570 used the same approach for 39 urine samples analyzed by LC×LC. Then,
571 they performed classification with SVM and k-NN, evaluating the
572 performance using cross-validation.

573 Reichenbach et al.[1] analyzed data from GC×GC with high-resolution
574 mass spectrometry (HRMS) of samples from breast cancer tumors. There
575 were eighteen samples each from different individuals, with six samples each
576 for grades one to three as determined by a cancer pathologist. They
577 followed the same approach as Cordero et al.[51] except that the process,
578 including drawing the regions around the peaks detected in the cumulative

579 chromatogram, was performed automatically by newer software. About
580 3300 peaks were detected in each of the eighteen individual chromatograms,
581 but only thirteen were reliable across all eighteen chromatograms. Note
582 that reliability was defined as bidirectional matching between all possible
583 pairs (more than 300 matches for each common peak). In the cumulative
584 chromatogram, more than 3300 peak-regions were defined. Figure 4 shows
585 the cumulative chromatogram overlaid with black ovals for the reliable
586 peaks used for registration and red outlines for the peak-regions. They
587 applied several machine learning methods with the peak-region features to
588 classify samples by tumor grade and to indicate potential biomarkers for
589 tumor grade which then were investigated using the high-resolution mass
590 spectra.

591 The peak-region approach is more comprehensive than using reliably
592 matched peak features and is more selective than region features. As with
593 the other feature methods, misalignment is a potential source of errors. As
594 with peak features, peak detection errors, such as unseparated coelutions
595 and incorrectly split peaks, are another source of errors for peak-region
596 features.

597 **7. Conclusion**

598 A common goal of chemical analysis is to compare samples, either for a
599 few specific compounds (targeted analysis), for groups of compounds
600 (group-type analysis), or for all compounds (i.e., non-targeted analysis).
601 The key to comparative analyses is to establish correspondences between
602 features of different data sets, e.g., recognizing that a peak in the data for

603 one sample and a peak in the data for another sample are induced by the
604 same compound. Establishing correspondences — *feature matching* — is
605 necessary before it is possible to perform comparisons and pattern
606 recognition across sample sets.

607 Targeted analyses and group-type analyses are more straightforward
608 than non-target analyses. In targeted analyses, the compounds of interest
609 are known, so the chromatography can be tailored to provide selectivity for
610 those compounds and the data processing methods can be refined for
611 detecting and recognizing the features for those compounds. For group-type
612 analysis, the method need not be selective of every individual analyte, so
613 many problems of feature generation (e.g., peak unmixing) and matching
614 can be avoided. Comprehensive non-target analyses are more difficult
615 because the most relevant compounds are unknown, so the chromatography
616 and data processing cannot be tuned specifically for individual compounds
617 or for groups of compounds.

618 Non-targeted cross-sample analysis is especially difficult because it
619 requires the analysis of all analytes in all chromatograms of a sample set.
620 Applications of non-targeted cross-sample analysis include sample
621 classification, chemical fingerprinting, monitoring, sample clustering, and
622 chemical marker discovery. Comprehensive two-dimensional
623 chromatography is a powerful technology for separating complex mixtures
624 and so is well suited for comprehensive non-targeted analysis, but fully
625 extracting chemical information from large and complex datasets is
626 challenging and the subject of ongoing research. And, the difficulty of
627 comparative analyses increases with the size of the sample set.

628 Feature matching for comprehensive two-dimensional chromatography
629 can be based on retention times, spectral signature, detected intensity,
630 and/or other characteristics of features. Past research on non-targeted
631 cross-sample analysis with comprehensive two-dimensional chromatography
632 has demonstrated the usefulness of qualitative visualization, individual
633 datapoints, detected peaks, chromatographic regions, and comprehensive
634 peak-regions.

635 Each type of feature has advantages and disadvantages. Visualization
636 is simple and intuitive, but is not quantitative, important differences may
637 not be visible, and working with large sample sets is difficult. Datapoint
638 features are highly selective and implicitly matched across aligned
639 chromatograms, but they are subject to misalignment errors and generate a
640 large number of features, many of which are redundant. Peak features
641 characterize individual analytes and so are especially consistent with
642 analytical goals, but peak matching is an intractable problem. Region
643 features are more attuned to meaningful analytical characteristics than
644 datapoint features and are easier to match across samples than peak
645 features, but they may not be as selective as datapoint or peak features.
646 Peak-regions define a region for each peak across chromatograms and so
647 aim for selectivity and accurate feature matching, but still are subject to
648 errors from misalignment and peak detection failures.

649 Future research will refine, compare, and combine these approaches.
650 There has been little research to deeply examine the variables that affect
651 feature generation and matching in the different approaches and to validate
652 performance in cross-sample analyses. Advances in instrument technologies

653 could contribute to improved feature generation and matching, e.g., with
654 increased repeatability and reproducibility, greater mass spectrometric
655 accuracy, and more effective column sets. Feature generation and matching
656 might be improved by better preprocessing methods, especially for detection
657 of coeluted peaks, but also for baseline correction and alignment. Likewise,
658 more research is needed to compare the performance of different approaches
659 for feature generation and matching in different applications. Ultimately, a
660 hybrid approach, using a combination of different approaches, may be most
661 effective e.g., peak features for peaks that can be reliably matched, and
662 peak-region, region, or datapoint features for other chromatographic data.
663 Again, such combined approaches require a better understanding of the
664 variables that affect the performance of the different approaches.

665 **Acknowledgements**

666 This work was supported in part by the U.S. National Science
667 Foundation under Award Number IIP-1013180 and by the Nebraska Center
668 for Energy Sciences Research.

669 **References**

- 670 [1] S. E. Reichenbach, X. Tian, Q. Tao, E. B. Ledford, Jr., Z. Wu,
671 O. Fiehn, *Talanta* 83 (2011) 1279.
- 672 [2] R. B. Gaines, G. S. Frysinger, M. S. Hendrick-Smith, J. D. Stuart,
673 *Environ. Sci. Technol.* 33 (1999) 2106.
- 674 [3] J. S. Arey, R. K. Nelson, C. M. Reddy, *Environ. Sci. Technol.* 41
675 (2007) 5738.

- 676 [4] V. G. van Mispelaar, Chromametrics, Ph.D. thesis, University of
677 Amsterdam, 2005.
- 678 [5] S. R. Sternberg, Computer 16 (1983) 22.
- 679 [6] S. E. Reichenbach, M. Ni, D. Zhang, E. B. Ledford, Jr., Journal of
680 Chromatography A 985 (2003) 47.
- 681 [7] Y. Zhang, H.-L. Wu, A.-L. Xia, L.-H. Hu, H.-F. Zou, R.-Q. Yu, J.
682 Chromatogr. A 1167 (2007) 178.
- 683 [8] S. E. Reichenbach, P. W. Carr, D. R. Stoll, Q. Tao, Journal of
684 Chromatography A 1216 (2009) 3458.
- 685 [9] J. Beens, H. Boelens, R. Tijssen, J. Blomberg, J. High Resolut.
686 Chromatogr. 21 (1998) 47.
- 687 [10] Q. Song, A. Savant, S. E. Reichenbach, E. B. Ledford, Jr., in: Visual
688 Information Processing, Proc. SPIE 3808, pp. 2.
- 689 [11] S. E. Reichenbach, M. Ni, V. Kottapalli, A. Visvanathan,
690 Chemometrics and Intelligent Laboratory Systems 71 (2004) 107.
- 691 [12] S. Peters, G. Vivó-Truyols, P. Marriott, P. Schoenmakers, J.
692 Chromatogr. A 1156 (2007) 14.
- 693 [13] E. J. C. van der Klift, G. Vivó-Truyols, F. W. Claassen, F. L. van
694 Holthoon, T. A. van Beek, J. Chromatogr. A 1178 (2008) 43.
- 695 [14] C. A. Bruckner, B. J. Prazen, R. E. Synovec, Anal. Chem. 70 (1998)
696 2796.

- 697 [15] B. J. Prazen, C. A. Bruckner, R. E. Synovec, B. R. Kowalski, J.
698 Microcolumn Sep. 11 (1999) 97.
- 699 [16] C. G. Fraga, B. J. Prazen, R. E. Synovec, J. High Resolut.
700 Chromatogr. 23 (2000) 215.
- 701 [17] C. G. Fraga, C. A. Bruckner, R. E. Synovec, Anal. Chem. 73 (2001)
702 675.
- 703 [18] B. J. Prazen, K. J. Johnson, A. Weber, R. E. Synovec, Anal. Chem. 73
704 (2001) 5677.
- 705 [19] A. E. Sinha, C. G. Fraga, B. J. Prazen, R. E. Synovec, J. Chromatogr.
706 A 1027 (2004) 269.
- 707 [20] A. E. Sinha, J. L. Hope, B. J. Prazen, C. G. Fraga, E. J. Nilsson, R. E.
708 Synovec, J. Chromatogr. A 1056 (2004) 145.
- 709 [21] C. G. Fraga, C. A. Corley, J. Chromatogr. A 1096 (2005) 40.
- 710 [22] H. Kong, F. Ye, X. Lu, L. Guo, J. Tian, G. Xu, J. Chromatogr. A 1086
711 (2005) 160.
- 712 [23] J. C. Hoggard, R. E. Synovec, Anal. Chem. 79 (2007) 1611.
- 713 [24] T. Skov, J. C. Hoggard, R. Bro, R. E. Synovec, J. Chromatogr. A 1216
714 (2009) 4020.
- 715 [25] Z.-D. Zeng, S.-T. Chin, H. M. Hugel, P. J. Marriott, J. Chromatogr. A
716 1218 (2011) 2301.
- 717 [26] C. G. Fraga, B. J. Prazen, R. E. Synovec, Anal. Chem. 72 (2000) 4154.

- 718 [27] K. J. Johnson, B. W. Wright, K. H. Jarman, R. E. Synovec, J.
719 Chromatogr. A 996 (2003) 141.
- 720 [28] V. G. van Mispelaar, A. C. Tas, A. K. Smilde, P. J. Schoenmakers,
721 A. C. van Asten, J. Chromatogr. A 1019 (2003) 15.
- 722 [29] K. J. Johnson, B. J. Prazen, D. C. Young, R. E. Synovec, J. Sep. Sci.
723 27 (2004) 410.
- 724 [30] K. M. Pierce, L. F. Wood, B. W. Wright, R. E. Synovec, Anal. Chem.
725 77 (2005) 7735.
- 726 [31] K. M. Pierce, J. L. Hope, K. J. Johnson, B. W. Wright, R. E. Synovec,
727 J. Chromatogr. A 1096 (2005) 101.
- 728 [32] B. V. Hollingsworth, S. E. Reichenbach, Q. Tao, A. Visvanathan, J.
729 Chromatogr. A 1105 (2006) 51.
- 730 [33] D. Zhang, X. Huang, F. E. Regnier, M. Zhang, Anal. Chem. 80 (2008)
731 2664.
- 732 [34] M. F. Almstetter, I. J. Appel, M. A. Gruber, C. Lottaz, B. Timischl,
733 R. Spang, K. Dettmer, P. J. Oefner, Anal. Chem. 81 (2009) 5731.
- 734 [35] J. Vial, H. Noçairi, P. Sassiati, S. Mallipatu, G. Cognon, D. Thiébaud,
735 B. Teillet, D. N. Rutledge, J. Chromatogr. A 1216 (2009) 2866.
- 736 [36] T. Gröger, R. Zimmermann, Talanta 83 (2011) 1289.
- 737 [37] K. J. Johnson, R. E. Synovec, Chemom. Intell. Lab. Syst. 60 (2002)
738 225.

- 739 [38] V. G. van Mispelaar, H.-G. Janssen, A. C. Tas, P. J. Schoenmakers, J.
740 Chromatogr. A 1071 (2005) 229.
- 741 [39] R. A. Shellie, W. Welthagen, J. Zrostliková, J. Spranger, M. Ristow,
742 O. Fiehn, R. Zimmermann, J. Chromatogr. A 1086 (2005) 83.
- 743 [40] R. E. Mohler, K. M. Dombek, J. C. Hoggard, E. T. Young, R. E.
744 Synovec, Anal. Chem. 78 (2006) 2700.
- 745 [41] K. M. Pierce, J. L. Hope, J. C. Hoggard, R. E. Synovec, Talanta 70
746 (2006) 797.
- 747 [42] K. M. Pierce, J. C. Hoggard, J. L. Hope, P. M. Rainey, A. N.
748 Hoofnagle, R. M. Jack, B. W. Wright, R. E. Synovec, Anal. Chem. 78
749 (2006) 5068.
- 750 [43] R. E. Mohler, K. M. Dombek, J. C. Hoggard, K. M. Pierce, E. T.
751 Young, R. E. Synovec, Analyst 132 (2007) 756.
- 752 [44] Y. Qiu, X. Lu, T. Pang, S. Zhu, H. Kong, G. Xu, JPharm. Biomed.
753 Anal. 43 (2007) 1721.
- 754 [45] T. Grögera, M. Schäffer, M. Pütz, B. Ahrens, K. Drew, M. Eschner,
755 R. Zimmermann, J. Chromatogr. A 1200 (2008) 8.
- 756 [46] X. Guo, M. E. Lidstrom, Biotechnol. Bioeng. 99 (2008) 929.
- 757 [47] R. E. Mohler, B. P. Tu, K. M. Dombek, J. C. Hoggard, E. T. Young,
758 R. E. Synovec, J. Chromatogr. A 1186 (2008) 401.

- 759 [48] E. Gaquerel, A. Weinhold, I. T. Baldwin, *Plant Physiol.* 149 (2009)
760 1408.
- 761 [49] X. Li, Z. Xu, X. Lu, X. Yang, P. Yin, H. Kong, Y. Yu, G. Xu, *Anal.*
762 *Chim. Acta* 633 (2009) 257.
- 763 [50] L. T. Vaz-Freire, M. D. R. Gomes da Silva, A. M. C. Freitas, *Anal.*
764 *Chim. Acta* 633 (2009) 263.
- 765 [51] C. Cordero, E. Liberto, C. Bicchi, P. Rubiolo, S. E. Reichenbach,
766 X. Tian, Q. Tao, *J. Chromatogr. Sci.* 48 (2010) 251.
- 767 [52] S. E. Reichenbach, X. Tian, Q. Tao, D. R. Stoll, P. W. Carr, *Talanta*
768 83 (2010) 1365.
- 769 [53] H.-G. Schmarr, J. Bernhardt, *J. Chromatogr. A* 1217 (2010) 565.
- 770 [54] H.-G. Schmarr, J. Bernhardt, U. Fischer, A. Stephan, P. Müller,
771 D. Durner, *Anal. Chim. Acta* 672 (2010) 114.
- 772 [55] S. Castillo, I. Mattila, J. Miettinen, M. Orešič, T. Hyötyläinen, *Anal.*
773 *Chem.* 83 (2011) 3058.
- 774 [56] M. M. Koek, F. M. van der Kloet, R. Kleemann, T. Kooistra, E. R.
775 Verheij, T. Hankemeier, *Metabolomics* 7 (2011) 1.
- 776 [57] G. T. Ventura, G. J. Hall, R. K. Nelson, G. S. Frysinger, B. R. A. E.
777 Pomerantz, O. C. Mullins, C. M. Reddy, *J. Chromatogr. A* 1218 (2011)
778 2584.

- 779 [58] J. Vial, B. Pezous, D. Thiébaud, P. Sassiati, , B. Teillet, X. Cahours,
780 I. Rivals, *Talanta* 83 (2011) 1295.
- 781 [59] J. Blomberg, P. J. Schoenmakers, U. A. T. Brinkman, *J. Chromatogr.*
782 *A* 972 (2002) 137.
- 783 [60] P. Marriott, R. Shellie, *Trends Analyt. Chem.* 21 (2002) 573.
- 784 [61] L. Mondello, A. C. Lewis, K. D. Bartle, *Multidimensional*
785 *Chromatography*, Wiley, Chichester UK, 2002.
- 786 [62] J. Dallüge, J. Beens, U. A. Brinkman, *J. Chromatogr. A* 1000 (2003)
787 69.
- 788 [63] S. Reichenbach, M. Ni, V. Kottapalli, A. Visvanathan, *Chemometrics*
789 *and Intelligent Laboratory Systems* 71 (2004) 107.
- 790 [64] A. E. Sinha, B. J. Prazen, R. E. Synovec, *Anal. Bioanal. Chem.* 378
791 (2004) 1948.
- 792 [65] P. Q. Tranchida, P. Dugo, G. Dugo, L. Mondello, *J. Chromatogr. A*
793 1054 (2004) 3.
- 794 [66] P. Q. Tranchida, P. Dugo, G. Dugo, L. Mondello, *Trends Analyt.*
795 *Chem.* 26 (2007) 191.
- 796 [67] M. Adahchour, J. Beens, U. Brinkman, *J. Chromatogr. A* 1186 (2008)
797 67.

- 798 [68] S. A. Cohen, M. R. Schure (Eds.), Multidimensional Liquid
799 Chromatography: Theory and Applications in Industrial Chemistry
800 and the Life Sciences, John Wiley and Sons, New York NY, 2008.
- 801 [69] P. Dugo, F. Cacciola, T. Kumm, G. Dugo, L. Mondello, J.
802 Chromatogr. A 1184 (2008) 353.
- 803 [70] L. Mondello, P. Q. Tranchida, P. Dugo, G. Dugo, Mass Spectrom. Rev.
804 27 (2008) 101.
- 805 [71] O. Amador-Muñoz, P. J. Marriott, J. Chromatogr. A 1184 (2008) 323.
- 806 [72] K. M. Pierce, J. C. Hoggard, R. E. Mohler, R. E. Synovec, J.
807 Chromatogr. A 1184 (2008) 341.
- 808 [73] H. J. Cortes, B. Winniford, J. Luong, M. Pursch, J. Sep. Sci. 32 (2009)
809 883.
- 810 [74] L. Ramos, Comprehensive Two Dimensional Gas Chromatography,
811 Elsevier, Oxford UK, 2009.
- 812 [75] J. C. Hoggard, R. E. Synovec, in: L. Ramos (Ed.), Comprehensive Two
813 Dimensional Gas Chromatography, Elsevier, Oxford UK, 2009, pp. 107.
- 814 [76] S. E. Reichenbach, in: L. Ramos (Ed.), Comprehensive Two
815 Dimensional Gas Chromatography, Elsevier, Oxford UK, 2009, pp. 77.
- 816 [77] S. E. Reichenbach, in: L. Ramos (Ed.), Comprehensive Two
817 Dimensional Gas Chromatography, Elsevier, The Netherlands, 2009,
818 pp. 77.

- 819 [78] M. M. Bushey, J. W. Jorgenson, *Anal. Chem.* 62 (1990) 161.
- 820 [79] J. Blomberg, P. J. Schoenmakers, J. Beens, R. Tijssen, J. High
821 Resolut. Chromatogr. 20 (1997) 539.
- 822 [80] C. M. Reddy, T. I. Eglinton, A. Hounshell, H. K. White, L. Xu, R. B.
823 Gaines, G. S. Frysinger, *Environ. Sci. Technol.* 36 (2002) 4754.
- 824 [81] H.-G. Janssen, W. Boers, H. Steenbergen, R. Horsten, E. Flöter, J.
825 Chromatogr. A 1000 (2003) 385.
- 826 [82] R. M. M. Perera, P. J. Marriott, I. E. Galbally, *Analyst* 127 (2002)
827 1601.
- 828 [83] J. L. Hope, B. J. Prazen, E. J. Nilsson, M. E. Lidstrom, R. E. Synovec,
829 *Talanta* 65 (2005) 380.
- 830 [84] R. K. Nelson, B. S. Kile, D. L. Plata, S. P. Sylva, L. Xu, C. M. Reddy,
831 R. B. Gaines, G. S. Frysinger, S. E. Reichenbach, *Environ. Forensics* 7
832 (2006) 33.
- 833 [85] G. D. Wardlaw, J. S. Arey, C. M. Reddy, R. K. Nelson, G. T. Ventura,
834 D. L. Valentine, *Environ. Sci. Technol.* 42 (2008) 7166.
- 835 [86] User's Guide to SIMCA-P, SIMCA-P+, Umetrics AB, version 11.0
836 edition, 2005.
- 837 [87] C. Oh, X. Huang, F. E. Regnier, C. Buck, X. Zhang, *J. Chromatogr. A*
838 1179 (2008) 205.

- 839 [88] S. E. Reichenbach, P. W. Carr, D. R. Stoll, Q. Tao, J. Chromatogr. A
840 1216 (2009) 3458.
- 841 [89] C. Cordero, E. Liberto, C. Bicchi, P. Rubiolo, P. Schieberle, S. E.
842 Reichenbach, Q. Tao, J. Chromatogr. A 1217 (2010) 5848.
- 843 [90] S. Reichenbach, Q. Tao, D. E. Hutchinson, S. B. Cabanban, H. A.
844 Pham, W. E. Rathbun, H. Wang, in: Pittcon, pp. 540.
- 845 [91] O. C. Mullins, G. T. Ventura, R. K. Nelson, S. S. Betancourt,
846 B. Raghuraman, C. M. Reddy, Energy Fuels 22 (2008) 496.
- 847 [92] S. S. Betancourt, G. T. Ventura, A. E. Pomerantz, O. Vilorio, F. X.
848 Dubost, J. Zuo, G. Monson, D. Bustamante, J. M. Purcell, R. K.
849 Nelson, R. P. Rodgers, C. M. Reddy, A. G. Marshall, O. C. Mullins,
850 Energy Fuels 23 (2008) 1178.
- 851 [93] G. T. Ventura, B. Raghuraman, R. K. Nelson, O. C. Mullins, C. M.
852 Reddy, Org. Geochem. 41 (2010) 1026.
- 853 [94] S. Luhn, M. Berth, M. Hecker, J. Bernhardt, Proteomics 3 (2003) 1117.

854 **List of Figures**

855 1 Top – A pseudocolored image of a chromatographic region
856 with BTEX peaks. Bottom – A pseudocolored image of the
857 differences between two aligned chromatograms with red indi-
858 cating a larger value in the reference image, green indicating
859 a smaller value, and grey indicating nearly equal values.[32] . . . 40

860 2 A pseudocolored image of an LC×LC chromatogram of a
861 urine sample. The open circles indicate the retention times
862 of the expected peaks recorded in the template. The outlines
863 indicate the detected peaks and the filled circles indicate the
864 retention times of the apexes of the detected peaks that are
865 matched by the template.[88] 41

866 3 A mesh of regions with automatically drawn vertical lines
867 at linear retention indices based on the *n*-alkanes and hand-
868 drawn crossing lines to separate compound groups.[90] 42

869 4 Cumulative chromatogram for eighteen breast-cancer tumor
870 samples overlaid with the feature template (registration peaks
871 shown with dark ovals and region features shown with red out-
872 lines). The color bar shows the logarithmic pseudocolorization
873 mapping.[1] 43

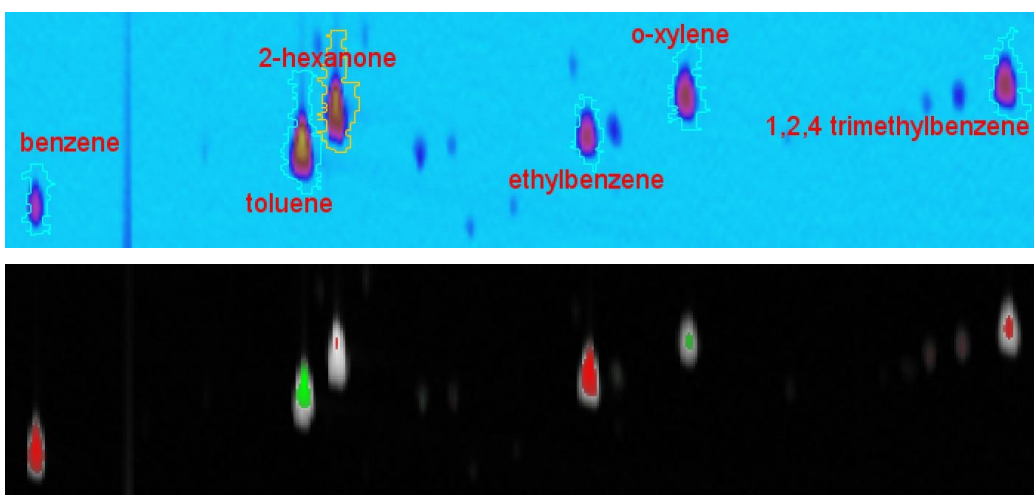


Figure 1: Top – A pseudocolorized image of a chromatographic region with BTEX peaks. Bottom – A pseudocolorized image of the differences between two aligned chromatograms with red indicating a larger value in the reference image, green indicating a smaller value, and grey indicating nearly equal values.[32]

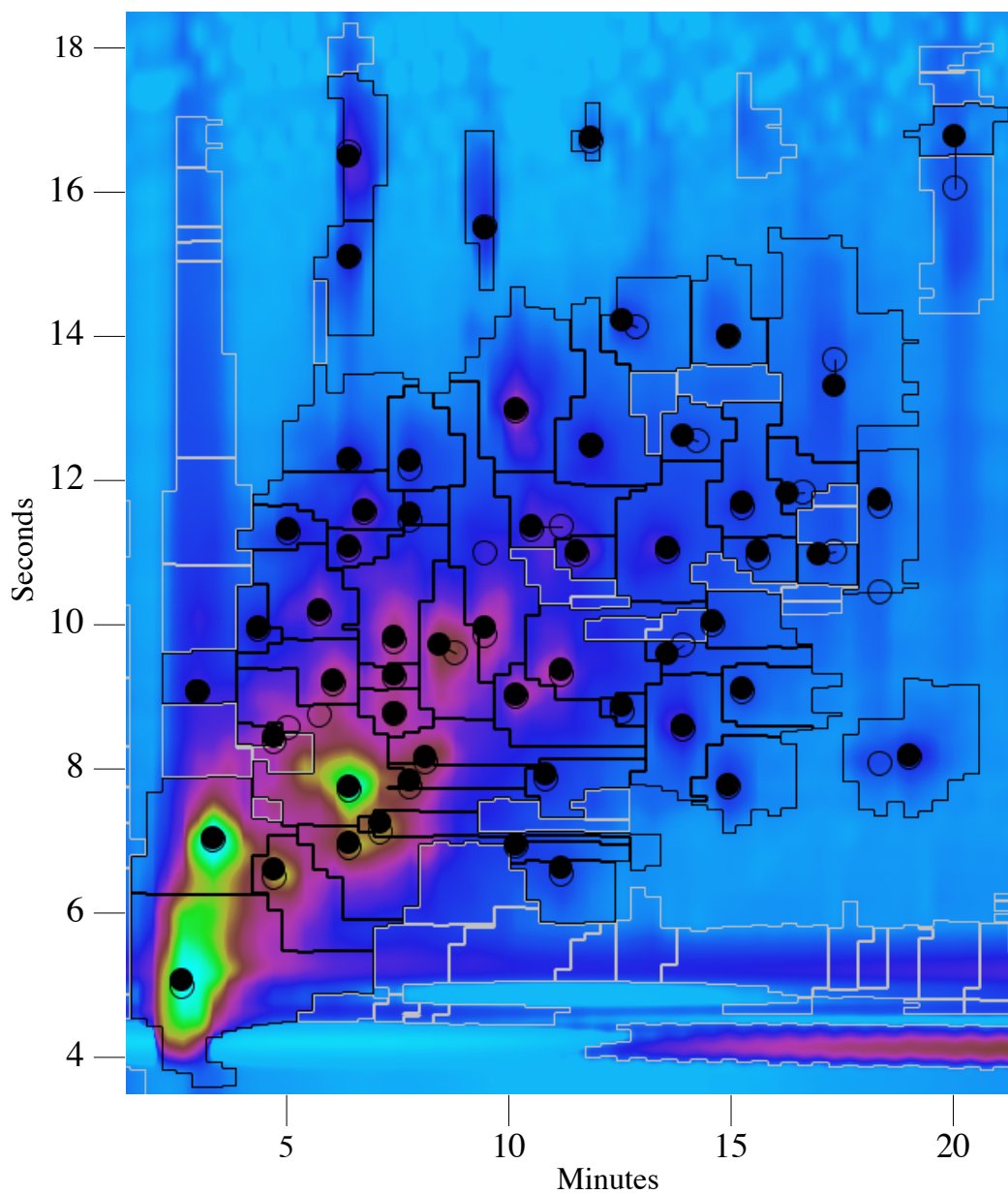


Figure 2: A pseudocolored image of an LC \times LC chromatogram of a urine sample. The open circles indicate the retention times of the expected peaks recorded in the template. The outlines indicate the detected peaks and the filled circles indicate the retention times of the apexes of the detected peaks that are matched by the template.[88]

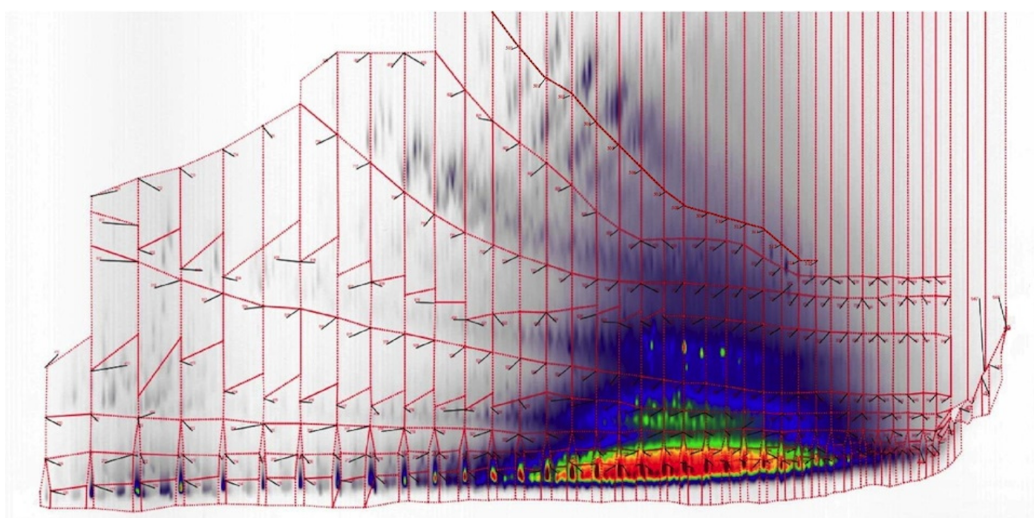


Figure 3: A mesh of regions with automatically drawn vertical lines at linear retention indices based on the *n*-alkanes and hand-drawn crossing lines to separate compound groups.[90]

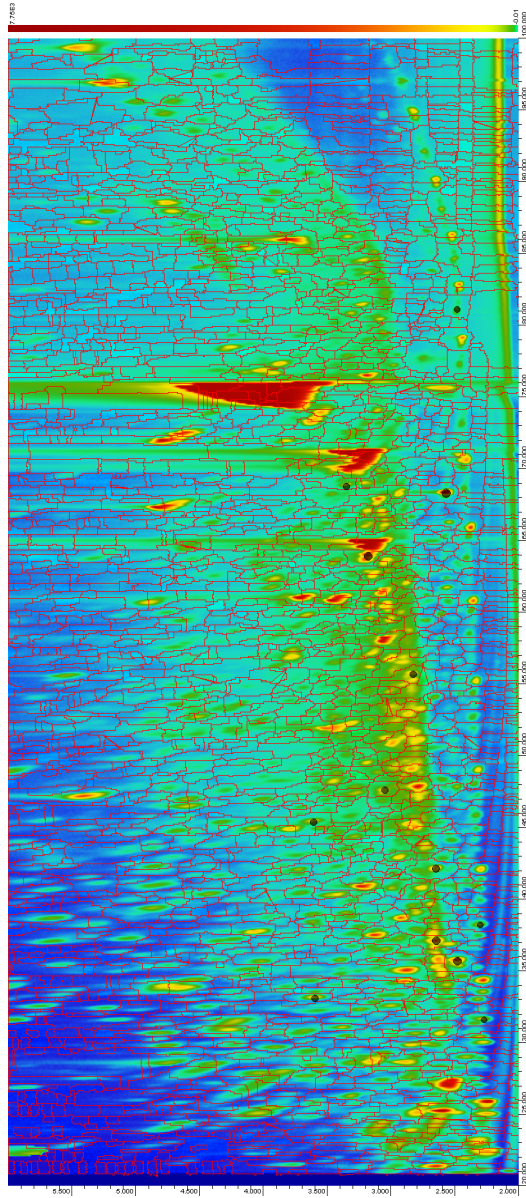


Figure 4: Cumulative chromatogram for eighteen breast-cancer tumor samples overlaid with the feature template (registration peaks shown with dark ovals and region features shown with red outlines). The color bar shows the logarithmic pseudocolorization mapping.[1]