



Buse, B., & Kearns, S. (2018). Evaluating X-Ray Microanalysis Phase Maps Using Principal Component Analysis. *Microscopy and Microanalysis*, 24(2), 116-125. <https://doi.org/10.1017/S1431927618000090>

Peer reviewed version

Link to published version (if available):
[10.1017/S1431927618000090](https://doi.org/10.1017/S1431927618000090)

[Link to publication record in Explore Bristol Research](#)
PDF-document

This is the author accepted manuscript (AAM). The final published version (version of record) is available online via Cambridge University Press at <https://www.cambridge.org/core/journals/microscopy-and-microanalysis/article/evaluating-xray-microanalysis-phase-maps-using-principal-component-analysis/E061505548B36F36CA25A18D643C4CC4> . Please refer to any applicable terms of use of the publisher.

University of Bristol - Explore Bristol Research

General rights

This document is made available in accordance with publisher policies. Please cite only the published version using the reference above. Full terms of use are available:
<http://www.bristol.ac.uk/pure/about/ebr-terms>

1 Evaluating x-ray microanalysis phase maps using principal component 2 analysis.

3 Ben Buse and Stuart Kearns

4 Abstract

5 Automated phase maps are an important tool for characterising samples but data quality must be
6 evaluated. Common options include overlay phases on BSE images and phase composition averages
7 and standard deviations. Both these methods have major limitations. We propose two methods of
8 evaluation involving principal component analysis. First, a RGB composite image of the first three
9 principal components, which comprise the majority of chemical variation, and provides a good
10 reference against which phase maps can be compared. Advantages over a BSE image include
11 discriminating between similar mean atomic number phases and sensitivity across the entire range
12 of mean atomic numbers present in a sample. Second, principal component maps for identified
13 phases, to examine for chemical variation within phases. This ensures the identification of
14 unclassified phases and provides the analyst with information regarding the chemical heterogeneity
15 of phases (e.g. chemical zoning within a mineral, or mineral chemistry changing across an alteration
16 zone). Spatial information permits a good understanding of heterogeneity within a phase and allows
17 analytical artefacts to be easily identified. These methods of evaluation were tested on a complex
18 geological sample. K-means clustering and K-nearest neighbour algorithms were used for phase
19 classification, with the evaluation methods demonstrating their limitations.

20 Key words

21 Electron probe microanalysis, scanning electron microscopy, phase mapping, k-means clustering,
22 wavelength dispersive spectroscopy (WDS), energy dispersive spectroscopy (EDS), elemental maps,
23 clustering

24 Introduction

25 Automated phase mapping is a widely available tool within software suites for energy dispersive
26 spectrometers (EDS) on scanning electron microscopes (SEM) and is now available for electron
27 probe microanalysis (EPMA) using wavelength dispersive spectrometers (WDS). Phase mapping uses
28 multi-dimensional data (spatially defined element intensities or concentrations) and identifies
29 chemically distinct phases to give their spatial distribution. The automated algorithms included in
30 the instrument software make this process straightforward for the operator (e.g. a form of principal
31 component analysis using rotation (Kotula et. al 2003), clustering algorithm in Oxford Instruments
32 AutoPhaseMap software (Statham et al. 2013), K-means clustering in Probe for Epma software
33 (www.probesoftware.com) and hierarchical cluster analysis in JEOL EPMA software (Mori et al.
34 2017)).

35 It is difficult to assess the quality of a phase classification, Munch et al. (2015) demonstrate some of
36 the limitations of phase algorithms and suggest the need for expert review. Common options to
37 evaluate phase classifications include overlay on BSE images and phase composition averages and
38 standard deviations. Liebske (2015) provides an improvement in the open source package iSpectra
39 which allows overlays on principal component maps or RGB elemental maps. Phase composition
40 averages are often difficult to interpret being affected by convoluted pixels at grain boundaries,
41 where the limits of analytical resolution results in measured intensities consisting of a convolution of
42 adjacent phases. van Hoek et al. (2011) and Liebske (2015) showed how these 'bad' pixels can be
43 eroded to give phase composition averages reflecting true compositions but this processing is not

44 available in most phase mapping packages. Algorithms can be independently verified for reference
45 samples against methods such as manual thresholding (Maloy & Treiman 2007) or EBSD phase maps
46 (Statham et al. 2013), but this does not ensure the algorithm works correctly for all samples and
47 operating conditions (Munch et al. 2015)

48 In this study k-means clustering and k-nearest neighbour (KNN) algorithms are used to demonstrate
49 some of the problems and show how, irrespective of the phase mapping algorithm, principal
50 component analysis (PCA) can be used to assess the quality of phase classification. PCA assigns new
51 dimensions which capture the variability of the dataset; the first dimension corresponds to
52 maximum variance; each subsequent dimension is orthogonal to the previous and captures the
53 maximum remaining variance (Tan et al. 2006). The merit of this technique is it provides an unbiased
54 method of reducing multiple dimensional systems (e.g. 10 chemical elements) to a small number of
55 dimensions which can be visualised graphically. PCA and various refinements are commonly used in
56 the generation of phase maps (e.g. Kotula et al. 2003, Parish 2011), here we demonstrate their
57 strength in evaluation of phase maps.

58 This study uses quantitative maps for the phase classification. Quantitative maps provide significant
59 advantages over raw count maps, allowing the interrogation of phase data and importantly average
60 phase compositions to be extracted.

61 [Methods and materials](#)

62 A complex sample was selected with multiple minerals of varying abundance and finely intergrown
63 minerals at or below the limits of analytical resolution (controlled by accelerating voltage and
64 pixel/step size) (see Figure 1c). The sample is a metamorphosed basalt xenolith within a kimberlite
65 (for details see Buse et al. 2010). The area mapped extends from the kimberlite into the basalt
66 xenolith (Figure 1a) with alteration of the basalt most intense adjacent to the kimberlite.
67 Quantitative element maps were collected using 5 WDS on a JEOL 8530F EPMA at the University of
68 Bristol. Elements collected were Si, Na, Ca, Fe and Ti in the first pass and Mg, Al, K, Mn in the second
69 pass. The operating conditions were 20 kV accelerating voltage, 40 nA beam current, 10 millisecond
70 dwell time and a 5 μm step size. The quantitative element maps are combined into a single array (X
71 coordinate, Y coordinate, Si, Na, Ca, Fe, Ti, Mg, Al, K, Mn) for processing.

72 The phase maps were generated and evaluated using R (General Public License software for
73 statistical computing), which includes k-means clustering, k-nearest neighbour and PCA packages.

74 K-means clustering requires the initial cluster centres to be specified or determined randomly for a
75 given number of clusters. Data is classified through a series of iterative loops in which all the points
76 (pixels) are assigned to the nearest cluster centre (centroid) and the centroid position is updated.
77 Here K-means clustering was run in three variants: (1) using randomly assigned initial cluster centres
78 for 15 clusters; (2) using specified cluster centres for discrete phases identified from the Red-Green-
79 Blue (RGB) composite image of principal components. 9 discrete phases were identified (Figure 1b
80 and Table 1). For each discrete phase a single area composition was extracted from the element
81 maps; (3) using maximum element intensities were used as the initial cluster centres.

82 K-nearest neighbour requires a reference dataset against which the pixel compositions are checked.
83 The reference dataset consisted of the compositions of the 9 discrete phases (Figure 1b and Table 1)
84 selected for K-means specified cluster centres. A normal distribution using the measured standard
85 deviation was applied to each composition to present a range of compositions for each phase. For
86 each pixel the 10 closest reference values in chemical space were examined, a pixel was assigned to
87 a phase if at least 7 of the reference values belonged to the same phase, otherwise it was rejected.

88 For PCA, the dataset was centred in principle component space so that the mean of each principle
89 component is 0 rather than the mean of the compositional data. This ensures that the first principle
90 component is not dominated by the position of the dataset with respect to the origin (Jolliffe 2002).
91 The dataset was not scaled (a covariance matrix was used); scaling (a correlation matrix) gives equal
92 weight to the variance of each component (here chemical element) and is important where
93 components of different units (e.g. length, weight etc.) are being analysed (Jolliffe 2002). It is not
94 desirable here, because the units of each of component (wt.%) are the same. By not scaling, the
95 variance is dominated by major element variance rather than giving trace elements equal weight.
96 This is desirable because the phase separation is based on major constituents and noise dominates
97 the trace element signals (van den Berg et al. 2006). PCA was conducted on the entire dataset of
98 multiple element intensities for the whole map to produce RGB composite images of the first three
99 principal components. PCA was also conducted on subsets of the data, consisting of the element
100 intensity pixels for a single phase to produce phase-PC (principal component) maps and scatter
101 graphs.

102 Results

103 Red-Green-Blue Principal Component (RGB-PC) images

104 Figure 2 compares several phase mapping algorithms to a backscattered electron (BSE) image and an
105 RGB composite image consisting of the first three principal components (PC), derived from principal
106 component analysis. The RGB-PC image provides a good tool to assess phase maps; similar to the
107 use of BSE images in some phase analysis software (e.g. Thermo Scientific NSS permits overlays of
108 phases onto a BSE image). Figure 2b shows that the RGB-PC image provides a greater phase
109 separation than the false-colour BSE image (Figure 2e). The BSE image has difficulty both separating
110 phases with similar mean atomic number (e.g. on Fig 2e colours of bultfonteinite [yellow-green],
111 clinopyroxence [green-blue], and serpentine [blue] overlap; see Table 1 for mineral compositions)
112 and covering the range of mean atomic numbers present in the image. The high atomic number
113 phases ilmenite, perovskite and barite cannot be separated (red on Fig 1e) with the detector
114 brightness and contrast set for sensitivity at the low mean atomic numbers. BSE images are more
115 sensitive to topography than most x-ray intensities. The RGB-PC image is derived from the same
116 dataset (x-ray intensities) as the phase maps. BSE images can still make a contribution in checking
117 phase maps; dependant on mean atomic number they can identify variations not measured by x-ray
118 intensities (e.g. H₂O in normalised EDS data; Munch et al. 2015) – in the case of figure 2 only the BSE
119 image differentiates between barite and holes – with sulphur and barium not measured. Whilst the
120 phase mapping system, using only WDS data for the measured elements, cannot be expected to
121 differentiate between the two, a comparison with the BSE image alerts the analyst. The RGB-PC
122 image extends the evaluation tools suggested by Liebske (2015) and accounts for most of the
123 chemical variation within the sample.

124 The RGB-PC image provides a reference for a visual assessment of the phase maps. The number of
125 phases and textural features (e.g. shape of grains) can be checked. In the examples given, K-means
126 with 15 random clusters (Fig. 2a) subdivided hydrogarnet and serpentine into numerous phases (on
127 Fig. 2a hydrogarnet is orange, dark-red and black, and serpentine is blue, pink and violet; see also
128 Table 2). The other phase maps (Fig. 2c,d & f) provide a close match to the RGB-PC image. In
129 addition K-means with 15 random clusters (Fig. 2a) identifies ilmenite and perovskite as a single
130 “oxide” phase, which from RGB-PC image can be seen to consist of chemically distinct phases. The
131 other phase maps correctly separate this “oxide” phase into ilmenite and perovskite. This distinction
132 is critical for interpreting the sample, with ilmenite absence from the margin of the basalt xenolith as
133 a result of alteration penetrating into the basalt from the kimberlite (Buse et al. 2010).

134 Serpentine and bultfonteinite form fine intergrowths below analytical resolution; on the RGB-PC
135 image (Fig. 2b & h) this intergrowth form a distinct phase (dark purple distinct from the bright pink
136 of serpentine). This phase is well characterised using the KNN algorithm, which on Fig 2i in
137 comparison with the RGB-PC image can be seen to faithfully reproduce the serpentine, serpentine-
138 bultfonteinite intergrowth, and hydrogarnet phases. Again K-means with 15 random clusters can be
139 seen to split the phases into many subdivisions (serpentine-bultfonteinite intergrowth is split into
140 dark brown and grey phases). The K-means, using specified clusters or maximum intensity, struggles
141 in the classification of serpentine and serpentine-bultfonteinite intergrowth. Serpentine and
142 serpentine-bultfonteinite intergrowth are under-represented, whilst a mixed serpentine phase (pink
143 on Fig 2l & j) and bultfonteinite are over represented (see black arrows on Fig2l in comparison with
144 Fig2i and Fig 2h).

145 The main distinction between KNN and K-means using specified clusters or maximum intensity, is
146 that in the latter the cluster centre can shift during the iterative process. The result of this is shown
147 in comparison with Table 1 where phase 2 has shifted and now represents an additional mixed
148 serpentine phase, which diminishes the serpentine-bultfonteinite intergrowth phase (pink phase;
149 Fig2l & j) and misclassifies convoluted boundary pixels. The latter is seen in Figure 3c where the
150 black arrow identifies pink boundary pixels, a convolution of serpentine and hydrogarnet, not visible
151 on either the RGB-PC image (Fig 3b) nor the KNN phase classification (Fig 3a). Another example of
152 iterative shifting of the cluster centre is phase 9 in Table 1 (lilac phase on phase maps) which has
153 shifted so that the phase includes both the initial apatite (see Fig 3c), mixed phases (see Fig 3c & 2l
154 circles) and pits and barite (Fig 2 d & f).

155

156 Phase analysis software commonly report phase composition averages and standard deviations.
157 Table 2 gives the values for K-means with 15 random clusters. The values are difficult to interpret as
158 averages may differ significantly from the true composition. Table 3 gives the composition of
159 clinopyroxene extracted from several pixels within a single clinopyroxene crystal, here the
160 composition closely matches stoichiometry. Poor phase composition averages are often the product
161 of analytical resolution (see van Hoek et al. 2011, Liebske 2015) as pixels at the margins of grains can
162 have convoluted x-ray intensities of multiple phases. The phase may also include bad pixels where
163 topography gives poor results again skewing the phase average. Table 2 also shows the problems of
164 identifying mixed phases resulting from small grains dominated by convoluted pixels of boundaries
165 or finely intergrown phases. To correctly identify phases, comparison with a RGB-PC image can be of
166 considerable help.

167 Phase Principal Component (phase-PC) maps

168 Phase-PC maps provide a useful tool to assess the homogeneity of each phase and determine
169 whether it includes multiple phases which the algorithm has failed to discriminate. Figure 4a shows a
170 phase-PC map of the "oxide" phase generated from the k-means algorithm using 15 random
171 clusters. The phase-PC map clearly distinguishes ilmenite (orange) from perovskite (purple). The
172 spatial separation and the magnitude of variance provides strong evidence for two distinct phases.

173 Phase-PC maps for perovskite and ilmenite, which the KNN algorithm correctly identifies as two
174 separate phases, shows the perovskite to be relatively homogenous whilst ilmenite contains two
175 spatially and chemically distinct phases. Both the perovskite and ilmenite phases include some
176 chemical variation from convolution at the margins of grains. Using PC-1 ilmenite can be subdivided
177 into two compositions (Table 4). The small purple grains in the kimberlite (identified on Fig. 4c)
178 represent Fe-Ti-Mg spinel are distinct from the ilmenite within the basalt xenolith.

179 The PC-1 Phase-PC map shows this distinction less clearly for ilmenite identified using k-means with
180 specified clusters. The ilmenite data contains more scatter than observed for the KNN ilmenite
181 phase. This is consistent with k-means not rejecting any 'bad' pixels, unlike KNN. For the low
182 abundance phase this scatter has significant influence and results in the rotation of the principal
183 components (see Figures 4g and 4h). In this case the PC-2 Phase-PC map most clearly distinguishes
184 ilmenite from Fe-Ti-Mg Spinel, whereas PC-1 shows variations within grains suggestive of
185 distinguishing convoluted pixels.

186 Figure 5 shows Phase-PC maps for clinopyroxene, serpentine and bultfonteinite, all of which are
187 relatively homogenous. A comparison with the PC scatter graphs illustrates the benefit of spatial
188 information. Both clinopyroxene and serpentine show small PC variations. On the maps it is evident
189 that for clinopyroxene this is uniformly distributed whereas for serpentine it varies between the
190 kimberlite and the xenolith. Variation within the clinopyroxene probably relates to pixel convolution
191 although could relate to chemical zonation within the clinopyroxene. Variations within the
192 serpentine suggest the chemistry of the serpentine differs spatially. The average compositions (Table
193 5) are inaccurate and difficult to interpret possibly due to convoluted pixels as discussed above.
194 Beam damage may also add to the reduced data quality. However, variations in Al, Si, Mg and Fe,
195 with Al enriched in the kimberlite are clearly apparent. The presence of spatial variations provides
196 important information about the sample, which should prompt further detailed investigations to
197 understand the cause. Element maps extracted for the serpentine phase (Figure 6) confirm the
198 variations suggested by the average compositions (Table 5) with Al substituting for Si and Mg for Fe.

199 Figures 5 c-d compare bultfonteinite from K-means using 15 clusters and from K-means specified
200 clusters. The difference can be explained as K-means specified clusters has fewer clusters resulting in
201 the bultfonteinite phase being less tightly constrained and containing marginal data. This
202 incorporation of marginal data is shown in Fig. 5d where the centre of the grains (purple)
203 corresponds closely to Fig 5c, whereas the rest of the data (orange) consists of increasingly mixed
204 compositions excluded from the more tightly constrained cluster of K-means using 15 clusters.

205 Discussion

206 Phase classification is complex and subject to the limitations of the algorithm used. PCA provides a
207 method of evaluating the quality of a given phase classification method. PCA is used in reference to
208 phase maps and element maps: checking that the phase maps represent the variation identified in
209 the RGB-PC image and checking for any variation within an individual phase. In the latter case,
210 variation within a phase is explained by the elemental data extracted for the discrete variations in PC
211 identified (e.g. Table 4 where extracted compositions allowed Fe-Ti-Mg spinel to be identified within
212 the ilmenite phase). This use of PCA in reference to phase maps and element maps avoids the
213 difficulties associated with interpreting principal components from their component weights (see
214 Kotula et al. 2003). PCA requires an orthogonal arrangement of components which may not
215 correspond to data variation (Kotula et al. 2003) as shown in the phase-PC maps and scatter graphs
216 in Figure 4 d, e and h. This problem is mitigated in the case of RGB-PC images for it is a composite of
217 3 principal components. In some cases the orthogonal requirement can obscure variations in phase-
218 PC maps suggesting in these cases phase-PC maps for each principal component are required, or
219 possibly scatter graphs or RGB-PC images for the phase. Improvements might be possible through
220 rotating PC or by removing the orthogonal constraint. Regardless the phase-PC maps show the value
221 of this or similar techniques in identifying variations in multi-dimensional space within individual
222 phases and displaying their spatial component.

223 The data presented shows how PCA can be used to identify incorrect phase classifications; here
224 exposing the limitations of K-means clustering using random clusters. K-means works best for phases
225 of similar abundance, which form spherical clusters in chemical space and where the initial allocated
226 centres reflect phase distribution (Tan et al. 2006). Both the RGB-PC image and the phase PC maps
227 identify the “oxide” phase and the phase PC maps show the “oxide” phase to actually consist of
228 perovskite and ilmenite. Due to their low abundance, these phases are not distinguished using
229 randomly allocated cluster centres, which only subdivides more abundant phases (Fig. 7, see also
230 Munch et al 2015).

231 Specifying initial cluster centers, either by identifying phases beforehand (K-means specified clusters,
232 Fig 2d, see also Munch et al. 2015) or by using maximum element intensities (K-means max element,
233 Fig 2f), to a large extent overcomes these limitations by ensuring the initial cluster centers represent
234 phase distribution. The use of maximum element intensities does not require prior knowledge of
235 phases but requires the number of phases to match the number of elements and for phases to be
236 discriminated to a large extent by a particular element. A variant on this is using the KNN algorithm
237 which does not iteratively shift cluster centres and allows pixels to be rejected (not classified). KNN
238 gives consistent results for spatially distinct regions of the same rock sample; the absence of
239 iteration means it is largely unaffected by the absence of a phase within an individual map. The
240 danger with these methods is, in the case of specifying phases, not all the phases present in the
241 sample may have been identified, and in the case of maximum element intensities, there may be
242 more phases than elements. In these cases, the phase PC maps work well at identifying aggregate
243 phases in which discrete phases have been classified together – as shown in the case of the “oxide”
244 phase, and also the Fe-Ti-Mg spinel phase which was identified subsequent to phase analysis.

245 An alternative approach to overcoming the tendency for k-means to subdivide high abundance
246 phases before distinguishing low abundance phases, is a set of criteria which recombine phases if
247 certain thresholds are exceeded (Munch et al. 2015, Statham et al. 2013). With this approach, similar
248 to using maximum element intensity, prior knowledge of phases is not required. However, it is still
249 important to evaluate the output classification (Munch et al. 2015).

250 Convolved pixels cause many problems in phase analysis and algorithms must ensure they are not
251 assigned to distinct “boundary” phases. To identify true phase compositions these pixels should be
252 rejected (van Hoek (2011), Liebske 2015) but for phase abundance, spatial distribution and textural
253 shape they must be considered (Liebske 2015). For the KNN algorithm it is important that the phases
254 within the reference dataset correspond approximately to the phases present in the sample. If the
255 number of phases in the reference dataset greatly exceeds that in the sample, there is a high
256 probability that the convolved boundary phase pixels will have a composition similar to a phase in
257 the reference dataset and be misidentified. When evaluating phases it is important to be able to
258 distinguish between variance due to convoluted pixels and actual variation due to chemical variation
259 within a phase or the presence of multiple phases. The example of Phase-PC maps for serpentine
260 and bultfonteinite (Figure 5 b-c) show the importance of spatial information in making this
261 assessment.

262 Principal component maps demonstrate how the generation of a phase map need not be the end of
263 the process. New phases may be identified allowing the initial phase map to be revised. A phase
264 could be subdivided based on principal components and its chemistry extracted or phase map
265 algorithms could be rerun with an additional specified cluster. Where Phase-PC maps suggest
266 variations within phases, e.g. as shown in the compositional difference between serpentine in the
267 kimberlite and the basalt xenolith, the user can further investigate thus improving the sample
268 characterisation.

269 Conclusions

270 In agreement with other work (e.g. Munch et al. 2015) the data presented illustrates the need for
271 phase maps to be subjected to critical analysis, exposing any limitations of the algorithm or
272 operating conditions resulting in incorrect classification. The performance of phase algorithms will
273 vary depending on the sample (Munch et al. 2015) and the input parameters (the number of phases
274 for K-means using random clusters; the phases specified for KNN and K-means using specified
275 clusters), making it important to check the data has been correctly classified. Principal component
276 maps provide an easy solution to evaluate phase classification. RGB PC images provide a good visual
277 reference for checking phase maps, more clearly discriminating between phases than BSE images.
278 Phase PC maps provide a good method of assessing variation within phases and identifying
279 unclassified phases with the spatial information important for discriminating real chemical variance
280 from convoluted pixels. The role of an operator in checking phase maps introduces subjectivity but
281 the provision of spatial information allows the operator to make high-quality decisions as to the
282 nature of variance, resulting in robust sample characterisation. This process of evaluation of phase
283 maps allows further refinement and can provide additional information about a sample prompting
284 further investigation.

285 KNN is potentially a very useful method of phase classification for geological samples, where the
286 analyst is familiar with the possible phases within the rock sample. It produces consistent results
287 similar to manual thresholding (e.g. Muir et al. 2012). K-means with specified clusters produces
288 similar results but is more affected by the absence of a particular phase, when shifting from area to
289 area within or between rock samples.

290 Acknowledgements

291 The authors would like to thank the reviewers for their insightful comments and revisions which
292 improved the manuscript.

293 References

- 294 Buse, B., Schumacher, J.C., Sparks, R.S.J. & Field, M. (2010). Growth of bultfonteinite and
295 hydrogarnet in metasomatized basalt xenoliths in the B/K9 kimberlite, Damtshaa, Botswana: insights
296 into hydrothermal metamorphism in kimberlite pipes. *Contrib Mineral Petrol* **160**, 533-550
- 297 Jolliffe, I.T. (2002). *Principal Component Analysis* Second Edition. New York, USA: Springer-Verlag
- 298 Kotula, P.G., Keenan, M.R. & Michael, J.R. (2003). Automated analysis of SEM X-ray spectral images:
299 A powerful new microanalysis tool. *Microsc Microanal* **9**, 1-17.
- 300 Liebske, C. (2015). iSpectra: An Open Source Toolbox for the analysis of spectral images recorded on
301 scanning electron microscopes. *Microsc Microanal* **21**, 1006-1016.
- 302 MALOY, A.K. & TREIMAN, A.H. (2007). Evaluation of image classification routines for determining
303 modal mineralogy of rocks from X-ray maps. *Am Mineral* **92**, 1781–1788.
- 304 Mori, N., Kato, N. & Morita, M. (2017). Automatic processing of element maps by automatic colour
305 map filter and high speed cluster analyses for EPMA. EMAS 2017 Conference abstract.
- 306 Muir, D. D., Blundy, J.D. & Rust, A.C. (2012). Multiphase petrography of volcanic rocks using element
307 maps: a method applied to Mount St Helens, 1980-2005. *Bull Volcanol* **74**, 1101-1120
- 308 Munch, B., Martin, L.H.J. & Leemann, A. (2015). Segmentation of elemental EDS maps by means of
309 multiple clustering combined with phase identification. *J Microsc* **260**, 411-426.

310 Statham, P., Penman, C., Chaldecott, J., Burgess, S., Sitzman, S. & Hyde, A. (2013). Validating a New
311 Approach to the Mapping of Phases by EDS by Comparison with the Results of Simultaneous Data
312 Collection by EBSD. *Microsc Microanal* **19 (S2)**, 752-753

313 Tan, P.N., Steinbach, M., Kumar, V. (2006). Introduction to Data Mining. Boston, USA: Pearson
314 Education, Inc.

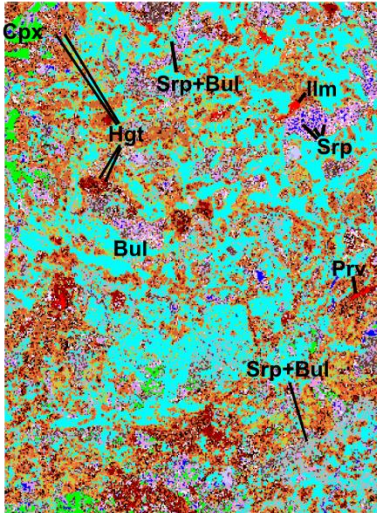
315 van den Berg, R.A, Hoefsloot, H.C.J., Westerhuis, J.A., Smilde, A.K. & van der Werf M.J. (2006).
316 Centering, scaling and transformations: improving the biological information content of
317 metabolomics data. *BMC Genomics* **7**, 142

318 VAN HOEK, C.J.G., DE ROO, M., VAN DER VEER, G. & VAN DER LAAN, S.R. (2011). A SEM-EDS study of
319 cultural heritage objects with interpretation of constituents and their distribution using PARC data
320 analysis. *Microsc Microanal* **17**, 656–660.

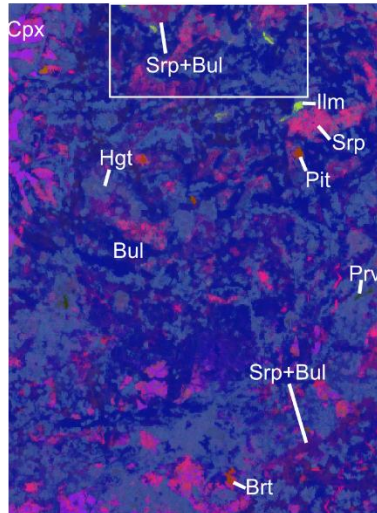
321

322 Figure 1. (a) BSE image showing mapped area from which the phase map was generated. Figures 2 and
323 3 correspond to expanded regions highlighting features within the phase maps. (b) RGB composite
324 image of the first three principle components. Locations from which the area compositions were
325 extracted for each identified phase are shown. (c) BSE Image showing the finely intergrown phase
326 hydrogarnet and serpentine (purple colour on 1b). Location of image is given by white rectangle marked
327 'c' on 1b. Mineral abbreviations are as follows: Bul - bultfonteinite, Hgt - hydrogarnet and Srp -
328 serpentine.
329
330

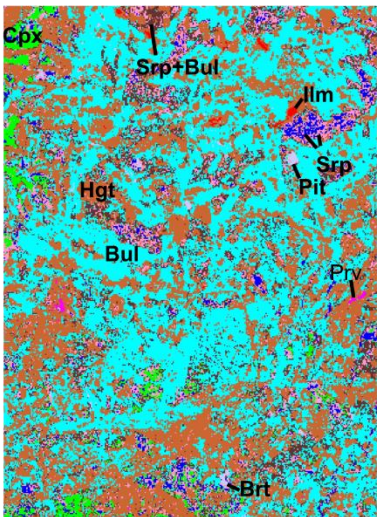
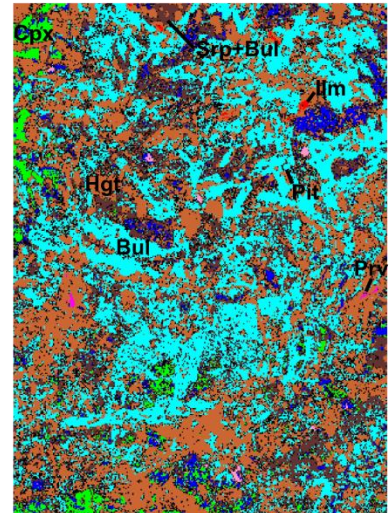
(a) K-means (15 random clusters)



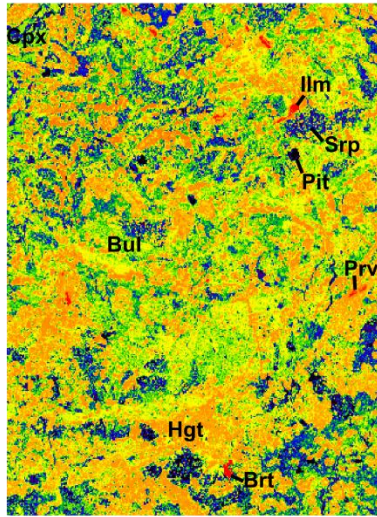
(b) RGB Principle components



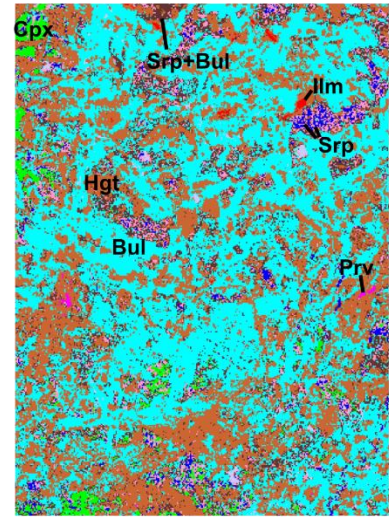
(c) KNN



(d) Kmeans (specified initial clusters)

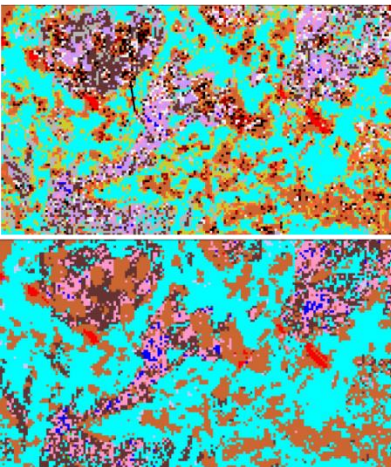


(e) BSE



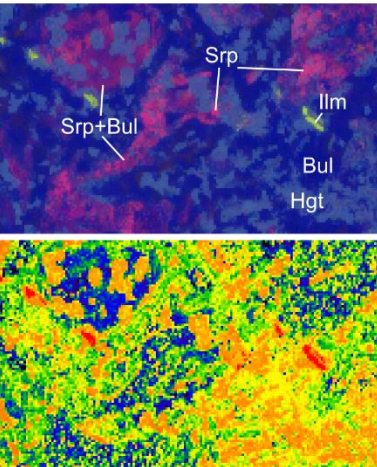
(f) Kmeans (max element for initial clusters)

(g) K-means (15 random clusters)



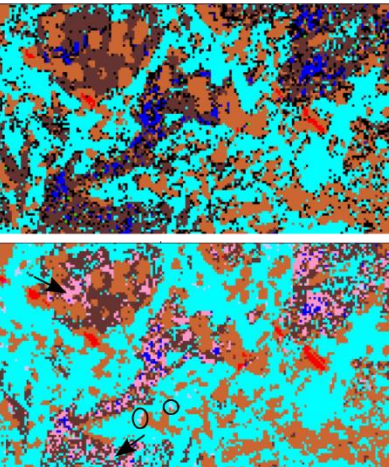
(j) Kmeans (specified initial clusters)

(h) RGB Principle components



(k) BSE

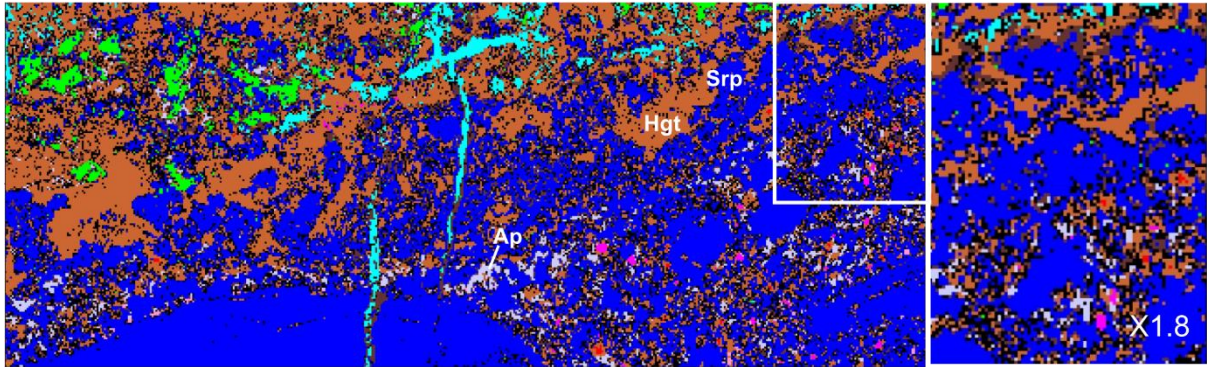
(i) KNN



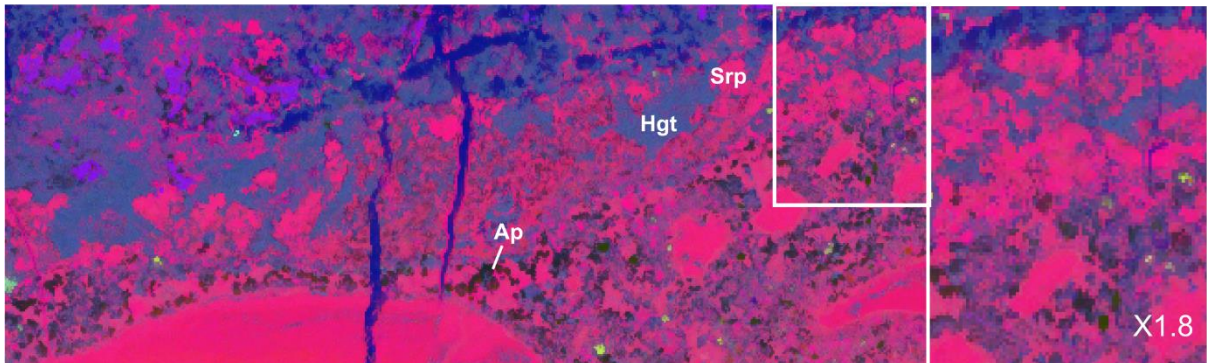
(l) Kmeans (max element for initial clusters)

333 Figure 2. Comparison of phase maps with a false-coloured BSE image and RGB-PC image. The area
334 shown is a small region of the map, part of the basalt xenolith. The location is given on Figure 1a. For
335 the KNN algorithm (c) black pixels are unclassified pixels. High magnification images (g-l) are shown for
336 the area represented by the white box in (b). Arrows in (l) correspond to areas where serpentine-
337 bultfonteinite is underrepresented in comparison to (h) and (i). Circles in (l) correspond to misclassified
338 convoluted pixels. Mineral abbreviations are as follows: Brt - barite, Bul - bultfonteinite, Cpx -
339 clinopyroxene, Hgt - hydrogarnet, Ilm - ilmenite, Prv - perovskite and Srp – serpentine.
340

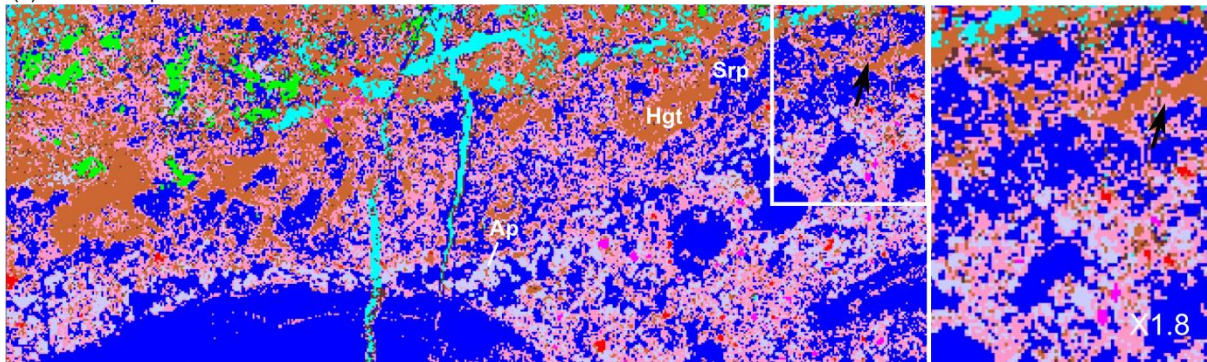
(a) KNN



(b) RGB-PC

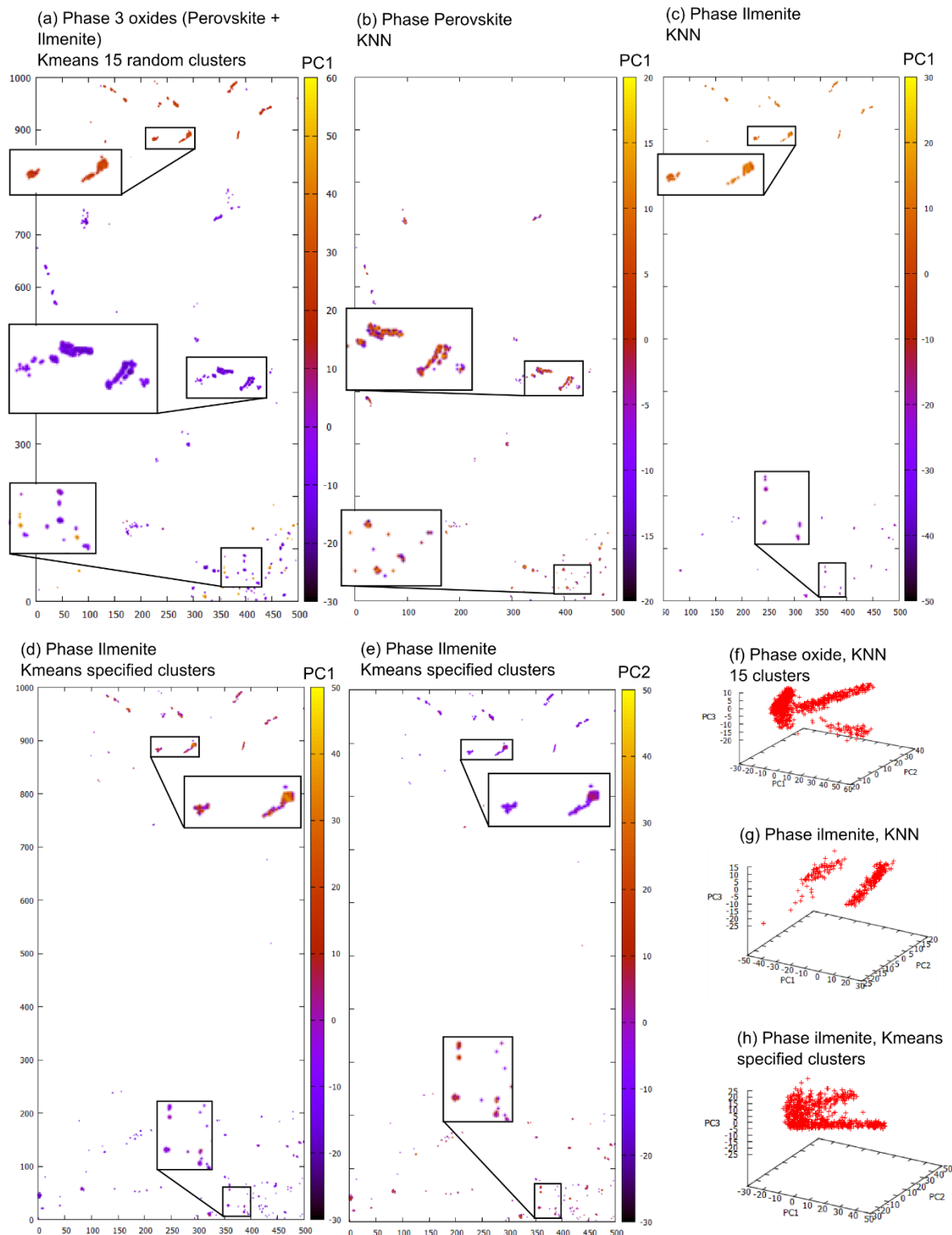


(c) K-means specified



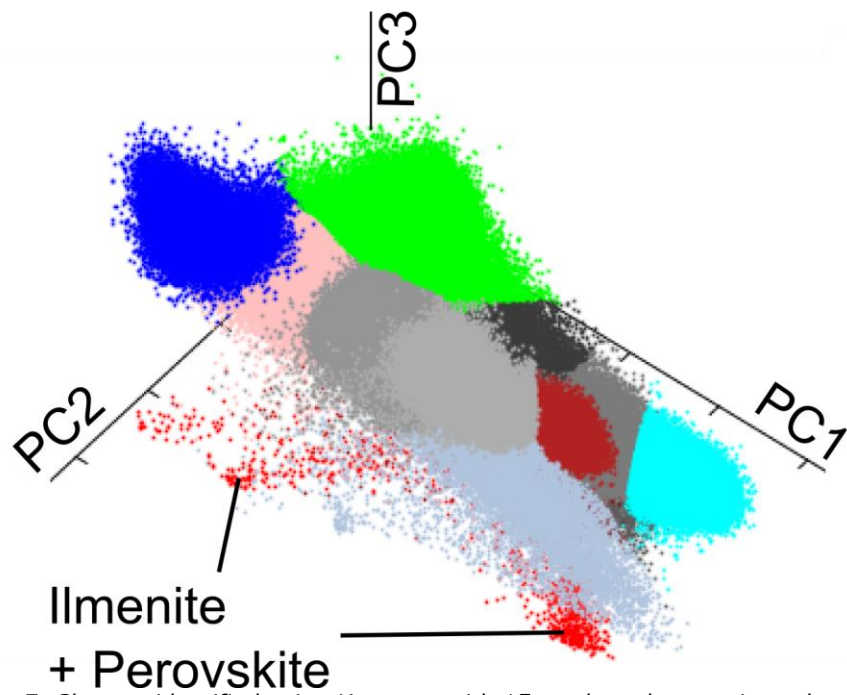
341
342
343
344
345

Figure 3. Comparison of phase maps with a RGB-PC image for the kimberlite region of the map. The location is show on Figure 1a. High magnification images are shown to the right for the area represented by the white box. For KNN algorithm (a) black pixels are unclassified pixels.



346
347
348
349
350
351

Figure 4. Phase-PC Maps (a-e) examining variation within the oxide, perovskite and ilmenite phases. High magnification insets are shown for the bottom, middle and top of the maps, which correspond to the kimberlite, xenolith margin and xenolith interior. PC scatter graphs (f-h) show the alignment of orthogonal principle components with data variance.



352
353
354
355
356

Figure 7. Clusters identified using K-means with 15 random clusters. Low abundance phases ilmenite and perovskite are identified as a single cluster.

357 **Table 1.** List of phases identified from RGB composite image of the first three principle components. Numbers
358 correspond to those given on Figure 1b.

Phase	Mineral	Ideal Formula
1	Ilmenite	FeTiO_3
2	Pit	
3	Bultfonteinite	$\text{Ca}_2\text{SiO}_2(\text{OH},\text{F})_4$
4	Clinopyroxene	$\text{Ca}(\text{Mg},\text{Fe})\text{Si}_2\text{O}_6$
5	Serpentine/Chlorite	$(\text{Mg},\text{Fe},\text{Mn},\text{Al})_{12}(\text{SiAl})_8\text{O}_{20}(\text{OH})_{16}$
6	Perovskite	CaTiO_3
7	Hydrogarnet	$\text{Ca}_3(\text{Fe},\text{Ti},\text{Al})_2\text{Si}_2\text{O}_8(\text{OH})_4$
8	Serpentine + Bultfonteinite intergrowth	
9	Sr-Apatite	$(\text{CaSrBa})_5(\text{PO}_4)_3(\text{OH},\text{F})$

359

360

361

362 **Table 2.** Average phase compositions and standard deviations for K-means clustering using 15 random
 363 clusters. Phase identification was made with reference to spatial distribution, BSE and RGB-PC images. High
 364 abundance elements are shown in bold. Standard deviation % is given in italic.

Identification	Na2O	MgO	Al2O3	SiO2	K2O	CaO	TiO2	FeO	MnO	Total
Hydrogarnet	0.08 <i>690.58</i>	4.06 <i>47.43</i>	7.03 <i>26.51</i>	29.97 <i>6.37</i>	0.04 <i>412.50</i>	31.50 <i>6.20</i>	1.36 <i>72.00</i>	13.18 <i>13.76</i>	0.22 <i>167.54</i>	87.44
Srp+Bul	0.20 <i>281.50</i>	11.61 <i>27.15</i>	2.12 <i>67.94</i>	32.84 <i>12.60</i>	0.04 <i>484.84</i>	29.15 <i>10.38</i>	0.37 <i>133.59</i>	4.32 <i>55.83</i>	0.09 <i>336.37</i>	80.74
Oxides	0.71 <i>157.87</i>	2.63 <i>151.71</i>	2.90 <i>75.29</i>	8.96 <i>66.45</i>	0.05 <i>377.52</i>	26.39 <i>39.64</i>	36.74 <i>32.41</i>	12.66 <i>113.57</i>	0.23 <i>219.36</i>	91.27
Hydrogarnet	0.11 <i>508.51</i>	1.54 <i>84.70</i>	2.57 <i>55.00</i>	26.71 <i>9.55</i>	0.01 <i>899.97</i>	40.14 <i>4.75</i>	0.43 <i>98.20</i>	5.97 <i>33.63</i>	0.11 <i>292.45</i>	77.59
Hydrogarnet	0.07 <i>743.48</i>	2.01 <i>85.81</i>	7.36 <i>34.04</i>	24.95 <i>8.50</i>	0.03 <i>464.74</i>	33.22 <i>6.44</i>	3.07 <i>104.75</i>	11.75 <i>16.55</i>	0.18 <i>199.97</i>	82.64
Serpentine	0.16 <i>310.02</i>	23.61 <i>13.17</i>	5.79 <i>62.19</i>	33.52 <i>12.43</i>	0.21 <i>261.26</i>	10.44 <i>36.97</i>	0.46 <i>110.00</i>	7.16 <i>38.96</i>	0.15 <i>214.63</i>	81.51
Hydrogarnet	0.09 <i>591.64</i>	2.26 <i>83.25</i>	4.56 <i>35.22</i>	27.84 <i>7.70</i>	0.02 <i>552.26</i>	36.47 <i>5.36</i>	0.81 <i>74.91</i>	10.55 <i>18.29</i>	0.19 <i>190.54</i>	82.80
Serpentine	0.08 <i>521.29</i>	31.68 <i>10.85</i>	4.69 <i>55.41</i>	36.12 <i>10.33</i>	0.12 <i>393.15</i>	3.04 <i>93.63</i>	0.13 <i>175.92</i>	4.68 <i>45.56</i>	0.11 <i>268.42</i>	80.65
Augite	0.67 <i>127.56</i>	15.47 <i>22.36</i>	1.45 <i>110.15</i>	49.23 <i>9.87</i>	0.02 <i>901.96</i>	20.79 <i>16.23</i>	0.56 <i>64.61</i>	6.48 <i>29.59</i>	0.13 <i>246.93</i>	94.80
Bul+Srp	0.14 <i>379.70</i>	6.39 <i>36.03</i>	1.37 <i>78.49</i>	28.78 <i>10.10</i>	0.02 <i>686.74</i>	36.60 <i>7.04</i>	0.19 <i>149.18</i>	2.41 <i>75.24</i>	0.06 <i>519.18</i>	75.95
Hydrogarnet	0.05 <i>964.25</i>	1.88 <i>82.19</i>	6.86 <i>25.98</i>	27.04 <i>7.49</i>	0.04 <i>416.93</i>	33.36 <i>6.79</i>	1.32 <i>82.74</i>	15.96 <i>13.95</i>	0.23 <i>161.73</i>	86.75
Apatite	0.76 <i>137.05</i>	7.74 <i>56.73</i>	3.32 <i>124.80</i>	14.92 <i>35.46</i>	0.08 <i>207.64</i>	28.55 <i>23.65</i>	0.73 <i>177.01</i>	4.19 <i>64.07</i>	0.09 <i>347.44</i>	60.38
Serpentine	0.18 <i>317.18</i>	16.98 <i>19.57</i>	5.14 <i>62.35</i>	30.61 <i>14.06</i>	0.11 <i>285.14</i>	19.10 <i>17.29</i>	0.83 <i>122.78</i>	8.16 <i>41.28</i>	0.16 <i>210.56</i>	81.27
?	0.10 <i>529.96</i>	9.70 <i>26.32</i>	6.02 <i>41.95</i>	28.58 <i>11.15</i>	0.06 <i>323.73</i>	26.35 <i>11.17</i>	1.42 <i>105.09</i>	11.94 <i>24.38</i>	0.20 <i>175.08</i>	84.38
Bulfonteinite	0.11 <i>466.91</i>	1.29 <i>105.34</i>	0.64 <i>120.89</i>	27.57 <i>8.64</i>	0.01 <i>1202.30</i>	44.12 <i>5.07</i>	0.10 <i>213.24</i>	1.16 <i>101.80</i>	0.03 <i>1110.73</i>	75.02

365

366

367

368 **Table 3.** Comparison of phase average and extracted composition for a selected area. Data is for the
 369 clinopyroxene phase, with phase average calculated from K-means clustering using 15 random clusters. High
 370 abundance elements are shown in bold. Standard deviation % is given in italic. Atoms per formula unit (apfu)
 371 are calculated on the basis of 6 oxygens.

Wt. %	Na2O	MgO	Al2O3	SiO2	K2O	CaO	TiO2	FeO	MnO	Total
Phase average	0.67 <i>127.56</i>	15.47 22.36	1.45 <i>110.15</i>	49.23 9.87	0.02 <i>901.96</i>	20.79 16.23	0.56 <i>64.61</i>	6.48 <i>29.59</i>	0.13 <i>246.93</i>	94.80
Spot extract	0.18 <i>410.93</i>	15.40 9.22	0.18 <i>226.55</i>	55.16 3.52	-0.04 <i>279.49</i>	23.49 9.70	0.30 55.15	4.70 <i>33.00</i>	0.19 <i>170.59</i>	99.57 3.85
apfu	Na	Mg	Al	Si	K	Ca	Ti	Fe	Mn	Total
Phase average	0.051	0.902	0.067	1.926	0.001	0.871	0.016	0.212	0.004	4.050
Spot extract	0.013	0.843	0.008	2.024	0.000	0.924	0.008	0.144	0.006	3.970

372

373

374

375 **Table 4** Extracted compositions for the discrete groupings identified from phase-PC1 map of ilmenite
 376 phase (KNN algorithm). For comparison the extracted composition for ilmenite used in the reference
 377 dataset is given. Standard deviation are given in italic.

PC1 threshold (location)	Na2O	MgO	Al2O3	SiO2	K2O	CaO	TiO2	FeO	MnO	Identification
<0 (In kimberlite)	0.00 <i>0.67</i>	14.01 <i>4.60</i>	4.05 <i>1.80</i>	3.89 <i>3.32</i>	0.00 <i>0.17</i>	4.68 <i>3.36</i>	16.42 <i>4.61</i>	48.37 <i>8.09</i>	0.92 <i>0.57</i>	Spinel
>0 (In xenolith)	0.00 <i>0.91</i>	2.13 <i>1.23</i>	2.25 <i>1.84</i>	7.70 <i>4.66</i>	0.00 <i>0.15</i>	9.45 <i>5.55</i>	39.58 <i>7.59</i>	33.62 <i>5.71</i>	0.71 <i>0.55</i>	ILM
reference composition (from xenolith)	1.97	3.43	0.00	0.00	0.00	0.92	46.93	43.15	0.10	ILM

378

379

380

381 **Table 5** Extracted compositions for the discrete groupings of serpentine identified in the phase-PC
 382 map. Olivine pseudomorphs are the large sub-rounded grains within the kimberlite which were
 383 original olivine but have been replaced by serpentine.

PC1 Threshold (location)	Na2O	MgO	Al2O3	SiO2	K2O	CaO	TiO2	FeO	MnO	Total
> 0 (xenolith)	0.08	33.20	3.35	38.86	0.10	3.89	0.15	4.02	0.08	83.74
< -4 (Kimberlite: rims of olivine pseudomorphs)	0.09	28.79	7.56	32.05	0.21	2.12	0.12	6.25	0.17	77.36
<0 and >-4 (Kimberlite: olivine pseudomorphs)	0.07	30.29	5.91	33.60	0.15	2.26	0.11	5.29	0.14	77.81

384

385

386