



HAL
open science

Detection of minor compounds in complex mineral samples from millions of spectra: A new data analysis strategy in LIBS imaging

Alessandro Nardecchia, Cécile Fabre, Jean Cauzid, Frederic Pelascini, Vincent Motto-Ros, Ludovic Duponchel

► To cite this version:

Alessandro Nardecchia, Cécile Fabre, Jean Cauzid, Frederic Pelascini, Vincent Motto-Ros, et al.. Detection of minor compounds in complex mineral samples from millions of spectra: A new data analysis strategy in LIBS imaging. *Analytica Chimica Acta*, 2020, 1114, pp.66-73. 10.1016/j.aca.2020.04.005 . hal-02616981

HAL Id: hal-02616981

<https://hal.science/hal-02616981>

Submitted on 20 May 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

1 **Detection of minor compounds in complex mineral samples from** 2 **millions of spectra: a new data analysis strategy in LIBS imaging**

3 Alessandro Nardecchia[†], Cécile Fabre[‡], Jean Cauzid[‡], Frédéric Pelascini[§], Vincent Motto-Ros^{‡*} and
4 Ludovic Duponchel^{†*}

5 [†] Univ. Lille, CNRS, UMR 8516 – LASIRE – Laboratoire de Spectroscopie pour les Interactions, la
6 Réactivité et l'Environnement, F-59000 Lille, France.

7 [‡] Université de Lorraine, Laboratoire GeoRessources, UMR CNRS 7359, France.

8 [§] Cetim Grand Est, Illkirch-Graffenstaden, France.

9 [‡] Institut Lumière Matière, UMR 5306, Université Lyon 1 - CNRS, Université de Lyon 69622 Vil-
10 leurbanne, France.

11

12 **Corresponding Authors**

13 * ludovic.duponchel@univ-lille.fr and vincent.motto-ros@univ-lyon1.fr

14 **Keywords :** Laser-Induced Breakdown Spectroscopy (LIBS), Big data, Hyperspectral imaging,
15 Clustering.

16

17 **ABSTRACT:** Today, Laser-Induced Breakdown Spectroscopy (LIBS) imaging is in full change.
18 Indeed, always more stable instrumentations are developed, which significantly increases the signal
19 quality and naturally the analytical potential of the technique for the characterization of complex and
20 heterogeneous samples at the micro-scale level. Obviously, other intrinsic features such as a limit of
21 detection in the order of ppm, a high field of view and high acquisition rate make it one of the most
22 complete chemical imaging techniques to date. It is thus possible in these conditions to acquire sev-
23 eral million spectra from one single sample in just hours. Managing big data in LIBS imaging is the
24 challenge ahead. In this paper, we put forward a new spectral analysis strategy, called embedded k-

25 means clustering, for simultaneous detection of major and minor compounds and the generation of
26 associated localization maps. A complex rock section with different phases and traces will be ex-
27 plored to demonstrate the value of this approach.

28 **INTRODUCTION**

29 Laser-induced breakdown spectroscopy (LIBS) imaging is actually becoming an essential
30 tool to characterize complex samples in many scientific domains [1–5]. In this spectroscopic tech-
31 nique, a pulse laser beam focused on the sample surface generates a plasma from a small amount of
32 vaporized material. Due to the electronic relaxation of excited atoms and ions, an emission spectrum
33 characteristic of the elemental composition of the sample can be acquired using an optical spectrom-
34 eter. In LIBS imaging experiments, the sample surface is explored in a raster scanning mode (i.e.
35 acquisition of one spectrum for each spatial position of a predefined grid) covering the region of
36 interest. An elemental image can then be generated from the acquired data set using a simple signal
37 integration of a given emission line. The richness of this imaging approach lies in its many ad-
38 vantages that cannot be observed simultaneously in any other spectroscopic technique. Indeed, LIBS
39 imaging has multi-elemental capabilities, a high acquisition rate (≥ 100 spectra/s), full compatibility
40 with optical microscopy and ease of use on samples without almost any size restriction (up to several
41 tens of cm^2), all under atmospheric conditions. On top of that, this technique has a high field of view
42 and a spatial resolution around $10\ \mu\text{m}$ coupled with a limit of detection in the order of weight ppm. It
43 is thus very convenient to explore a sample at the micronic scale by acquiring several million spectra
44 in just hours.

45 Concerning data analysis in LIBS, we see today big differences between the two frameworks
46 of bulk analysis and imaging. Indeed, researchers have quickly learned that multivariate data analy-
47 sis could bring valuable tools for qualitative and quantitative explorations of samples at the bulk
48 level, for instance by developing regression or classification models [6–10]. At the imaging level,
49 there is a relatively limited number of papers dealing with the use of multivariate data analysis in the
50 LIBS community. Indeed, elemental images are, in general, generated from single emission wave-

51 lengths, even though the whole spectral domain could be used. The application of chemometric ap-
52 proaches to imaging data sets is in fact more complex, both from a conceptual and practical point of
53 view. Although a large part of the LIBS community is increasingly sensitive to the use of chemomet-
54 ric tools, understanding the concept of hyperspectral imaging, finding appropriate tools for data ex-
55 ploration, and finally interpreting their outputs represent a big task for non-expert researchers. In
56 addition, it is clear that managing millions of spectra increases the difficulty of this task even if they
57 know the great potential of chemometrics. This is not just about the availability of computational
58 resources, but also, the development of new data exploration tools able to manage such big data
59 structures.

60 In this paper, the idea is obviously not to systematically apply a well-known unsupervised
61 classification method to a LIBS imaging data set. Indeed, it would be totally inefficient in detecting
62 minor compounds because most chemometric algorithms exploit explained variances. As a conse-
63 quence, we will introduce a new data processing strategy, that we call embedded k-means clustering,
64 in order to detect and localize simultaneously major and minor compounds in a complex mineral
65 sample from a data set of more than 2 million spectra.

66 **EXPERIMENTAL SECTION**

67

68 **Sample description and preparation**

69 In order to demonstrate the potential of our strategy of spectroscopic exploration, we have se-
70 lected a complex mineral sample from the polymetallic W–Au–Pb–Zn–Ag (Sb–Ba) district of
71 Tighza (Central Morocco). More specifically, it is related to the Sidi Ahmed hydrothermal event
72 [11]. This district has been mined for centuries for Pb and Ag, Pb–Zn–Ag mineralization being
73 formed of sulfides in gangues of carbonates. Naturally, we can expect the simultaneous presence of
74 major and minor compounds but also traces in such mineralization [3]. The size of the selected rock
75 section is approximately 3.2 cm x 1.6 cm and 1 cm thick. Prior to LIBS analysis, the surface of the
76 sample has been finely polished using polisher as it is usually done in other techniques such as
77 Scanning Electron Microscopy (SEM) and Energy Dispersive X-Ray Spectroscopy (EDS) .

78 **Experimental setup and spectral data acquisition**

79 The LIBS instrumental setup used in this work is based on a homemade optical microscope
80 and a Nd:YAG laser (Centurion GRM, Quantel by Lumibird) with an 8 ns pulse duration operating
81 at 100 Hz. The laser beam is focused on the sample surface using a 15x magnification objective
82 (LMM-15X-P01, Thorlabs). The rock section is placed on a three axes XYZ motorized stage in order
83 to move precisely the sample during the mapping experiment. Atomic force microscopy (AFM) has
84 been used in order to check that the crater size after ablation was smaller than the distance between
85 two consecutives acquisition positions on the sample which is 15 μm . An autofocus system is also
86 used during the analysis in order to keep the objective-to-sample distance from changing. Thus, we
87 always have the same distance between the objective and the plasma emission regardless of the sam-
88 ple flatness. Every spectra in the data set have been acquired from single laser pulses at each spatial
89 position of the sample. The plasma emission has been collected by a quartz lens and focused onto the
90 entrance of a round-to-linear fiber bundle (19 fibers with a 200- μm core diameter) connected to a
91 Czerny-Turner spectrometer (Shamrock 500, Andor Technology). This spectrometer is equipped
92 with a 600 l/mm grating blazed at 300 nm and an intensified charge-coupled device (ICCD) camera
93 (iStar, Andor Technology). The camera is synchronized with the Q-switch of the laser, and spectra
94 are acquired with a delay of 500 ns and a gate of 3000 μs , in full vertical binning mode. Moreover, a
95 servo control loop based on a power meter and a computer-controlled attenuator (ATT1064, Quan-
96 tum Composers) is used to control the laser power. A homemade software, developed under Lab-
97 VIEW® environment, has allowed the automation of scanning sequence as well as the spectral ac-
98 quisition. All measurements have been performed at room temperature under ambient pressure con-
99 ditions.

100 The hyperspectral LIBS data set has been acquired considering a 15 μm spatial resolution
101 and a 0.15 nm spectral resolution. The 251.38 - 339.99 nm spectral domain (2048 spectral channels)
102 has been selected to cover the main emission lines of all elements of interest. In these conditions, we
103 have obtained a data cube of size 2100 pixels x 1090 pixels x 2048 wavelengths (i.e. 2.289.000 ac-
104 quired spectra for a 515 mm^2 field of view). The total acquisition time was approximately six hours,

105 which is finally not so long regarding the richness of the chemical information. It is then easy to un-
106 derstand that a specific data analysis strategy must be implemented if we really want to extract in-
107 formation about major and minor compounds from such a big data set.

108 **Multivariate data exploration**

109 In this work, the main idea is to propose a method able to explore megapixel LIBS data set
110 without prior knowledge about the sample composition and to highlight simultaneously the presence
111 of major, minor compounds, and even traces. In the multivariate data analysis framework, this task
112 corresponds to the development of an unsupervised classification model. In other words, such tech-
113 niques try to find natural groupings of spectra in the considered data set, which will represent differ-
114 ent chemical compounds. Even if the chemometric community has developed different tools for un-
115 supervised classification of spectra, we can say without hesitation that the well-known k-means [12]
116 clustering (KM) is certainly the most popular one. Indeed, behind the apparent simplicity of this
117 method, it has been proved effective for many different kinds of data sets and spectroscopies. To the
118 best of our knowledge, as this algorithm has never been used in the framework of LIBS imaging, a
119 short description of the algorithm is provided below. Like any other chemometric method, a spec-
120 trum is considered as a point (denoted \mathbf{x}_i) in a multidimensional space. Let $X = \{\mathbf{x}_i, i = 1, \dots, n\}$ be a
121 dataset composed of n points (i.e. spectra) with $\mathbf{x}_i \in \mathbb{R}^w$, w being the number of spectral variables in
122 a spectrum. For illustrative purposes, let's consider a small LIBS imaging data set. This data cube of
123 size 5 pixels x 5 pixels x 2 wavelengths consists of 25 pseudo-spectra. Figure 1a illustrates the suc-
124 cessive steps of the k-means algorithm applied to this toy example. In a first step, k initial points
125 called centroids (in this example $k=3$) are randomly generated within the data domain (shown in col-
126 or in figure 1a). In the second step, one calculate distances between all points of the data set and the
127 generated centroids. In fact, the distance is used as a measure of similarity between spectra. In this
128 work, the cosine distance has been preferred to the Euclidean one, the latter being sensitive to global
129 intensity changes in spectra. However if the Euclidean distance had been selected, then it would
130 have been necessary to use a signal normalization commonly used in the LIBS community. The co-

131 sine distance $d_{i,j}$ between spectra \mathbf{x}_i and \mathbf{x}_j is given in equation 1 considering a point as a vector in a
132 multidimensional space:

$$133 \quad d_{i,j} = 1 - \frac{\vec{x}_i \cdot \vec{x}_j}{\|\vec{x}_i\| \cdot \|\vec{x}_j\|}$$

134 (1)

135 As we can see, this distance corrects for global intensity variations by dividing each spectrum i and j
136 by its norm. Given all the distances, each point (i.e. spectrum) is associated with the nearest centroid
137 and now belongs to one of the k classes. In a third step, the mean spectrum of each class is calculated
138 and will represent the k new centroids. In the fourth step, spectra in the data set are again unassigned.
139 Then steps 1–3 are repeated in a loop in order to refine the position of the k centroids. Calculations
140 are stopped when convergence is observed, i.e. when no further changes are observed in the spectra
141 class memberships. In the last step, the knowledge of the class membership of each spectrum and its
142 localization in the pixel space allow us to generate a clustering map using a color-coding. At the
143 same time, the centroid corresponding to each class is a spectrum used for chemical interpretation.

144 Behind the simplicity and ease of use of KM, there is an important issue which we have to
145 address, namely, how to select the optimal number of clusters or classes. Unfortunately far too often
146 in the literature, authors select with *a priori* this value of k , which is definitely the ultimate negative
147 choice. Indeed, no one can know the whole chemical complexity of the considered sample. In gen-
148 eral, the most reasonable way is to use a criteria called index in order to automatically choose this
149 value. This index is a mathematical function that measures the quality of a partition. The idea is then
150 to perform a KM clustering for different values of k ($2 \leq k \leq k_{max}$) and to calculate this index for each
151 partition. The highest index value indicates the optimal number of clusters for the considered data
152 set. One of the best index in the literature is PBM (Pakhira–Bandyopadhyay–Maulik) [13]. It is de-
153 fined as the square ratio between the largest normalized inter-cluster distance ER and the normalized
154 sum of intra-cluster distances RA:

$$155 \quad PBM(k) = \left(\frac{ER}{RA} \right)^2 \quad (2)$$

156 with $ER = \frac{\max_{l,m=1,\dots,k} \|c_l - c_m\|}{k}$, $RA = \frac{\sum_{i=1}^k \sum_{j=1}^{n_i} \|x_{i(j)} - c_i\|}{\sum_{i=1}^n \|x_i - \bar{x}\|}$, c_i the centroid of the i^{th} cluster ($i=1\dots k$),
157 $x_{i(j)}$ the j^{th} spectrum of the cluster i , n_i the total number of spectra in the cluster i , and \bar{x} the mean
158 point of the considered dataset. The PBM index will be used in this work in order to select the opti-
159 mal number of clusters.

160 We could obviously explore directly the proposed data set with KM in these conditions, but
161 we should not lose sight of our main goal, which is the simultaneous detection of major and minor
162 compounds. Indeed, this inquiry about the intrinsic data structure is very important because KM al-
163 gorithm (and most of the clustering methods) can fall into a trap under two specific conditions (fig-
164 ure 1b). The first problematic situation is observed when classes in the data set are unbalanced, that
165 is to say when a big difference in the number of spectra between classes is observed. This is precise-
166 ly the case for major and minor chemical compounds present in an imaging data set. As a conse-
167 quence, small populations of spectra would not be detected and wrongly associated with the nearest
168 big clusters. The second problematic situation arises when subpopulations of spectra are observed in
169 a given cluster. In this case, only a global cluster is generated and small spectroscopic details are lost
170 during this exploration. To address these issues, we have developed a new strategy, which we call
171 embedded k-means clustering (EKM). We were inspired by the way our brain works when we are
172 looking at a picture. We first extract the main features of the image (i.e. the main classes of objects)
173 and, then, we extract details about sub-zones of it. Thus, in the EKM strategy, a first k-means clus-
174 tering will be applied to the whole data set and the second round of clusterings will be applied to
175 each previously calculated cluster (figure 1c). Obviously, the PBM index will be used at each step of
176 the way.

177 All calculations in this work have been performed under the Matlab 2016b environment (The
178 Mathworks, Inc., Natick, Massachusetts) using homemade codes.

179 **RESULTS AND DISCUSSION**

180 To better understand the strengths of our data analysis strategy, it is essential to open this sec-
181 tion with the exploration of the considered imaging data set using the state-of-the-art method to gen-

182 erate chemical maps [4,14]. First, a single emission line is selected for an element of interest. Then a
183 baseline correction is applied on every single spectrum of the data set in order to extract correspond-
184 ing net intensities at the given wavelength. Lastly, color-coding is used in order to generate a colored
185 elemental map from these extracted values, the intensity of the chosen color being correlated with
186 abundance. Of course, this procedure can be successively repeated for all elements of interest in the
187 sample, with the possibility to observe them simultaneously in overlay mode on the same image.
188 Nevertheless, despite this operational simplicity, this traditional method imposes two constraints
189 which should be considered for the generation of unbiased chemical maps. First, each selected emis-
190 sion line should be the strongest one in the spectral domain for each element. But what is more im-
191 portant, a selected emission line should not present potential interferences with other lines. Due to
192 the natural complexity of the samples we usually explore, we quickly see that it is a strong hypothe-
193 sis, which, for each element of interest, could be difficult to hold in relation to the very high number
194 of lines in a spectrum. Figure 2 illustrates the use of this conventional approach to the rock section.
195 More specifically, figure 2a shows the mean spectrum calculated from all spectra of the imaging data
196 set. From this spectrum, it is always simple and fast to identify major elements by matching the ob-
197 served emission lines with an atomic spectra database. Thus it is easy to see, without being exhaus-
198 tive, the presence of different elements such as Pb, Ag, Fe, Ca, Mg, Mn, Cu, and Si. Figure 2b pre-
199 sents the global intensity image of the sample generated from the integration of the emission signal
200 for each pixel on the whole spectral domain. Of course, we are losing elemental information with
201 this observation but different zones of the samples can nevertheless be highlighted in this image. It is
202 even possible to observe different levels of homogeneity, textures, and sub-structures on the sample.
203 By contrast, figures 2c and 2d give elemental images generated with the conventional approach us-
204 ing single integrations described above. At first glance, we notice that many elements are localized
205 in specific areas. Although it is possible to observe the colocalization of element pairs such as Ag/Pb,
206 Si/Al, Si/Ti, and Zn/Cu for example, finding a correlation between all elements in this data set is a
207 hard task. Yet, we have to remember that such correlations should allow a trace-back to molecular
208 information i.e. mineral phases in this particular case. A further point concerns the detection of po-

209 tential anti-correlation between elements, which is especially difficult to achieve by just comparing
210 elemental images. It is indeed very interesting to know if a specific element is present in a zone
211 when another one is systematically absent or has a low concentration, and vice versa. In conclusion
212 of this section, while the usual procedure allows us to generate consistent elemental images most of
213 the time, we can clearly see that we are still not harnessing all the information contained in the data
214 set, minor compounds and minor phases not being particularly highlighted.

215 In this new section, the idea is to apply the strategy of embedded k-means clustering on the
216 considered data set and assess its interest for the simultaneous detection of major and minors com-
217 pounds. As explained previously, the initial step of this approach consists of the application of a first
218 k-mean clustering on the whole data set (i.e. all spectra). Figure 3a shows the evolution of the PBM
219 index according to the number of cluster k used in this first partitioning of pixels. Here it can be seen
220 clearly that an optimal number of five clusters has to be considered. Using this consideration as a
221 starting point, figure 3b provides a classification map from which we can observe the localization of
222 the five compounds. The percentage of pixels in a class for the total number of pixels in the data set
223 is also given. We can see, therefore, that classes 1,4 and 5 correspond to major compounds with
224 37%, 27% and 23% of pixels respectively. Nevertheless, at this point, we cannot say that classes 2
225 and 3 correspond to minor compounds with 3% and 10% of pixels respectively. In fact, they are only
226 somewhat less present. As regards the dispersion of compounds in the sample, classes 1 et 2 are
227 strictly observed in well-delimited and continuous areas. It is almost the case for class 5, which is
228 nevertheless also located around the area of class 1. More heterogeneous distributions are observed
229 for classes 3 and 4. Figure 3c gives the corresponding spectra of the centroids for each class. These
230 representative spectra are naturally used for chemical interpretation. Despite the fact that LIBS spec-
231 troscopy is an elemental one, the use of the whole spectral domain and some prior knowledge about
232 the genesis of rocks allow us to identify potential mineral phases. Thus, class 1 is associated with
233 galena (PbS) with traces of copper, silver, antimony, and tin. The mineral phase corresponding to
234 class 2 is calcite (CaCO_3) with traces of manganese, magnesium, silicon, and aluminum. Class 3 is
235 linked to quartz (SiO_2) with traces of magnesium, aluminum, calcium, titanium and iron. The next

236 mineral phase with class 4, is potentially an aluminosilicate ($\text{SiO}_2 / \text{Al}_2\text{O}_3$) or kinds of clays with
237 traces of magnesium, calcium, iron, manganese and titanium. Finally, class 5 is associated with an-
238 kerite ($\text{Ca}(\text{Fe}, \text{Mg}, \text{Mn})(\text{CO}_3)_2$) with traces of titanium.

239 To go deeper into the exploration of previous mineral phases, we shall apply the second step
240 of the embedded k-means strategy. Therefore, for each class, a new k-means clustering is applied
241 only to associated spectra. In other words, five k-means clustering are calculated in parallel consid-
242 ering the five different sub-populations of spectra contained in the five classes. Obviously, the PBM
243 index is used again to optimize the number of clusters of each k-means clustering. The five graphs
244 representing the evolution of the PBM index according to the number of clusters k are supplied in
245 the supplementary material (Figure S1). We then discover that all mineral phases exhibit sub-
246 populations of spectra. The galena (class 1) contains 3 sub-classes of compounds, the calcite (class
247 2) has 4, the quartz (class 3) has 5, the aluminosilicate phase (class 4) has 3 and ankerite (class 5)
248 has 6. Figure 4 gives classification maps for each phase and corresponding spectra of sub-classes.
249 For galena, classes 1.1 and 1.3 (in blue and red respectively) are the two major compounds of the
250 galena phase with 64% and 26% of pixels respectively. These two sub-classes exhibit different ratios
251 of elements such as Cu, Sb, Ag, and Sn. In this case, it is difficult to see any particular geographic
252 locations of the two. Class 1.2 (in yellow) constitutes the minor compound of the phase with 10% of
253 pixels for the total number of pixels in class 1. It takes the form of fine veins containing the highest
254 concentrations of Cu, Fe and Al compared to the two other sub-classes. For the calcite phase, classes
255 2.4 and 2.1 (respectively in blue and red) are the most abundant with 50% and 30% of pixels respec-
256 tively. They are distributed rather homogeneously and are very close in terms of element concentra-
257 tions except for Y and La. They form the purest calcites, Ca and Mn being their major elements. The
258 situation is very different for classes 2.3 and 2.2 (respectively in yellow and green), which are con-
259 centrated in small areas mainly at the borders of class 2. These minor compounds correspond to 14%
260 and 6% of pixels respectively. It is also remarkable that class 2.2 has the highest concentration of
261 Mg, Si, Fe, and Mn. Moreover, very small contributions of Y and La are now particularly detected in
262 the class 2.3, while being almost totally indetectable from the raw data set. The quartz phase is

263 slightly more complex with 5 sub-classes. However, a more balanced split can be observed between
264 the percentage of pixels of sub-classes. Classes 3.1 and 3.4 (respectively in red and blue) are the
265 most abundant. They are regularly distributed over a trapezoidal area such as class 3.3 (in yellow).
266 For its part, class 3.5 (in grey) is spread all over the class 3 area mostly in the form of tiny clusters.
267 This quartz is really particular because it has by far the highest concentration of Mg, Ca, Fe, Al, and
268 Ti. Class 3.2 (in green) is a minor compound with 9% of pixels. It is mainly observed along a vein
269 through the trapezoidal area. It contains less Si than the classes 3.1, 3.2, 3.3 and 3.4 but more Mg,
270 Ca, Fe, Al, and Ti. The aluminosilicate phase seems less complex with 3 sub-classes. However, from
271 a spectroscopic point of view, they are well-contrasted. Class 4.3 (in red) is the major compound
272 with 71% of pixels, followed by class 4.1 (in blue) with 25%. They are both spread all over the class
273 4 area. They show high concentrations of Si, Mg, Fe, and Al but also different ratios between them.
274 Class 4.2 (in yellow) is the minor compound of this phase with only 4% of pixels. It is spread all
275 over the area in the form of small clusters. At the same time, it has by far the highest Ti concentra-
276 tion and the lowest concentrations for all other elements. The fifth and last phase i.e. ankerite is cer-
277 tainly the most complex case with six sub-classes and the most contrasted element concentrations.
278 Classes 5.4 (in blue), 5.1 (in green) and 5.5 (in pink) are the most abundant with 34%, 33%, and
279 21% of pixels respectively. They are distributed rather homogeneously with rather high concentra-
280 tions of Mg, Ca, and Fe. The last three sub-classes are minor compounds. Class 5.6 (in grey) with
281 7% of pixels is mainly located at the border of the rock section. It has medium concentrations of Ca
282 and Si, a medium one for Mg and contains neither Fe nor Zn. Class 5.3 (in yellow) with 4% of pixels
283 is only located on one side of the area defined by classes 5.1, 5.4, and 5.5. It has also concentrations
284 of Fe, Mg, Ca and Si comparable to those three previous classes. However, small variations of con-
285 centration ratios are observed between them. For its part, class 5.2 (in red) is the less abundant com-
286 pound with 0.2% of pixels. It is presented in the form of a single cluster. It is the only compound
287 containing Zn and a small concentration of Fe. The other elements are absent. Readers interested in a
288 global representation of the 21 sub-classes in overlay mode should refer to figure S2 in the supple-
289 mentary material. As we have just seen, our strategy allows us to deeply explore LIBS data sets of

290 complex samples providing simultaneously the localization and the identification of major and minor
291 compounds. Class 5.2 is certainly the perfect example of the potential of this approach because it
292 corresponds to the detection of only 730 specific spectra of a given compound over the 2.289.000
293 present in the considered data set. In a natural way, the PBM index was also used on each cluster of
294 the second levels of clustering demonstrating that there was no more possible discrimination at this
295 level thus ending the exploration of this megapixel LIBS imaging data set.

296 **CONCLUSION**

297 The main objective of this work was to evaluate an original strategy called embedded k-
298 means clustering in order to explore a big LIBS imaging data set acquired from a complex mineral
299 sample. More specifically, the idea was to propose a simultaneous identification and localization of
300 both major and minor compounds. From the very start of this work, we have quickly observed that
301 while the traditional signal integration method generates unbiased elemental images most of the
302 time, it remains especially tricky if the objective is to obtain information at the phase level, for the
303 highest as well as the lowest concentrations. Generally speaking, we have demonstrated that multi-
304 variate data analysis is an efficient complementary tool to explore LIBS imaging data sets in this
305 particular framework. Indeed, the k-means algorithm has allowed us to group similar pixels (i.e.
306 spectra) without any prior knowledge of class memberships. We have also highlighted the im-
307 portance of using an index in order to select the right number of clusters, with no *a priori* about the
308 considered sample, which to our knowledge has never been done in the LIBS framework. Lastly, we
309 have shown that our approach based on successive k-means clustering provides a deeper exploration
310 of the sample from major to minor compounds with great sensitivity, without compromise on the
311 detection of both.

312 **AUTHOR INFORMATION**

313 **ORCID**

314 Ludovic Duponchel: 0000-0002-7206-4498

315 **Cécile Fabre: 0000-0001-8627-4050**

316 **Jean Cauzid: 0000-0001-5587-9874**

317 **Frédéric Pelascini: 0000-0002-3779-6685**

318 **Vincent Motto-Ros: 0000-0002-6955-3098**

319 **Author Contributions**

320 The manuscript was written through the contributions of all authors. Moreover, they have
321 given approval to the final version of the manuscript.

322 **ACKNOWLEDGMENT**

323 This work was partially supported by Pulsalys (#L0978-L1294), the French region Grand
324 Est, and the French region Rhône Alpes Auvergne (Optolyse, CPER2016). In addition, we grateful-
325 ly acknowledge Tristan Mantoy and the CMT company (Compagnie Minière de Touissit) for
326 providing the sample. Lastly, we would like to thank Ms Astrid Marissal for proofreading.

327

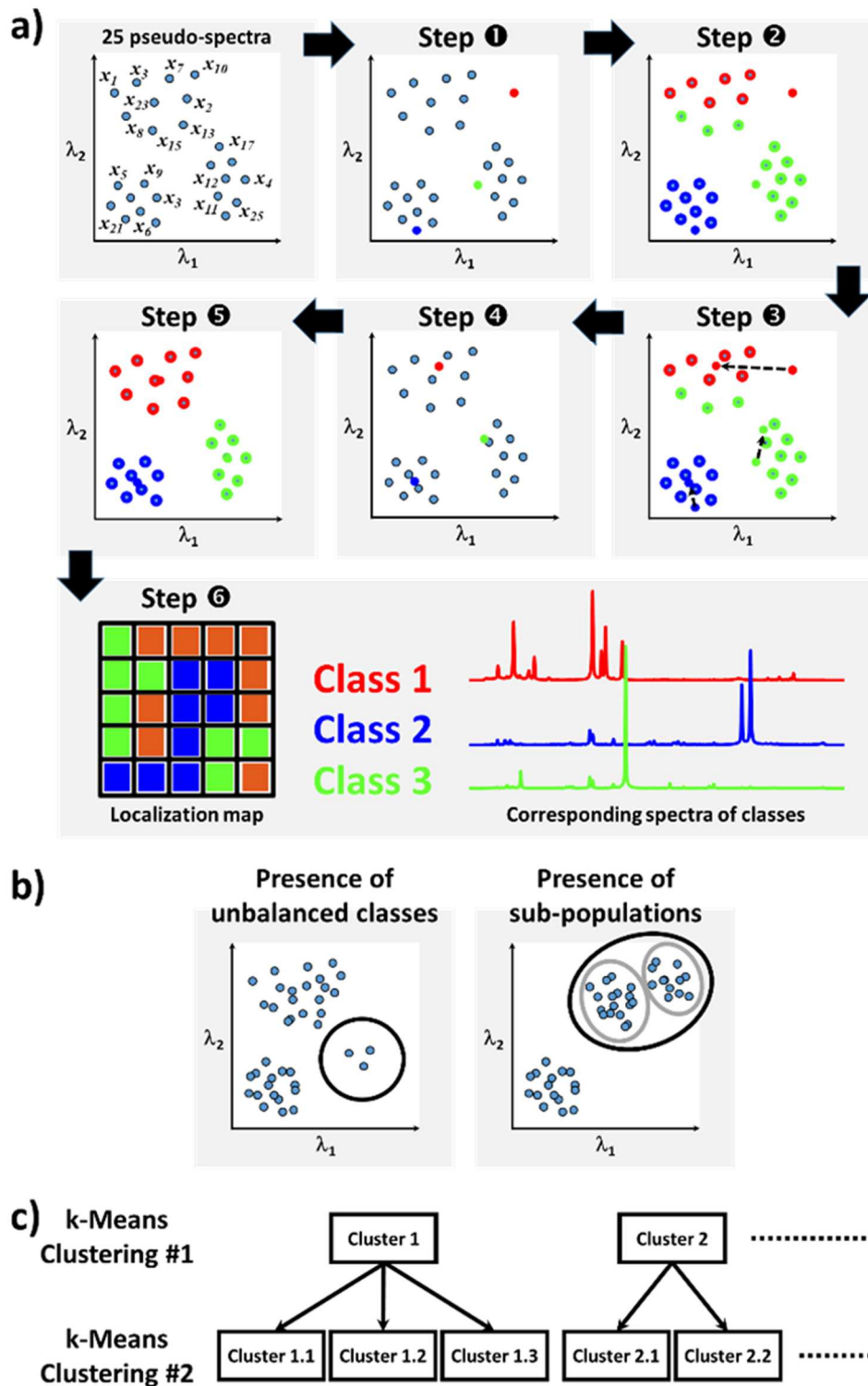
328 **REFERENCES**

- 329 [1] L. Jolivet, M. Leprince, S. Moncayo, L. Sorbier, C.-P. Lienemann, V. Motto-Ros, Review of
330 the recent advances and applications of LIBS-based imaging, *Spectrochim. Acta Part B At.*
331 *Spectrosc.* 151 (2019) 41–53. <https://doi.org/10.1016/j.sab.2018.11.008>.
- 332 [2] R. Gaudiuso, N. Melikechi, Z.A. Abdel-Salam, M.A. Harith, V. Palleschi, V. Motto-Ros, B.
333 Busser, Laser-induced breakdown spectroscopy for human and animal health: A review, *Spec-*
334 *trochim. Acta Part B At. Spectrosc.* 152 (2019) 123–148.
335 <https://doi.org/10.1016/j.sab.2018.11.006>.
- 336 [3] C. Fabre, D. Devismes, S. Moncayo, F. Pelascini, F. Trichard, A. Lecomte, B. Bousquet, J.
337 Cauzid, V. Motto-Ros, Elemental imaging by laser-induced breakdown spectroscopy for the
338 geological characterization of minerals, *J. Anal. At. Spectrom.* 33 (2018) 1345–1353.
339 <https://doi.org/10.1039/C8JA00048D>.
- 340 [4] J.O. Cáceres, F. Pelascini, V. Motto-Ros, S. Moncayo, F. Trichard, G. Panczer, A. Marín-
341 Roldán, J.A. Cruz, I. Coronado, J. Martín-Chivelet, Megapixel multi-elemental imaging by La-

- 342 ser-Induced Breakdown Spectroscopy, a technology with considerable potential for paleocli-
343 mate studies, *Sci. Rep.* 7 (2017) 1–11. <https://doi.org/10.1038/s41598-017-05437-3>.
- 344 [5] F. Trichard, F. Gaulier, J. Barbier, D. Espinat, B. Guichard, C.-P. Lienemann, L. Sorbier, P.
345 Levitz, V. Motto-Ros, Imaging of alumina supports by laser-induced breakdown spectroscopy:
346 A new tool to understand the diffusion of trace metal impurities, *J. Catal.* 363 (2018) 183–190.
347 <https://doi.org/10.1016/j.jcat.2018.04.013>.
- 348 [6] J. El Haddad, L. Canioni, B. Bousquet, Good practices in LIBS analysis: Review and advices,
349 *Spectrochim. Acta Part B At. Spectrosc.* 101 (2014) 171–182.
350 <https://doi.org/10.1016/j.sab.2014.08.039>.
- 351 [7] J.-B. Sirven, B. Bousquet, L. Canioni, L. Sarger, Laser-Induced Breakdown Spectroscopy of
352 Composite Samples: Comparison of Advanced Chemometrics Methods, *Anal. Chem.* 78
353 (2006) 1462–1469. <https://doi.org/10.1021/ac051721p>.
- 354 [8] I. Gaona, J. Serrano, J. Moros, J.J. Laserna, Range-Adaptive Standoff Recognition of Explosive
355 Fingerprints on Solid Surfaces using a Supervised Learning Method and Laser-Induced Break-
356 down Spectroscopy, *Anal. Chem.* 86 (2014) 5045–5052. <https://doi.org/10.1021/ac500694j>.
- 357 [9] N.C. Dingari, I. Barman, A.K. Myakalwar, S.P. Tewari, M. Kumar Gundawar, Incorporation of
358 Support Vector Machines in the LIBS Toolbox for Sensitive and Robust Classification Amidst
359 Unexpected Sample and System Variability, *Anal. Chem.* 84 (2012) 2686–2694.
360 <https://doi.org/10.1021/ac202755e>.
- 361 [10] T. Zhang, H. Tang, H. Li, Chemometrics in laser-induced breakdown spectroscopy, *J.*
362 *Chemom.* 32 (2018) e2983. <https://doi.org/10.1002/cem.2983>.
- 363 [11] M. Bouabdellah, J.F. Slack, *Mineral Deposits of North Africa*, Springer, 2016.
- 364 [12] J. MacQueen, Some methods for classification and analysis of multivariate observations, in:
365 *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Vol-*
366 *ume 1: Statistics*, 1967.
- 367 [13] M.K. Pakhira, S. Bandyopadhyay, U. Maulik, Validity index for crisp and fuzzy clusters, *Pat-*
368 *tern Recognit.* 37 (2004) 487–501. <https://doi.org/10.1016/j.patcog.2003.06.005>.

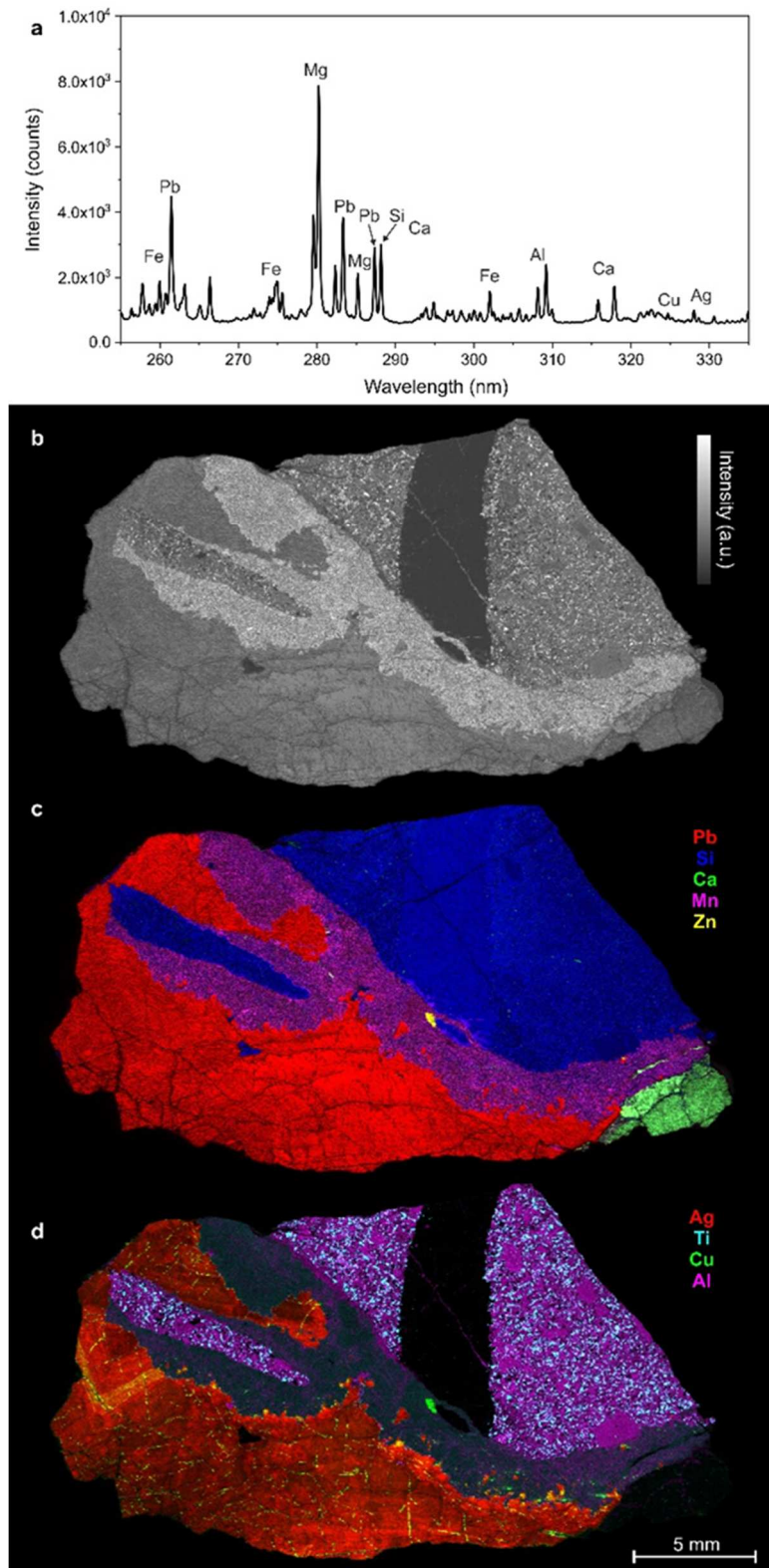
369 [14] Y. Gimenez, B. Busser, F. Trichard, A. Kulesza, J.M. Laurent, V. Zaun, F. Lux, J.M. Benoit,
370 G. Panczer, P. Dugourd, O. Tillement, F. Pelascini, L. Sancey, V. Motto-Ros, 3D Imaging of
371 Nanoparticle Distribution in Biological Tissue by Laser-Induced Breakdown Spectroscopy, *Sci.*
372 *Rep.* 6 (2016) 29936. <https://doi.org/10.1038/srep29936>.

373



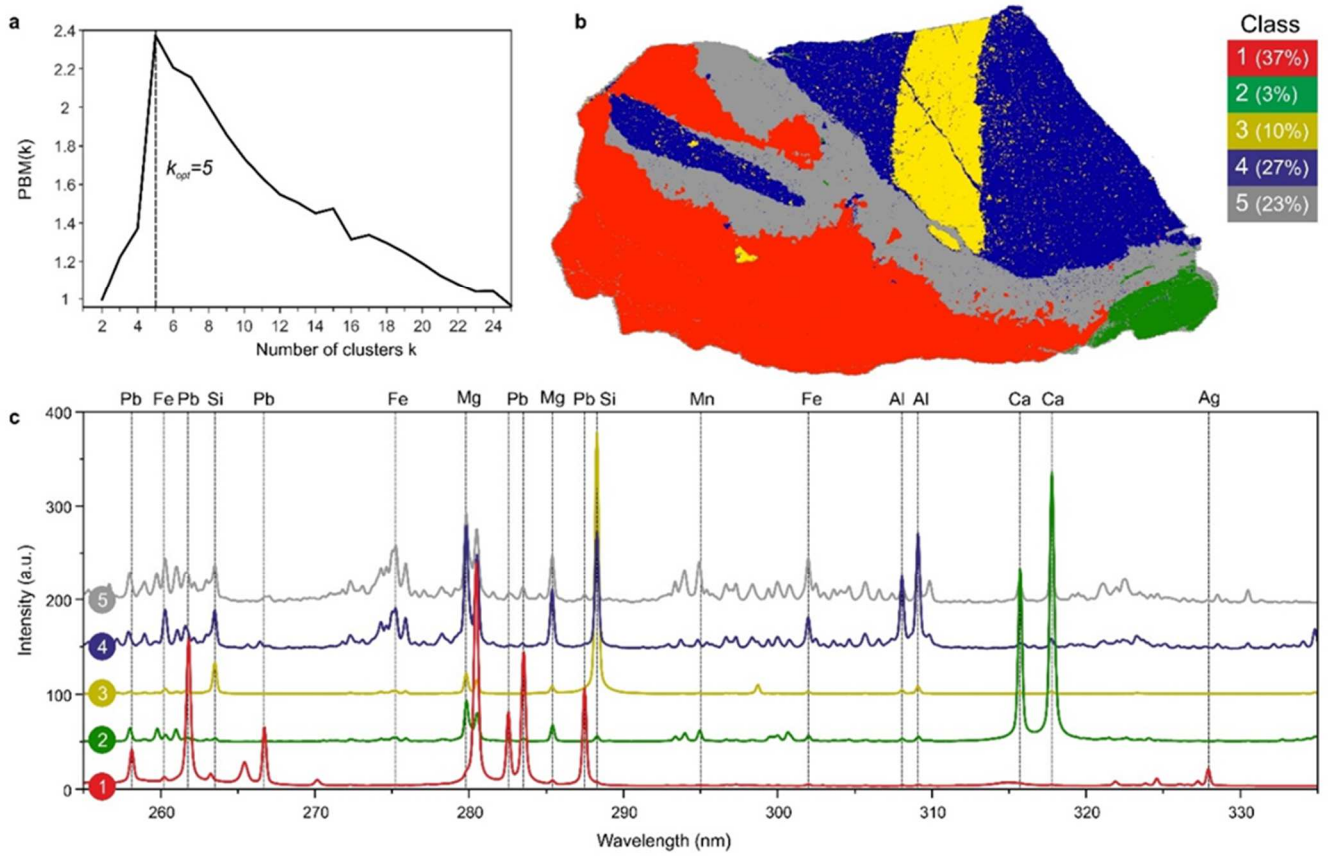
377 **Figure 1.** a) The k-means algorithm applied to spectroscopic imaging. b) Problematic data structures

378 hardy managed by k-means. c) The proposed method called embedded k-means clustering.



379

380 **Figure 2.** a) The mean spectrum of the LIBS data set. b) The global intensity image. c) and d) Ele-
 381 mental images generated with the conventional approach. A high-resolution version of this image
 382 can be downloaded from supplementary materials.

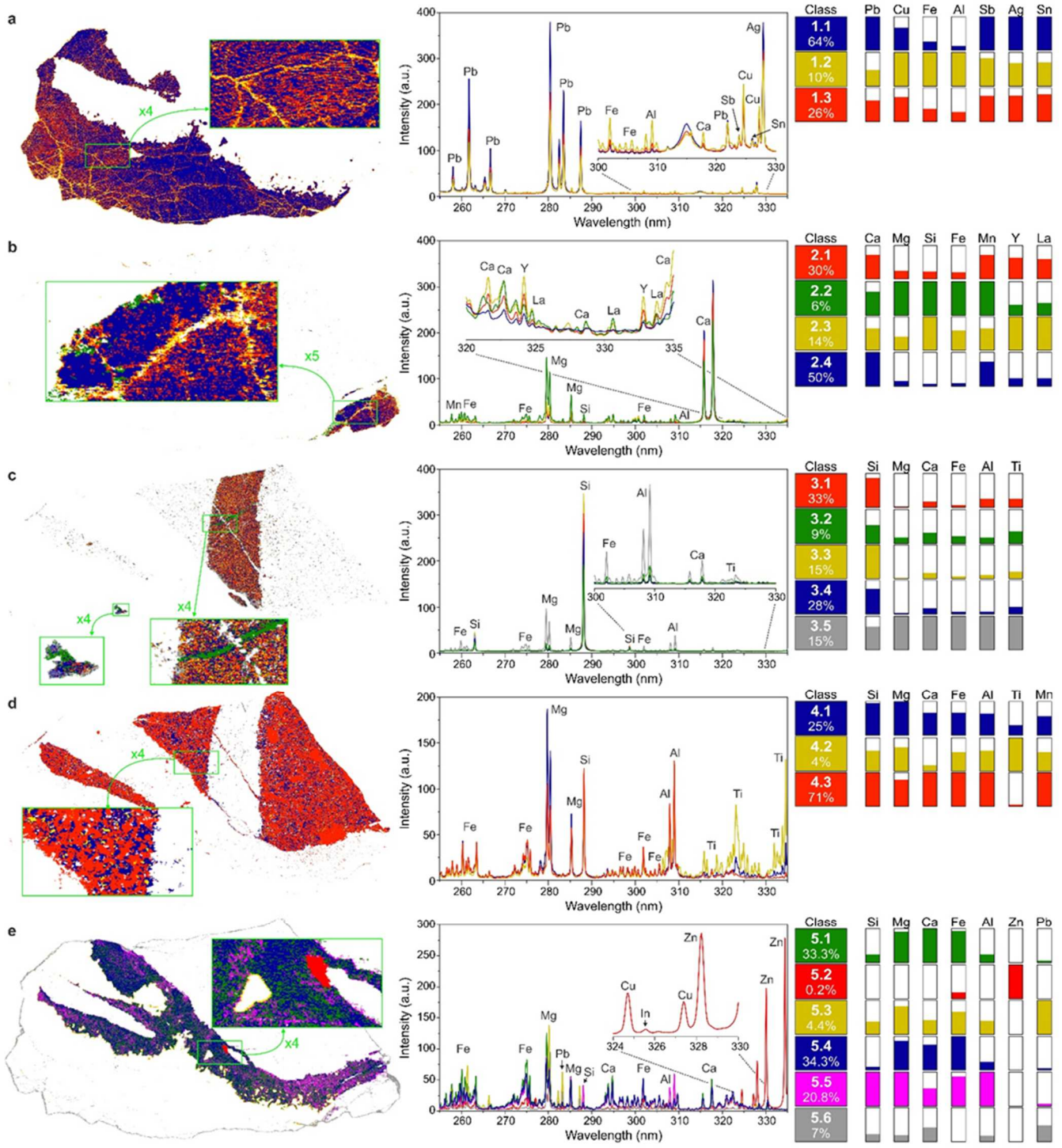


383

384 **Figure 3.** a) Evolution of the PBM index according to the number of clusters k . b) The classification

385 map considering an optimal number of clusters equal to 5. c) Representative spectra of each class.

386



387

388 **Figure 4.** Classification maps obtained for each phase with corresponding spectra of sub-populations

389 and relative concentrations of elements.

390

391

392

**Deep diving
in**

**2 million
of spectra!**

