## SUPPLEMENT

# Supplement for: Automated fragment identification for electron ionisation mass spectrometry: application to atmospheric measurements of halocarbons

Myriam Guillevic<sup>1\*</sup>, Aurore Guillevic<sup>2</sup>, Martin K. Vollmer<sup>1</sup>, Paul Schlauri<sup>1</sup>, Matthias Hill<sup>1</sup>, Lukas Emmenegger<sup>1</sup> and Stefan Reimann<sup>1</sup>

\*Correspondence:

myriam.guillevic@empa.ch <sup>1</sup>Laboratory for Air Pollution /Environmental Technology, Empa, Swiss Federal Laboratories for Materials Science and Technology, Ueberlandstrasse 129, 8600 Dübendorf, Switzerland Full list of author information is available at the end of the article

### 1 Training set and validation set: presence of molecular ion

We report here, for each substance used in the training and validation sets, the presence of the molecular ion, and if a mass spectrum is documented in the NIST spectral database [1]. Note that each substance was measured from a real air sample. Both low atmospheric molar fraction and fragmentation pathways may explain the absence of the molecular ion.

Table 1 Known compounds used as training set: presence of the molecular ion. If the molecular ion is absent, we give the detected maximal fragments instead. Note that several maximal fragments may be detected for one substance. The last column indicates if a mass spectrum can be downloaded from the NIST chemistry webbook [1].

	Chemical		Molecular	NIST
Compound	formula	CAS number	ion present	spectrum
C <sub>2</sub> H <sub>6</sub>	C <sub>2</sub> H <sub>6</sub>	74-84-0	Ves	ves
C <sub>3</sub> H <sub>8</sub>	C <sub>3</sub> H <sub>8</sub>	74-98-6	ves	ves
CH₃CI	CH₃CI	74-87-3	ves	ves
COS	COS	463-58-1	ves	ves
NF <sub>3</sub>	NF <sub>3</sub>	7783-54-2	ves	ves
Benzene	$C_6 H_6$	71-43-2	yes	yes
$CH_2CI_2$	$CH_2CI_2$	75-09-2	yes	yes
HCFC-22	$HCF_2CI$	75-45-6	yes	yes
CF <sub>4</sub>	$CF_4$	75-73-0	CF <sub>3</sub>	yes
Toluene	$C_7H_8$	108-88-3	yes	yes
CH₃Br	$CH_3Br$	74-83-9	yes	yes
HCFC-142b	$H_3C_2F_2CI$	75-68-3	$H_2C_2F_2CI$ , $H_3C_2FCI$ , $H_3C_2F_2$	yes
$SO_2F_2$	$SO_2F_2$	2699-79-8	yes	yes
CFC-13	CF <sub>3</sub> CI	75-72-9	$CF_2CI$ , $CF_3$	yes
HCFC-141b	$H_3C_2FCl_2$	1717-00-6	$H_2C_2CI_2$ , $H_3C_2FCI$	yes
CHCl <sub>3</sub>	CHCl <sub>3</sub>	67-66-3	yes	yes
CFC-12	$CF_2CI_2$	75-71-8	yes	yes
$C_2HCI_3$	$C_2HCI_3$	79-01-6	yes	yes
CFC-11	CFCI <sub>3</sub>	75-69-4	yes	yes
HCFC-124	$HC_2F_4CI$	2837-89-0	yes	yes
PFC-116	$C_2F_6$	76-16-4	$C_2F_5$	yes
CH <sub>3</sub> I	CH <sub>3</sub> I	74-88-4	yes	yes
SF6	SF <sub>6</sub>	2551-62-4	SF <sub>5</sub>	no
Halon-1301	CF <sub>3</sub> Br	75-63-8	yes	no
		56-23-5		yes
CFC-115	C <sub>2</sub> F <sub>5</sub> Cl	76-15-3	$C_2F_4CI, C_2F_5$	yes
$CCI_2 = CCI_2$	$C_2CI_4$	127-18-4	yes	yes
Halon-1211	CF <sub>2</sub> CIBr	353-59-3	$CFCIBr, CF_2Br, CF_2CI$	yes
CFC-114	$C_2F_4Cl_2$	76-14-2	$C_2F_3Cl_2, C_2F_4Cl$	yes
CH <sub>2</sub> Br <sub>2</sub>	CH <sub>2</sub> Br <sub>2</sub>	74-95-3	yes	yes
CFC-113		/0-13-1 76 10 7	yes	yes
		10-19-1		yes
$SF_5CF_3$		3/3-8U-8		yes
FFC-C318		110-20-3		yes
		124-13-2		yes
C6F14	C6F14	JJJ-42-U	U5F9	ves

Table 2 Known compounds used as validation set: presence of the molecular ion. If the molecular ion is absent, we give the detected maximal fragments instead. Note that several maximal fragments may be detected for one substance. The last column indicates if a mass spectrum can be downloaded from the NIST chemistry webbook [1].

Compound	Chemical	CAS number	Molecular	NIST			
Compound	formula	CAS number	ion present	spectrum			
Kigali Amendment to the Montreal Protocol							
HFC-41	CH₃F	593-53-3	yes	yes			
HFC-32	$CH_2F_2$	75-10-5	yes	yes			
HFC-152	$C_2H_4F_2$	624-72-6	yes	yes			
HFC-152a	$C_2H_4F_2$	75-37-6	yes	yes			
HFC-23	$CHF_3$	75-46-7	$CF_3$ , $HCF_2$	yes			
HFC-143	$C_2H_3F_3$	430-66-0	yes	yes			
HFC-143a	$C_2H_3F_3$	420-46-2	yes	yes			
HFC-134	$C_2H_2F_4$	359-35-3	yes	yes			
HFC-134a	$C_2H_2F_4$	811-97-2	yes	yes			
HFC-125	$C_2HF_5$	354-33-6	$C_2F_5$ , $HC_2F_4$	yes			
HFC-245ca	$C_3H_3F_5$	679-86-7	$H_2C_3F_3$ , $H_3C_2F_2$ , $HC_3F_4$	yes			
HFC-245fa	$C_3H_3F_5$	460-73-1	yes	no			
HFC-365mfc	$C_4H_5F_5$	406-58-6	$H_4C_4F_3$ , $H_2C_3F_5$	no			
HFC-236cb	$C_3H_2F_6$	677-56-5	$H_2C_3F_5$ , $HC_3F_6$	no			
HFC-236ea	$C_3H_2F_6$	431-63-0	$H_2C_3F_5$	yes			
HFC-236fa	$C_3H_2F_6$	690-39-1	$H_2C_3F_5$	yes			
HFC-227ea	$C_3HF_7$	431-89-0	$HC_3F_6$	no			
HFC-43-10mee	$C_5H_2F_{10}$	138495-42-8	$H_2C_4F_7$ , $HC_5F8$ , $H_2C_5$	no			
		HFOs					
HFO-1234yf	$H_2C_3F_4$	754-12-1	yes	no			
HFO-1234ze(E)	$H_2C_3F_4$	29118-24-9	yes	no			
HCFO-1233zd(E)	$H_2C_3F_3CI$	102687-65-0	yes	no			

#### 2 Mass calibration procedure

A time-of-flight (ToF) instrument measures a time, elapsed between two events: the extraction and the detection, when the ions hit the detector plate. To convert this time measurement into a mass measurement in the most possible accurate manner, internal mass calibration proves to be a good strategy. Known masses detected during a measurement are used to establish the calibration function between ToF and mass. In our cases, we use known masses produced by fragmentation of a mass calibration substance, perfluoroperhydrophenanthrene (PFPHP), of chemical formula  $C_{14}F_{24}$ . Only two types of atoms are present in this molecule, carbon (mass 12.000000) and fluorine (mass 18.99840316). The most abundant peaks can therefore be associated to a unique molecular formula. For example, the peak at integer mass 69 can be associated to  $CF_3^+$  only, of exact mass 68.99466108.

Based on the NIST mass spectrum and our measured mass spectrum for PFPHP, we have chosen a list of xx masses present with a sufficient abundance to be used for mass calibration. In addition, we use the masses of N<sub>2</sub>, O<sub>2</sub>, Ar, which are the most abundant air components and may slightly leak in our system, as well as Cl, also observed to be always present in our detector. The masses are chosen to be evenly distributed between the minimum and maximum masses, covering a range from m/z = 28.0055994348 (N<sub>2</sub><sup>+</sup>) to m/z = 292.98188618 (C<sub>7</sub>F<sub>11</sub><sup>+</sup>).

We have observed that our ToF detector is subject to mass drift of up to 100 ppm during a run. This drift is possibly due to temperature variation of our preceding GC. To correct for this drift, we perform a mass calibration every four minutes, using average data of the preceding and following two minutes. For each set of four-minutes-averaged data, all peaks in the mass domain near the exact masses of the selected list are detected. Note that several mass peaks can be detected where only one exact mass from the calibrant is expected. Each detected peak is fitted using a pseudo-Voigt function, which is a combination of a Gaussian and a Lorentzian function:

$$f(x; A, \mu, \sigma, \alpha, b) = (1 - \alpha) \frac{A}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right) + \alpha \frac{A}{\pi} \left(\frac{\sigma_L}{(x - \mu)^2 + \sigma_L^2}\right) + b \quad (1)$$

where

$$\sigma_L = \sigma \sqrt{2 \ln(2)} \tag{2}$$

and the full-with-at-half-maximum (FWHM)

$$FWHM = 2\sigma\sqrt{2\ln(2)} = 2\sigma_L .$$
(3)

The parameter FWHM is used later on (Section xx) to generate candidate mass peaks with the appropriate peak broadness.

Then, all obtained centres of ToF are associated to the closed expected exact mass. The entire set of pairs (ToF; expected exact mass) is used to fit a calibration function of the form:

$$i_{\rm ToF} = p_1 m^{p_3} + p_2 \tag{4}$$

with  $i_{\text{ToF}}$  the time of flight index (ref TofDaq manual) and m the exact mass. The parameters  $p_1$ ,  $p_2$  and  $p_3$  are optimised using the Python lmfit package. If one or several pairs are further away from the fit than a set maximum value (20 ppm in our case), the furthest away pair is eliminated and the optimisation routine is repeated. This one-by-one pair elimination improves the robustness of the algorithm. This was proven necessary as when measuring real air from the industrial area where Empa is located, once in while a prominent pollution event occurs, producing outstanding mass peaks that may be in the vicinity of the expected peak, even masking it, therefore disturbing the mass calibration function.

Once all residuals are below the set value, the mass calibration is complete. Any ToF value can then be converted to a m/z value using:

$$m = \left(\frac{i_{\text{ToF}} - p_2}{p_1}\right)^{\frac{1}{p_3}} \tag{5}$$

where  $p_1$ ,  $p_2$  and  $p_3$  are calculated at any given specific time as linear interpolation using their time-bracketing optimised values.



#### **3** Uncertainty of the mass calibration

After the optimisation, the obtained fit parameters are used to calculate the reconstructed m/z values; these values are then compared to the expected exact m/z, for each time slice of four minutes. The obtained offsets, expressed in ppm over the mass domain, are displayed in Fig. 1. With our instrument, the observed mass accuracy is better for larger masses, with residuals usually below 5 ppm for masses higher than 100 m/z, while the accuracy can deteriorate to 20 ppm below 50 m/z. This is potentially due to the mass resolution of our instrument that is around 3000 for masses below 50 m/z but around 5000 for masses above.

To reflect this varying mass accuracy over the mass domain, for each used exact mass, we use as uncertainty the maximum observed offset at this mass. Then for any measured mass, its uncertainty is calculated as a linear interpolation between uncertainty at bracketing masses. This constitutes the mass calibration uncertainty.

### 4 Validation set: preparation of qualitative standards for compounds newly regulated by the Kigali amendment to the Montreal Protocol

Eighteen substances listed under the Kigali Amendment to the Montreal Protocol were part of the validation set.

First, substances were separated in two groups, with the aim that each group should not contain isomers, to make sure each substance could be identified by its mass spectrum only. Group A contained: HFC-41, HFC-143a, HFC-134a, HFC-227ea, HFC-236ea, HFC-245fa, HFC-43-10-mee, HFC-152 and HFC-236cb. Group B contained: HFC-32, HFC-23, HFC-125, HFC-152a, HFC-365mfc, HFC-143, HFC-236fa, HFC-245ca and HFC-134.

The pure substances were bought from Synquest Laboratories (Florida, USA). For each group, the pure substances were spiked one after the other into synthetic air, and the mixture was pressurised into a flask. The two obtained mixtures were prepared at approximately  $6.5 \text{ nmol} \cdot \text{mol}^{-1}$ .

Then, each mixture was measured by our preconcentration, gas chromatography, time-of-flight mass spectrometry instrumentation. Data analysis then followed the same procedure as explained in the main article.

#### 5 Algorithmic improvements

We describe our algorithmic improvements to speed-up the running-time of some critical steps.

#### 5.1 Organisation of the sum-formulae in a directed acyclic graph (DAG)

After running the knapsack algorithm, many candidate sum-formulae are obtained. We organise them in a graph. In this graph, a node  $n_j$  is a descendant of a node  $n_i$ if the node's fragment  $s_j$  is a sub-fragment of the fragment  $s_i$  of the node  $n_i$ . For example, CCl is a descendant of CCl<sub>3</sub>. Conversely, a node  $n_i$  is an ancestor of a node  $n_j$  if its fragment  $s_i$  is a sup-fragment of the fragment  $s_j$ . This define a *partial* order on the chemical formulas that we formally define in Definition 1.

**Definition 1** (Partial order) We define the following partial order on the chemical formulas. Let  $s_i$  and  $s_j$  be two sum-formulas encoded as vectors of non-negative integers.

- The sum formula s<sub>j</sub> is smaller than the sum-formula s<sub>i</sub>, denoted s<sub>j</sub> ≤ s<sub>i</sub>, if s<sub>j</sub> is a sub-fragment of s<sub>i</sub>;
- The sum formula  $s_j$  is greater than the sum-formula  $s_i$ , denoted  $s_j \ge s_i$ , if  $s_j$  is a sup-fragment of  $s_i$ ;
- otherwise  $s_i$  and  $s_j$  are *incomparable*.

Organising a set S of n items (here fragments) in a DAG (directed acyclic graph) can takes as many as  $n^2$  comparisons of items, but since S can be quite large (e.g., n = 10000), it makes sense to reduce the number of comparisons. Moreover, the graph should have as less edges as possible, that is, two sum formulae  $s_i \ge s_j$  are binded with an edge if and only if there is no other sum-formula  $s_{i'}$  that could be inserted between them like  $s_i \ge s_{i'} \ge s_j$ . For example with CCl, CCl<sub>2</sub> and CCl<sub>3</sub>,



we will define the graph with minimal edges on the left, not the one of the right (Fig. 2).

We now explain how we reduced the number of comparisons between fragments to set the edges, thus improving the complexity of building the graph. First the target masses are sorted in decreasing order before the knapsack step. Hence the output of the knapsack is made of batches of chemical formulae, each one for a given target mass, in decreasing order. Because the mass uncertainty is far thinner than 0.5m/z, all sum-formulae for one target mass are incomparable: it is not possible to find that a sum-formula is a sub-fragment of another, for the same target mass, otherwise the mass difference between the two would be at least 1m/z.

Therefore we have a list of sum-formulae  $\{s_i\}_{1 \le i \le \#S}$  such that for any sum-formula  $s_i$ , the sum-formulae in the preceeding batches weight more and are either incomparable or contain  $s_i$  as a sub-fragment, and the sum-formulae in the forthcoming batches are lighter and either incomparable or sub-fragments of  $s_i$ . The maximal fragments will be the nodes at the "roots" of the DAG. They are made of the sum-formulae of the first (heaviest) batch, and some other sum-formulae from other batches.

We maintain a list of the root nodes (maximal fragments) of the graph. They have no ancestor, and they are incomparable to each other. To add a new sum-formula in the graph, because of the ordering of the sum-formulae, we know that it is either incomparable, or a subfragment of any node of the graph. It cannot be a sup-fragment (a parent) of any node of the graph. We compare the new sum-formula to each of the maximal fragments. If it is incomparable to any of them, we add it as as new maximal fragment. Otherwise, for each maximal fragment that has the new sum-formula as sub-fragment, we compare the new one to its children. If it is incomparable to any of the children, we add it as a new child of the maximal fragment (we add an edge toward it). Otherwise, we recursively explore the children of the children that have this new sum-formula as sub-fragment. In this way, we avoid many useless comparisons: all children of incomparable nodes are omitted.

Thanks to the list of maximal fragments, the singletons are identified right away: they are the maximal fragments without any child. Figure 2n the main article shows the graph obtained for CCl<sub>4</sub>.

#### 5.2 Removing a node and updating the edges

A node n to be removed has parents (closest sup-fragments) connected with one edge, and children (closest sub-fragments) connected with one edge. We wrote this

procedure to remove a node and update the edges and list of maximal fragments (Alg. 1, example in Fig. 3).

Algorithm 1: remove the node n and update the edges for each parent p<sub>i</sub> of n do remove the edge from p<sub>i</sub> to n; for each child c<sub>i</sub> of n do lists its own parents q<sub>j</sub> (without n); if none of q<sub>j</sub> is a sub-fragment of p<sub>i</sub> then add an edge from p<sub>i</sub> to c<sub>i</sub> (if there is one q<sub>j</sub> which is a sub-fragment of p<sub>i</sub>, do nothing: the edges are already fine) for each child c<sub>i</sub> of n do remove the edge from n to c<sub>i</sub>; if c<sub>i</sub> has no more parent after removing n then add it in the list of maximal fragments if n is a maximal fragment (it has no parent) then remove it from the list of maximal fragments Remove n



#### 5.3 Enumeration of isotopocules

We now recall the computation of the relative intensities of the rare isotopocules (see [2] for a detailed computation). The abundant formula has proportion (of the set of all isotopocules)  $pr = \prod_{\{\text{element } e\}} (a_{e,0})^{n_e}$ . The product is over all the distinct atoms (denoted e),  $a_{e,0}$  is the abundance of the most abundant isotope of an atom, and  $n_e$  is the number of occurrences of that atom in the chemical formula. For example, with CCl<sub>4</sub> one computes  $pr = a_{\text{C}}a_{\text{Cl}}^4 = 0.326$ . An isotopocule with only one element e and i rare isotopes of abundance  $a_{e,i}$  has proportion

$$pr_{e} = a_{e,0}^{n_{e,0}} \binom{n_{e}}{n_{e,0}} a_{e,1}^{n_{e,1}} \binom{n_{e-n_{e,0}}}{n_{e,1}} a_{e,2}^{n_{e,2}} \binom{n_{e-n_{e,0}-n_{e,1}}}{n_{e,2}} \cdots a_{e,i}^{n_{e,i}} \binom{n_{e-n_{e,0}-n_{e,1}-\cdots-n_{e,i-1}}}{n_{e,i}} = a_{e,0}^{n_{e,0}} a_{e,1}^{n_{e,1}} a_{e,2}^{n_{e,2}} \cdots a_{e,i}^{n_{e,i}} \frac{n_{e}!}{n_{e,0}!n_{e,1}!n_{e,2}!\cdots n_{e,i}!}$$

where the terms  $\binom{n}{m} = \frac{n!}{m!(n-m)!}$  are binomal coefficients with  $n \ge m$ , denoting the number of ways to choose m items in a set of size n. Their product simplifies in  $n_e!/(n_{e,0}!n_{e,1}!n_{e,2}!\cdots n_{e,i}!)$ . An isotopocule has proportion

$$pr = \prod_{e} n_{e}! \prod_{i} (a_{e,i})^{n_{e,i}} / (n_{e,i}!)$$
(6)

where e ranges over the elements, i ranges over the isotopes of an element,  $n_e$  is the total number of an element (with all isotopes),  $n_{e,i}$  is the number of occurrences of

one isotope. The relative intensity of a rare formula is the ratio (see [2])

$$p = \frac{\prod_{e} n_{e}! \prod_{i} (a_{e,i})^{n_{e,i}} / (n_{e,i}!)}{\prod_{e} (a_{e,0})^{n_{e}}} = \prod_{e} n_{e}! \prod_{i>0} \left(\frac{a_{e,i}}{a_{e,0}}\right)^{n_{e,i}} \frac{1}{n_{e,i}!} .$$
(7)

One needs to enumerate all the possible combinations of isotopes of a given element. This is a classical problem in combinatorics. Denote by i the number of isotopes (the abundant one included). Enumerate all the ways to sum at most i positive integers to obtain  $n_e$ . It can be seen as a knapsack-like problem: given all isotopes each of "weight" 1, finds how to sum to  $n_e$ . In particular, each isotope is allowed at most  $n_e$  times.

All ratios are relative to the abundant sum-formula, whose maximum possible intensity is known: this is the intensity of the corresponding measured mass. To improve the running-time of the enumeration, we do not list the rare isotopocules whose intensity would be below the detection threshold of the ToF-MS. For this purpose, we consider the elements one after the other. We maintain a list of partial isotopocules with their partial relative intensity, made of the elements processed so far. The list is ordered by decreasing relative intensity. Given a new element eand its occurence  $n_e$ , we generate its isotopic patterns, the abundant one included, and sort them in decreasing order of relative intensity. Reading the two lists in decreasing order of relative intensity, we combine the new isotopes to the partial solutions and multiply together the intensities. A loop over a list stops as soon as the product of intensities is below the partial threshold. Once the new list of partial solutions is computed, it is sorted in decreasing order of relative intensity. Then, the next element is processed in the same manner, until all elements are done. The first item of the resulting list is the isotopocule of highest relative intensity (it can be greater than one). For implementation purpose, we scale the list and divide all numbers by this highest relative intensity, so that everything is in the interval [0, 1].

#### 6 Full Numerical Example with carbon tetrachloride

For  $CCl_4$  found at a retention time of 1708.27 s, nineteen masses are observed, listed in Table 7 with uncertainty and intensity.

# 6.1 Knapsack algorithm with two lists but without considering separately multi-valent and mono-valent atoms

In this paragraph we present an example of a knapsack algorithm with two sets of atoms, and two lists of intermediate masses. The candidate atoms are H, C, N, O, F, S, Cl, Br, I. We arbitrarily define two subsets: {C, N, O, S, Br} and {H, F, Cl, I}. The knapsack algorithm is run with input the masses of the atoms of each set, and the minimal mass is set to 1. One obtains two lists: A and B given in Table 5. The lists are sorted in increasing order of mass. One obtains  $A = [(12.0, C), (14.0030740074, N), (15.9949146223, O), (24.0, C_2), (26.003074007400002, CN), (27.9949146223, CO), (28.0061480148, N_2), (29.9979886297, NO), (31.97207073, S), (31.9898292446, O_2)] and <math>B = [(1.0078250319, H), (2.0156500638, H_2), \dots (all H_{3...17}), (18.140850574199998, H_{18}), (18.99840316, F), (19.1486756061, H_{19}), (20.0062281919, HF), \dots (all H_{20...33} and H_{2...14}F), (34.1157786385, H_{15}F), (34.2660510846, H_{34}),$ 

(34.96885271, Cl)]. Then pairing the masses of the two lists, one obtains masses within the target interval (if any): there is one solution (34.96885271, Cl). Finally, the DBE is computed and solutions with a negative DBE are discared. This first approach has a major drawback: many impossible sub-fragments are enumerated, im particular because of the Hydrogen which as a very small mass, and is mono-valent.

Here is a detailed example to compute the lists A and B for the first target mass  $m_{\rm max} = 34.97006625322406$ . First one computes multiples of each mass up to  $m_{\rm max}$ . One obtains Tables 3 and 4. Then one combines at most one mass per column, starting the enumeration with the lightest mass of each column (the first row value). One obtains the fragments of the first row of Table 5.

Table 3 Initial partial list of masses before computing the list A made of fragments of atoms {C, N, O, S, Br} and masses up to  $m_{\rm max}=$  34.97006625322406.

#	C	N	0	S	Br
1	12	14.0030740074	15.9949146223	31.97207073	
2	24	28.0061480148	31.9898292446		

Table 4 Initial partial list of masses before computing the list B made of fragments of atoms {H, F, Cl, I} and masses up to  $m_{\rm max} =$  34.97006625322406. There are 34 multiples of Hydrogen: H, H<sub>2</sub>, up to H<sub>34</sub>.

	Н	F	CI	
1	1.0078250319	18.99840316	34.96885271	
2	2.0156500638			
•				
34	34.2660510846			

**Table 5** Intermediate lists in the knapsack algorithm with two sets of atoms {C, N, O, S, Br} for the list A and {H, F, Cl, I} for the list B. The notation  $H_{1-34}$  means all the 34 fragments made of one to 34 atoms of Hydrogen. The notation  $C_{1-2}O$  means CO and  $C_2O$ .

target mass interval	#A	A	#B	В	all	sol.	DE	$BE \ge 0$
34.96751070677594,	10	S, O <sub>1-2</sub> , NO, CO,	51	CI, H <sub>1-15</sub> F, F,	1	CI	1	CI
34.97006625322406		N <sub>1-2</sub> , CN, C <sub>1-2</sub>		H <sub>1-34</sub>				
35.974137795964566,	10	S, O <sub>1-2</sub> , NO, CO,	54	HCI, CI,	1	HCI	1	HCI
35.97778836403543		N <sub>1-2</sub> , CN, C <sub>1-2</sub>		H <sub>1-16</sub> F, F,				
				H <sub>1-35</sub>				
36.96406557296814,	11	S, O <sub>1-2</sub> , NO, CO,	56	HCI, CI,	0		0	
36.967505987031856		N <sub>1-2</sub> , CN, C <sub>1-3</sub>		H <sub>1-17</sub> F, F,				
				H <sub>1-36</sub>				
46.96648952357635,	21	NS, CS, S, $NO_2$ ,	95	H <sub>1-11</sub> Cl, Cl,	1	CCI	1	CCI
46.970287436423654		CO <sub>2</sub> , O <sub>1-2</sub> , N <sub>2</sub> O, NO,		H <sub>1-8</sub> F <sub>2</sub> , F <sub>1-2</sub> ,				
		CNO, C <sub>1-2</sub> O, N <sub>1-3</sub> ,		H <sub>1-27</sub> F, H <sub>1-46</sub>				
		CN <sub>2</sub> , C <sub>1-2</sub> N, C <sub>1-3</sub>						
48.962558174089345,	24	OS, NS, CS, S, O <sub>1-3</sub> ,	103	H <sub>1-13</sub> Cl, Cl,	0		0	
48.96840118591066		$NO_2$ , $CO_2$ , $N_2O$ , $NO_2$		$H_{1-10}F_2, F_{1-2},$				
		CNO, C <sub>1-2</sub> O, N <sub>1-3</sub> ,		H <sub>1-29</sub> F, H <sub>1-48</sub>				
		CN <sub>2</sub> , C <sub>1-2</sub> N, C <sub>1-4</sub>						

#### 6.2 Faster knapsack: avoiding enumerating fragments of negative DBE value

With two arbitrary lists, many impossible partial fragments are enumerated, in particular with too many hydrogens. It would speed-up the process to know an upper bound on the number of mono-valent atoms. To be able to compute such value, the first list, denoted M, is now made of the multi-valent atoms and the second list, denoted m, made of mono-valent atoms only. In this way, after computing the list M, one can compute the DBE value of each partial fragment made of multi-valent atoms only. An upper bound on the number of mono-valent atoms is two times

**Table 6** Intermediate lists in the knapsack algorithm with two sets of atoms {C, N, O, S, Br} for the list A and {H, F, Cl, I} for the list B, then with a set of multi-valent atoms {C, N, O, S} giving the list M, and a set of mono-valent atoms {H, F, Cl, Br, I} giving the list m. After enumerating the list M, the maximum possible valence is computed and used as an upper bound on the number of mono-valent atoms to generate the list m. One observes that this technique allows to divide by more than two the length of the second list.

target mass interval	#A	#B	#	#M	2 max	#m	DBE
	,,,,,	<i>11</i> –	sol.		DBE	<i>\_</i>	$\geq 0$
34.96751070677594, 34.97006625322406	10	51	1	10	6	13	1
35.974137795964566, 35.97778836403543	10	54	1	10	6	14	1
36.96406557296814, 36.967505987031856	11	56	0	11	8	18	0
46.96648952357635, 46.970287436423654	21	95	1	21	8	31	1
48.962558174089345, 48.96840118591066	24	103	0	24	10	39	0
59.960220186945, 59.971315773055004	37	156	1	37	10	48	1
81.9330864132061, 81.9395331467939	90	315	1	89	14	110	1
82.93831758560036, 82.95111397439963	94	324	3	93	14	115	2
83.93024931474109, 83.9372426452589	94	333	0	93	14	117	0
84.92634272134954, 84.97112963865045	101	342	7	100	16	135	4
85.92419323227152, 85.93924312772849	101	351	1	100	16	137	1
97.9236485334006, 97.9389656265994	154	477	3	150	18	189	2
99.90729357057015, 99.94127718942985	161	501	4	157	18	197	3
116.90200574284233, 116.90847941715768	274	735	2	262	20	293	1
117.89455759263474, 117.92205636736527	275	751	5	262	20	300	3
118.89897942315969, 118.9056775368403	287	766	0	274	20	305	0
119.89648859190275, 119.91785116809724	290	782	2	275	20	311	1
120.89523104367791, 120.90302931632209	306	798	1	291	22	345	0
122.8875535007597, 122.90537265924031	323	830	1	305	22	357	1

Table 7 Measured masses at RT = 1708.27s. In blue, the correct guess made by knapsack. In orange, the identified isotopocules.

measured	mana un ma unit	h	intensity	luna nanalı	id a matifie of
mass (m/z)	mass range with	n uncertainty	intensity	кпарѕаск	Identified
34.96878848	[ 34.96751070677594,	34.97006625322406]	2722.2042	CI	CI
35.97596308	35.974137795964566,	35.97778836403543	1051.6898	HCI	HCI
36.96578578	36.96406557296814,	36.967505987031856]	914.6638	-	[ <sup>37</sup> CI]
46.96838848	46.96648952357635,	46.970287436423654	3784.4981	CCI	CCI 1
48.96547968	48.962558174089345,	48.96840118591066]	1192.8077	-	C[ <sup>37</sup> Cl]
59.96576798	59.960220186945,	59.971315773055004]	120.657	COS	-
81.93630978	81.9330864132061,	81.9395331467939]	6160.9695	CCI <sub>2</sub>	CCl <sub>2</sub>
82.94471578	[ 82.93831758560036,	82.95111397439963]	319.247	S <sub>2</sub> F, HCCl <sub>2</sub>	HCCl <sub>2</sub>
83.93374598	83.93024931474109,	83.9372426452589]	3956.1947	-	CCI[ <sup>37</sup> CI]
84.94873618	84.92634272134954,	84.97112963865045]	140.2337	$H_2S_2F$ , $H_2OSCI$ ,	HCCI <sup>37</sup> CI
	-	_		$HNCI_2, CF_2CI$	
85.93171818	[ 85.92419323227152,	85.93924312772849]	564.31	OCl <sub>2</sub>	$C[^{37}CI]_2$
97.93130708	[ 97.9236485334006,	97.9389656265994]	134.1543	H <sub>2</sub> S <sub>3</sub> , COCl <sub>2</sub>	COCl <sub>2</sub>
99.92428538	99.90729357057015,	99.94127718942985]	106.7792	$HS_2CI$ , $HO_2SCI$ ,	COCI[ <sup>37</sup> CI]
	-	_		NOCl <sub>2</sub>	
116.90524258	[116.90200574284233, 1	.16.90847941715768]	28974.7117	CCI <sub>3</sub>	CCl <sub>3</sub>
117.90830698	[117.89455759263474, 1	.17.92205636736527]	189.527	$S_2FCI$ , $OSCI_2$ ,	[ <sup>13</sup> C]Cl <sub>3</sub>
				HCCI <sub>3</sub>	
118.90232848	[118.89897942315969, 1	.18.9056775368403]	29078.7276	-	CCl <sub>2</sub> [ <sup>37</sup> Cl]
119.90716988	[119.89648859190275, 1	.19.91785116809724]	182.0316	$C_2S_3$	[13C]Cl <sub>2</sub> [ <sup>37</sup> Cl]
120.89913018	[120.89523104367791, 1	.20.90302931632209]	9220.1959	-	$CCI[^{37}CI]_2$
122.89646308	[122.8875535007597, 1	.22.90537265924031]	886.6747	CSBr	C[ <sup>37</sup> Cl] <sub>3</sub>
	-	-			

the maximum DBE value obtained over the list M. This upper bound allows to constraint the enumeration of the list m, hence reducing the running-time and the length of the second list. Table 6 presents a comparison of the lengths of the lists A, B, M and m for the target masses of  $CCl_4$ . We observed that the list m is at least two times smaller than the list B.

The chosen possible atoms are H, C, N, O, F, S, Cl, Br, I. At start, only the abundant atoms are considered. The multi-valent atoms are C, N, O, S, the mono-valent are H, F, Cl, Br, I.

Consider now the target mass m = 116.90524258 m/z, with uncertainty range  $m_{\min} = 116.90200574284233$ ,  $m_{\max} = 116.90847941715768$ . Our knapsack algorithm first lists all possible sum-formulae made of any number of the multi-valent atoms and so that the mass of the fragment is at most  $m_{\max}$ . There are 263 combinations whose DBE ranges from to 2 to 20. For each fragment, its DBE value n is computed and a second knapsack algorithm is run to find a complement fragment made of at most n mono-valent atoms so that the total mass fits within the bounds  $m_{\min}, m_{\max}$  and the total DBE is positive or zero. For the mass m = 116.90524258 m/z, there is one solution: CCl3.

To account for sum-formulae made of mono-valent atoms only, a final knapsack algorithm is run to find the sum-formulae with only one or two mono-valent atoms whose mass fits in the uncertainty range (this gives the solution Cl for the mass 34.96878848).

isotopocule	mass (m/z)	proportion	relative intensity
С	12.00000000	0.988922	1.000000
[ <sup>13</sup> C]	13.00335484	0.011078	0.011202
CI	34.96885271	0.757647	1.000000
[ <sup>37</sup> Cl]	36.96590260	0.242353	0.319876
CCI <sub>4</sub>	151.87541084	0.325859	1.000000
CCl <sub>3</sub> [ <sup>37</sup> Cl]	153.87246073	0.416938	1.279504
CCl <sub>2</sub> [ <sup>37</sup> Cl] <sub>2</sub>	155.86951062	0.200052	0.613923
$CCI[^{37}CI]_3$	157.86656051	0.042661	0.130920
C[ <sup>37</sup> Cl] <sub>4</sub>	159.86361040	0.003412	0.010470
[ <sup>13</sup> C]Cl <sub>4</sub>	152.87876568	0.003650	0.011202
[ <sup>13</sup> C]Cl <sub>3</sub> [ <sup>37</sup> Cl]	154.87581557	0.004671	0.014333
[ <sup>13</sup> C]Cl <sub>2</sub> [ <sup>37</sup> Cl] <sub>2</sub>	156.87286546	0.002241	0.006877
[ <sup>13</sup> C]CI[ <sup>37</sup> CI] <sub>3</sub>	158.86991535	0.000478	0.001467
[ <sup>13</sup> C][ <sup>37</sup> Cl] <sub>4</sub>	160.86696524	0.000038	0.000117

Table 8 isotopocules of  $CCl_4$  and relative intensity. See also Fig. 3n the main article.

#### 7 Results for the validation set

Results for the validation set are given in Fig. 4 and Fig. 5.

#### Author details

<sup>1</sup>Laboratory for Air Pollution /Environmental Technology, Empa, Swiss Federal Laboratories for Materials Science and Technology, Ueberlandstrasse 129, 8600 Dübendorf, Switzerland. <sup>2</sup>Université de Lorraine, CNRS, Inria, LORIA, 54000 Nancy, France.

#### References

- 1. NIST: NIST EPA NIH Mass Spectral Library. Online database (2020).
- https://www.nist.gov/srd/nist-standard-reference-database-1a-v17 Accessed 11.03.2020
- Yergey, J.A.: A general approach to calculating isotopic distributions for mass spectrometry. International Journal of Mass Spectrometry and Ion Physics 52, 337–349 (1983). doi:10.1002/jms.4498. e4498 JMS-20-0003







