

1 **Machine Learning Investigation of Clinopyroxene Compositions to Evaluate and Predict**

2 **Mantle Metasomatism Worldwide**

3 Ben Qin^{1,2*}, Fang Huang³, Shichun Huang⁴, Andre Python⁵, Yunfeng Chen^{1,2}, J Zhang Zhou^{1,2*}

4 1. School of Earth Sciences, Zhejiang University, Hangzhou 310027, China.

5 2. Key Laboratory of Geoscience Big Data and Deep Resource of Zhejiang Province,
6 Hangzhou 310027, China.

7 3. CSIRO Mineral Resources, Kensington, WA 6151, Australia.

8 4. Department of Geoscience, University of Nevada, Las Vegas, 4505 S. Maryland Pkwy, Las
9 Vegas 89154, USA.

10 5. Center for Data Science, Zhejiang University, Hangzhou, China.

11

12 Corresponding author: Ben Qin (charlesbenq@zju.edu.cn); J Zhang Zhou

13 (zhangzhou333@zju.edu.cn)

14

15 **Key points:**

16 • Mantle clinopyroxene major and trace element data compiled and evaluated with machine
17 learning models

18 • Accuracy comparisons between low- and high-dimensional dataspaces reveal the most
19 important features for classification

20 • Machine learning models identify clusters of mantle metasomatism worldwide

21

22 **Plain Language Summary**

23 Clinopyroxene is a major mineral in Earth's upper mantle. Previous studies have
24 attempted to discriminate between reactions modifying the mantle by plotting clinopyroxene
25 major and trace element compositions in two-dimensional (2-D) diagrams. However, these 2-

26 D methods show poor accuracy when applied to global datasets. Therefore, we suggest a
27 machine learning approach to evaluate clinopyroxene compositional data in higher dimensions.
28 Our results demonstrate that machine learning can significantly improve the accuracy of
29 clinopyroxene compositional predictions over classical methods utilizing elemental ratios.
30 Furthermore, the application of our algorithm to a global clinopyroxene dataset suggests that
31 mantle metasomatism is globally widespread.

32

33 **Abstract**

34 Clinopyroxene major and trace element compositions document their physicochemical
35 evolution and have been widely used to detect mantle metasomatism. Classical methods
36 typically rely on one or several elemental ratios such as Ca/Al, Mg/Fe, La/Yb, and Ti/Eu to
37 determine whether rocks or minerals have been metasomatized. These methods have proven
38 useful at specific sites, but not globally. In this study, we used machine learning methods to
39 classify the chemical compositions of clinopyroxenes from mantle xenoliths and examine their
40 relationship with mantle metasomatism. We compiled major element data from 8,713
41 clinopyroxene samples (21,605 analyses) and trace element data from 1,235 clinopyroxene
42 samples (2,967 analyses). Samples were labeled “positive” if clearly affected by patent
43 metasomatism based on petrographic evidence, “negative” if clearly unaffected by
44 metasomatism, or were left unlabeled if neither case applied. We then trained an XGBoost
45 machine learning model, which achieved higher accuracy than traditional methods using a
46 limited number of elemental ratios. Our results identify numerous locations with high mean
47 probabilities of mantle metasomatism and show variability in the probability distributions
48 observed across locations worldwide. These results indicate that metasomatism may be
49 globally widespread, but the probability of metasomatism is not correlated with geophysical

50 parameters such as crustal thickness, lithospheric thickness, or mantle *S*-wave velocity. Hence,
51 the spatial distribution of metasomatism appears mainly driven by unobserved factors.

52

53 **1. Introduction**

54 Metasomatism modifies the mineralogy and composition of pre-existing rocks through
55 reaction with melt/fluid at high temperature. This important process produces geochemical and
56 isotopic heterogeneities within Earth's mantle (Aiuppa et al., 2021; Roden & Murthy, 1985;
57 Wang et al., 2022; Zhang et al., 2009), which in turn affect chemical differentiation in the
58 mantle, craton stability, and the physical properties of the lithosphere (Araújo et al., 2009;
59 Dawson, 1984; Liu et al., 2021; Lloyd & Bailey, 1975; Menzies & Murthy, 1980; O'Reilly &
60 Griffin, 2013; Pearson et al., 2021; Peng et al., 2021; Rudnick et al., 1993). Therefore,
61 evaluating mantle metasomatism at the global scale is essential to understanding mantle
62 heterogeneity.

63 The composition of the upper mantle has been investigated through various approaches.
64 The chemical compositions of basalts and petrological models of peridotite-basalt melting
65 relationships have provided crucial information about the mantle (McDonough & Sun, 1995;
66 Ringwood, 1962). Mantle xenoliths brought to the surface as inclusions in basalts and peridotite
67 massifs directly sample the upper mantle and thus provide direct insights into upper mantle
68 processes and compositions. Knowledge of the petrography and geochemistry of these
69 peridotite samples is therefore likely to provide important information, particularly regarding
70 the nature of upper-mantle partial melting, fractional crystallization, and metasomatism
71 processes (Frey & Green, 1974; Frey & Prinz, 1978).

72 The occurrence of metasomatism can be revealed by the petrography and geochemistry of
73 mantle xenoliths. Patent metasomatism is identified straightforwardly based on the presence of
74 secondary minerals in xenoliths. In comparison, cryptic metasomatism and stealth

75 metasomatism are not so easy to identify because the former does not produce new phases
76 (Dawson, 1984) whereas the latter only produces minerals that are indistinguishable from
77 common mantle minerals (e.g., clinopyroxene or garnet; Griffin et al., 2009). Therefore, it is
78 often extremely difficult to identify mantle metasomatism based on petrographic observations
79 alone, and researchers instead turn to the geochemistry of xenolithic minerals. Here, we focus
80 on clinopyroxene compositions for three reasons: 1) clinopyroxene is abundant in the upper
81 mantle, 2) it is a substantial reservoir for various minor and trace elements, and 3) its
82 composition is sensitive to mantle metasomatism. Many studies have proposed using the
83 chemical composition of clinopyroxene to identify metasomatism, including the ratios Ca/Al
84 (Rudnick et al., 1993; Wang et al., 2010), Mg/(Mg + Fe) (Mg#) (Yaxley & Green, 1998), La/Yb
85 (Coltorti et al., 1999; Zong & Liu, 2018), and Ti/Eu (Coltorti et al., 1999; Zong & Liu, 2018).
86 These elemental ratios are based on petrological models describing reactions between silicate
87 minerals and metasomatic agents (e.g., carbonate or silicate melts), which enrich clinopyroxene
88 in incompatible elements (Dalou et al., 2009; Green et al., 1992; Klemme et al., 1995; Rudnick
89 et al., 1993; Sweeney et al., 1995).

90 Many major and trace element analyses have been reported for clinopyroxenes in xenolith
91 samples worldwide (Figure 1; see section 2 for details on data selection). This highlights the
92 challenge of directly estimating the scale and extent of mantle metasomatism, which requires
93 the development of effective analytical tools capable of exploiting information from indirect
94 indicators of major/trace element compositions. Although elemental ratios have been
95 successfully used to identify metasomatism at specific sites, they have proven inaccurate when
96 applied globally (Figure 2; see section 4 for more details), suggesting that the spatial
97 distribution of metasomatism worldwide is driven by processes that cannot be captured by
98 elemental ratios alone. A chemical model quantifying mantle metasomatism at the global scale
99 has yet to be developed.

100 Whereas elemental ratios do not seem to capture the main drivers of the spatial distribution
101 of metasomatism worldwide, relevant information may be obtained from large-volume, high-
102 dimensional geochemical data. In recent years, machine learning (ML) approaches have been
103 applied to mineralogy, petrology, and geochemistry datasets to provide new insights and
104 identify trends and patterns that would otherwise be unobservable (e.g., Petrelli & Perugini,
105 2016; Petrelli et al., 2020; Thomson et al., 2021; Ueki et al., 2018; Valetich et al., 2021; Wang
106 et al., 2021; Zhao et al., 2019), demonstrating their potential to quantify mantle metasomatism
107 worldwide.

108 Here, we compiled a global dataset of major and trace element compositions of
109 clinopyroxenes from mantle xenoliths and trained a supervised ML algorithm (XGBoost; Chen
110 & Guestrin, 2016) to classify metasomatism in high-dimensional space. We also trained
111 unsupervised machine learning models to ensure that the labeled training and testing dataset
112 and unlabeled application dataset had similar distributions. Finally, we applied our trained ML
113 model to predict the probability of the occurrence of metasomatism at the global scale.

114

115 **2. Data compilation and labeling**

116 We downloaded compositional data for clinopyroxenes from mantle xenoliths from 972
117 locations worldwide (Figure 1) from the GEOROC database ([http://georoc.mpch-
118 mainz.gwdg.de/georoc/](http://georoc.mpch-mainz.gwdg.de/georoc/); accessed 14 July 2020) (Sarbas, 2008). Each location includes
119 multiple samples and analyses (Figure S1). To exclude unreliable samples, we used only
120 clinopyroxenes with 40–60 wt.% SiO₂, <40 wt.% MgO, <30 wt.% FeO^T (the superscript ‘T’
121 indicates total iron), <26 wt.% CaO, and oxide totals of 98.5–101.5 wt.%. Elements missing
122 from >60% of the entire dataset were not considered.

123 After this initial filtering, our “Parent” dataset contained 21,605 observations (rows)
124 corresponding to clinopyroxene major element analyses (SiO₂, TiO₂, Al₂O₃, Cr₂O₃, FeO^T, CaO,

125 MgO, MnO, and Na₂O), and 2,967 rows of trace element analyses (including Sc, Ti, V, Cr, Ni,
126 Rb, Sr, Y, Zr, Nb, Ba, La, Ce, Pr, Nd, Sm, Eu, Gd, Tb, Dy, Ho, Er, Tm, Yb, Lu, Hf, Ta, Pb,
127 Th, and U) (Figure S2). In general, the proportion of missing values varies among elements,
128 with major element data being rather complete, e.g., 5% missing analyses for Cr₂O₃ (1,024 of
129 21,605), and a higher proportion of missing data for trace elements, e.g., 32% missing analyses
130 for Ti (961 of 2,967) (Figure S2).

131 The Parent dataset was further divided into a labeled training and testing subset (“Labeled
132 dataset”) and an unlabeled application subset (“Application dataset”). The Labeled dataset was
133 used to classify clinopyroxene as being affected or unaffected by metasomatism. We classified
134 clinopyroxenes in the Labeled dataset as being affected (“positive”) or unaffected by
135 metasomatism (“negative”) based on the petrographic descriptions provided in the original
136 literature. A sample was labeled “positive” if its petrographic description contains evidence of
137 metasomatic phases, including silicate glass, calcite, hornblende, phlogopite, and apatite. In
138 contrast, a sample was labeled “negative” if it shows a monotonous increase in chondrite-
139 normalized (McDonough & Sun, 1995) light rare earth element (LREE) concentrations (i.e.,
140 La, Ce, Pr, Nd, Sm, and Gd). Based on these criteria, 1,650 major and 539 trace element
141 analyses were labeled “positive” and 439 and 333 respective analyses were labeled “negative”.
142 In the Parent dataset, most examples (>70%; i.e., 19,516 major and 2,095 trace analyses) did
143 not contain petrographic descriptions of metasomatic minerals in the original literature; this
144 unlabeled Application dataset was used to test the unsupervised ML algorithm (see Section
145 3.2).

146

147 **3. Methods**

148 We employed a three-step modeling process (Figure 3). First, supervised learning models
149 were trained based on the Labeled dataset. Second, unsupervised learning models were

150 implemented on the Parent dataset to verify that the overall data distributions of the Labeled
151 and Application datasets were similar. Third, the optimal classification model obtained in step
152 one was applied to the unlabeled clinopyroxene compositions in the Application dataset. In this
153 study, all models were implemented using the scikit-learn Python package (Pedregosa et al.,
154 2011).

155

156 **3.1 Training the supervised learning models**

157 In the first step, we trained supervised ML models to classify clinopyroxenes as a binary
158 variable (1 if affected by metasomatism, 0 otherwise). We tested several ML algorithms,
159 including Random Forest (Breiman, 2001) and Support Vector Machines (Boser et al., 1992),
160 and eventually chose XGBoost (Chen & Guestrin, 2016) due to its flexibility, predictive
161 performance, computational efficiency, and interpretability. Importantly, the Random Forest
162 and Support Vector Machines algorithms are not designed to handle missing values (Boser et
163 al., 1992; Breiman 2001), which are frequent in our dataset. In contrast, XGBoost can
164 accommodate sparse feature formats and can automatically identify the best imputation value
165 for missing values based on reduction on training loss (Chen & Guestrin, 2016). Furthermore,
166 in addition to its high predictive capability and computational efficiency, the tree structure of
167 XGBoost facilitates interpretation of the results (Azodi et al., 2020), which is important for
168 identifying features associated with the occurrence of metasomatism. We directly used the
169 elemental data without any preliminary transformation as the input into the XGBoost
170 classification algorithm.

171 XGBoost is based on a gradient-boosting decision tree method (Friedman, 2001) and has
172 been recently applied in a wide range of applications aiming to predict complex spatial
173 phenomena at the global scale (e.g., Cook-Patton et al., 2020; Python et al., 2021; Zheng et al.,
174 2021). XGBoost uses a gradient-descent algorithm to minimize the loss when adding new

175 models. In practice, it continuously adds trees to fit the residuals of the previous prediction,
176 and the predictions are computed as the sum of the effects of all trees. For a dataset with
177 n observations, label element $\mathbf{y}_i \in R$ with $i = \{1, \dots, n\}$, and m features composed of
178 feature element $\mathbf{x}_i \in R^m$ with $j = \{1, \dots, m\}$, the predictions \hat{y}_i are obtained by summing
179 the scores obtained in the corresponding leaves, which is expressed as (Chen & Guestrin, 2016):

180
$$\hat{y}_i = \phi(\mathbf{x}_i) = \sum_{k=1}^K f_k(\mathbf{x}_i), \text{ with } f_k \in \mathcal{F}, \quad (\text{Eq. 1})$$

181 where $\mathcal{F} = \{f(\mathbf{x}) = w_{q(\mathbf{x})}\} (q : R^m \rightarrow T, w \in R^T)$ is the space of the regression trees, additive
182 function tree $k = \{1, \dots, K\}$, and each f_k corresponds to an independent tree structure q and
183 leaf weight w . Here, q represents the structure of each tree that maps an observation to a
184 corresponding leaf, with T the total number of leaves in the tree. $w_{q(\mathbf{x})}$ represents the set of
185 scores computed in all leaf nodes in a tree. We used XGBoost within a classification framework
186 since the label y_i is binary (1 if affected by metasomatism, 0 otherwise). For each observation
187 i , the output of the classification \hat{y}_i represents the probability that metasomatism is present.
188 In this classification framework, \hat{y}_i is calibrated as a probability by taking only values
189 between 0 and 1. To compute the elements of the confusion matrix, we dichotomize \hat{y}_i as
190 equal to 1 if $\hat{y}_i \geq 0.5$ and 0 otherwise.

191 To minimize bias and variance in the predictive scores, we performed a ten-fold cross-
192 validation procedure (Kohavi, 1995) by randomly splitting the Labeled dataset into training
193 (70%) and testing subsets (30%). Therefore, overfitting and randomness can be mitigated by
194 cross-validation and the splitting of the training and testing subsets utilized to evaluate the
195 classifier performance. We applied Grid Search Cross-Validation (from the scikit-learn
196 package), which aims to find an optimal hyperparameter combination (eta, gamma, max depth,

197 and alpha) through an iterative grid-search process. The procedure trains 9,000 candidate
198 models and selects the model with the best predictive performance via ten-fold cross-validation.

199 To evaluate the performance of XGBoost models, several classification metrics can be
200 defined based on the confusion matrix (Stehman, 1997), a specific table layout that visualizes
201 model performance. Each row of the confusion matrix represents the instances in an actual
202 class, whereas each column represents the instances in a predicted class. We use Accuracy and
203 the F1 score (Dice, 1945; Sørensen, 1948) described below:

204 Accuracy is the ratio of the total number of correct “positive” and “negative” predictions
205 to the total number of known “positive” and “negative” cases:

$$206 \quad \text{Accuracy} = \frac{\text{Correct positive + negative predictions}}{\text{Known positive + negative cases}}. \quad (\text{Eq. 2})$$

207 The F1 score is the harmonic mean of Precision and Recall:

$$208 \quad \text{F1 score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}, \quad (\text{Eq. 3})$$

209 where Precision is a measure of accuracy provided that a specific class (here, “positive”) has
210 been predicted:

$$211 \quad \text{Precision} = \frac{\text{Correct positive predictions}}{\text{All positive predictions}}. \quad (\text{Eq. 4})$$

212 Recall is a measure of the ability of a model to select instances of a certain class (again,
213 “positive” here) from a dataset:

$$214 \quad \text{Recall} = \frac{\text{Correct positive predictions}}{\text{Known positive cases}}. \quad (\text{Eq. 5})$$

215 Although Accuracy is a common and direct way to evaluate and improve models, the F1
216 score can give a better measure of the incorrectly classified cases than Accuracy. The F1 score
217 is more suitable when the classes are imbalanced. Given the class imbalance observed in our
218 dataset, we favor the F1 score to compare the predictive performance of our models.

219 To test the impacts of feature selection on the ML results, we trained the XGBoost model
220 on different major and trace element subsets selected from the Labeled dataset (Table 1). For
221 major elements, we compared the relative feature importance when using two traditional
222 element pairs (CaO and Al₂O₃; MgO and FeO^T), all four of those elements (i.e., both pairs
223 combined), and all nine major elements. For trace elements, we also used two traditional
224 element pairs (Eu and Ti; La and Yb) and both pairs combined, but also considered 13 elements
225 after dimensional reduction by trace element correlation, 25 non-fluid-mobile trace elements
226 (i.e., excluding Rb, Sr, Ba, Pb, and U), and all 30 considered trace elements.

227

228 **3.2 Unsupervised machine learning models**

229 We used unsupervised machine learning models to assess the degree of similarity between
230 the data distributions of the Labeled and Application datasets. The input data (all major
231 elements) used in the unsupervised machine learning models was centered log-ratio
232 transformed to prevent data closure (Aitchison, 1982).

233 We applied k -means clustering to the Parent dataset to measure the similarity of the
234 Labeled and Application datasets. k -means clustering is an unsupervised machine learning
235 algorithm that classifies a given dataset into k clusters. It defines k (an *a priori* fixed number)
236 centroids, or mean points, one for each cluster S_i that minimizes a norm of the kind:

$$237 V = \sum_{i=1}^K \sum_{X_j \in S_i} (X_j - \bar{X}_i), \quad (\text{Eq. 6})$$

238 where \bar{X}_i are the mean points of all $X_j \in S_i$, and V is the objective function. In other words,
239 k -means clustering can divide a dataset into several clusters. As a result, data in the same cluster
240 have similar information that is different from data in other clusters. For example, after
241 clustering of the Parent dataset, if the Labeled dataset is distributed in all clusters, then the
242 Labeled dataset contains the same data distribution as that the Application dataset. Therefore,

243 we can utilize this result to predict the Application dataset. To visualize the k -means clustering
244 result, we used principal component analysis (PCA) (Smith, 2002) to reduce the dimensionality
245 to 2 dimensions.

246

247 **3.3 Application of trained models to unlabeled data**

248 Finally, we applied the best-performing supervised training model (see section 3.1) to the
249 Application dataset. For an unlabeled clinopyroxene analysis, the XGBoost model classifies it
250 as “positive” (metasomatized) or “negative” (unmetasomatized) based on its chemical
251 composition. For each xenolith sampling site l , we defined the mean probability of occurrence
252 of metasomatism (\bar{p}_l) by averaging the predictive probabilities obtained in all n analyses at
253 that site as:

$$254 \bar{p}_l = \frac{\sum_{i=1}^n p_{l,i}}{n}. \quad (\text{Eq. 7})$$

255 Based on the predicted probabilities of metasomatism computed for each sampling location,
256 we mapped the mean predicted probability of metasomatism at the global scale within $1^\circ \times 1^\circ$
257 bins, for a total of 599 bins that represent our study area.

258

259 **4. Results and Discussion**

260 **4.1 Limitations of traditional models**

261 Previous studies have proposed that metasomatism be identified based on elemental ratios
262 such as $\text{CaO}/\text{Al}_2\text{O}_3$, $\text{Mg}\#$, Ti/Eu , or $(\text{La}/\text{Yb})_{\text{N}}$ (Brey et al., 2008; Coltorti et al., 1999; Klemme
263 et al., 1995; Rudnick et al., 1993; Yaxley & Green, 1998; Zong & Liu, 2018). Although these
264 parameters effectively characterize metasomatism at specific sampling locations, they perform
265 poorly when used to predict metasomatism at the global scale (Figure 2).

Classification results using elemental ratios are inconsistent with the petrographically confirmed metasomatized natures of the samples. For example, among the 2,089 samples labeled according to the presence of metasomatic minerals, the classification of metasomatized and unmetasomatized xenoliths using clinopyroxene $\text{CaO}/\text{Al}_2\text{O}_3$ has an accuracy rate of 47% and 72%, respectively (Figure 2a). Other elemental ratios similarly show low accuracy in predicting metasomatism (17%, 15%, and 53% for $\text{Mg}\#$, Ti/Eu , and $(\text{La}/\text{Yb})_N$, respectively), and higher accuracy in predicting unmetasomatized observations (69%, 100%, and 100%, respectively; Figure 2b–d). Therefore, it is complicated to predict the occurrence of metasomatism. Classic bi-variate plots are suitable for predicting unmetasomatized samples because the broad data range of those samples overlaps that of metasomatized samples: unmetasomatized samples are typically taken as having $\text{Ca}/\text{Al} < 5$, $\text{Mg}\# < 92$, $\text{Ti}/\text{Eu} > 1500$, or $(\text{La}/\text{Yb})_N < 3$, but these ranges also describe numerous metasomatized samples. Our results show that globally, the accuracy of traditional bi-variate plots is relatively poor at only 43–76.5%.

Identifications of metasomatism using different traditional element ratios are also inconsistent (Figure S3). To demonstrate this, we labeled data with $\text{CaO}/\text{Al}_2\text{O}_3 > 5$ as “positive” and <5 as “negative” (Figure S3a). However, plotting these data points, labeled by their $\text{CaO}/\text{Al}_2\text{O}_3$ ratios, on the MgO vs. FeO^T diagram (and taking $\text{Mg}\# > 92$ as also indicating “positive”) exhibits a mean accuracy of 67% (Figure S3b). Similarly for trace elements, applying labels based on Ti/Eu values to $(\text{La}/\text{Yb})_N$ shows 75% accuracy (Figure S3c, d). Therefore, the traditional elemental ratios $\text{CaO}/\text{Al}_2\text{O}_3$, MgO/FeO^T , Ti/Eu , and $(\text{La}/\text{Yb})_N$ cannot effectively and accurately classify metasomatism across different sampling sites worldwide.

288

289 **4.2 Classification results and geochemical explanation**

290 Figure 4a presents the confusion matrix (see section 3.1) of the classification results based
291 on the XGBoost model trained using data for all nine major elements from the analyses in the
292 Labeled dataset. Based on the confusion matrix, we obtained a F1 score and accuracy of 0.968
293 and 0.949, respectively. Consistently, the results based on all 30 trace elements produce an F1
294 score and accuracy of 0.957 and 0.947, respectively (Figure 4c). These high F1 scores and
295 accuracies suggest that both major and trace element models may outperform traditional ion-
296 pair classification methods.

297 Based on the results of the XGBoost algorithm, we calculated the relative importance of
298 each feature to the metasomatism classification. As shown in Figure 4b, Na_2O , FeO^T , MnO ,
299 and CaO have the highest relative importance scores among the major elements, indicating that
300 they are important for discriminating whether a sample has been metasomatized. Indeed, the
301 presence of melt affects clinopyroxene compositions, producing clinopyroxene with lower
302 $\text{Mg}^\#$ and higher Na/Ca (Yaxley & Green, 1998). Furthermore, Mn's redistribution among
303 garnet, orthopyroxene, clinopyroxene, and olivine is affected by metasomatism (Achterbergh
304 et al., 2001; Norman, 1998). As shown in Figure 4d, Ho, Ce, U, Sr, Yb, and Ba are the most
305 important trace elements for classifying metasomatism, consistent with previous studies
306 evidencing that LREEs preferentially enter the mineral phase compared to HREEs during
307 interactions between peridotite and melts enriched in incompatible elements (Green et al., 1992;
308 Klemme et al., 1995; Rudnick et al., 1993; Sweeney et al., 1995). We note that most elements
309 have positive but small (<17%) feature importance values, suggesting that they may not play a
310 major role in classification of metasomatism, and that metasomatism cannot be effectively
311 identified by using those elements alone.

312

313 **4.3 Feature correlation and selection**

314 The Parent dataset includes 21,605 observations (rows) corresponding to 9 major elements
315 and 2,967 observations corresponding to 30 trace elements, and we calculated Pearson's
316 correlation coefficients (ρ) between major (Figure 5) and trace elements (Figure S4), where ρ
317 = 1 (-1) indicates a perfect positive (negative) correlation, and $\rho = 0$ indicates no correlation.
318 Several major element pairs are highly correlated (e.g., $\rho = -0.81$ for SiO_2 and TiO_2) or
319 moderately correlated (e.g., $\rho = -0.61$ for MgO and TiO_2), but most are poorly correlated ($|\rho|$
320 < 0.40 ; Figure 5). In comparison, less than a quarter of all trace element pairs are highly
321 correlated ($|\rho| > 0.75$; Figure S4).

322 Our PCA results (Figure S5) for the Parent dataset show that only 64% and 56% of the
323 variance in the major and trace element data, respectively, can be explained by two dimensions
324 (Figure S5). Therefore, the correlation matrix (Figures 5 and S4) and PCA results further
325 evidence that most elemental ratios provide distinct information and may independently
326 contribute to identifying metasomatism.

327 In general, XGBoost provides better classification results when it is trained on more
328 elements (Table 1). For example, the respective F1 scores and accuracies of models trained on
329 major element data from the Labeled dataset improved from 0.891–0.899 and 0.821–0.833
330 when only two elements were used to 0.941 and 0.910 for four elements and 0.968 and 0.949
331 for all nine major elements. For models trained on trace element data from the Labeled dataset,
332 the respective F1 scores and accuracies improved from 0.818–0.933 and 0.771–0.916 for two
333 elements to 0.945 and 0.931 for four elements and 0.960 and 0.950 for 13 elements, but do not
334 improve markedly when using 25 (0.954 and 0.943) or 30 elements (0.957 and 0.947). Our
335 results show that XGBoost performs optimally when trained on 13 features (elements) and does
336 not improve when more features are used. Despite that each feature shows a relatively low
337 variable importance value, the best predictive performance is achieved when most features are

338 included. Given the data and within the limitations of the models, our results suggest that each
339 feature may contribute to partially explain metasomatism.

340

341 **4.4 Evaluating ML model performance and applicability**

342 The ten-fold cross-validation procedure we performed on the Labeled dataset before it
343 was randomly split into training (70%) and testing subsets (30%) resulted in a F1 score of 0.871
344 with standard deviation (s.d.) = 0.073 for the major element data and 0.918 (s.d. = 0.122) for
345 the trace element data. These results demonstrate that the Labeled dataset is relatively balanced.

346 In Table 1, the mean F1 score of the best model as determined by Grid Search Cross-
347 Validation on the major element training data was 0.950 (MajorI-9) and that for the trace
348 element data was 0.973 (TraceI-25). The XGBoost models can then be further evaluated by
349 applying these best models to the testing set (30% of the Labeled data). The best major and
350 trace element models achieved accuracies of 0.949 (MajorI-9) and 0.950 (TraceI-13),
351 respectively, when applied to the testing set (Table 1).

352 Unsupervised learning is useful for discerning patterns from the characteristics of the data
353 itself (Figure 6). In our k -means unsupervised model trained on the Parent dataset, the highest
354 silhouette coefficient (a measure of how similar an object is to its cluster compared to other
355 clustered, with high values indicating objects are well-matched to their clusters and poorly
356 matched to neighboring clusters; Rousseeuw, 1987) corresponds to two major element clusters
357 and two trace element clusters. The Labeled dataset is distributed across all clusters in which
358 the Application dataset is distributed, indicating that both the Application and Labeled datasets
359 have similar distributions. These results indicate that the model trained on the Labeled dataset
360 can be confidently applied to the Application dataset.

361

362 **4.5 Probability of mantle metasomatism at the global scale**

363 When applied to the Application dataset, our model predictions are presented as the
364 probability of metasomatism, which ranges from 0 to 1 by definition. The global map of the
365 predicted mean probability of metasomatism identifies locations with high probabilities of
366 metasomatism (Figure 7). Here, we computed the mean for each location because multiple
367 analytical points were available at each location.

368 The map highlights variations in the distribution of the probability of metasomatism.
369 Figure 8 shows the predicted probability distributions at four localities. These locations were
370 chosen because they cover four continents and because a sufficient number of samples (>100)
371 were available at each to accurately estimate metasomatism. The results suggest bimodal
372 distributions at Hannuoba (North China Craton) and Zealandia (South Pacific Ocean), and
373 unimodal distributions with high probabilities of metasomatism at Pulpwood Harbour (South
374 Canadian Shield) and Finsch (Kaapvaal Craton). The variability observed in the results may
375 indicate that mantle metasomatism occurs widely but heterogeneously. Therefore, the
376 probability of metasomatism at the global scale is generally high, and melt heterogeneity may
377 reduce the likelihood of metasomatism in some locations. Alternatively, the machine learning
378 algorithms work well for moderately metasomatized samples from the classic stable cratons
379 (e.g., South Canadian Shield, Kaapvaal Craton), yet misclassify extensively metasomatized
380 samples affected by multiple metasomatic agents that first fertilized and later depleted mineral
381 chemical compositions (Zhang, 2009).

382 In addition, we also compared the probability of metasomatism to xenolith rock type. The
383 results show that no correlation exists between the probability of metasomatism and rock type,
384 including clinopyroxenite, dunite, harzburgite, lherzolite, peridotite, pyroxenite, and wehrlite
385 (Figure S6). Indeed, it has been suggested that metasomatism may occur in various tectonic
386 settings (Aiuppa et al., 2021; Dawson, 1984; Liu et al., 2021; Menzies & Murthy, 1980; Roden
387 & Murthy, 1985; Wang et al., 2022). In particular, carbon and water lower the melting

388 temperatures of peridotites, and carbonated and hydrous silicate melts have been suggested as
389 effective metasomatic agents (Hirschmann, 2000; Dasgupta and Hirschmann, 2006; Sarafian
390 et al., 2017; Sun and Dasgupta, 2019; Thomson et al., 2016). However, considering only mantle
391 xenoliths might present a bias because they are preferentially affected by melt when brought
392 to the surface by eruptions, but cannot represent the average mantle composition (Artemieva,
393 2009).

394 To assess whether the probability of metasomatism is related to local lithospheric
395 structures, we compared our results with geophysical observations of crustal thickness,
396 lithospheric thickness, and *S*-wave velocity (at 50–200 km depth in 25-km depth intervals) and
397 parameterized the globe into a $1^\circ \times 1^\circ$ grid (Figure S7). Within each cell, we averaged the
398 probabilities of metasomatism for each location and compared them with geophysical
399 observations (Figures S7, S8). We also used unsupervised machine learning to search for
400 clustering of metasomatism probabilities and geophysical parameters, but did not observe any
401 correlations (Figure S8). We identified three possible reasons for this. First, mantle
402 metasomatism may not necessarily be related to specific tectonic settings. Second,
403 metasomatism produces only secondary effects on geophysical parameters such as seismic
404 wave velocity, and a full separation of compositional from thermal factors is required to
405 identify potential metasomatic modifications to the lithospheric mantle. Third, our dataset does
406 not provide information on the depth or age distributions of the clinopyroxenes, making it
407 difficult to relate the predicted metasomatic probabilities to the lithospheric mantle at a specific
408 spatiotemporal location. Therefore, further efforts are required to reconcile the effects of
409 metasomatism on both the chemical and physical properties of the lithospheric mantle at the
410 global scale.

411

412 **5. Conclusions**

413 We developed a model to predict whether xenolithic clinopyroxenes have been
414 metasomatized by applying a multidimensional approach using the XGBoost machine learning
415 algorithm. Our model can predict whether a given sample has been metasomatized with better
416 accuracy (95%) than traditional approaches using elemental ratios (43–77%). Our results
417 indicate that models trained on clinopyroxene compositions, including all major and at least 13
418 trace elements, achieve the best prediction accuracy compared to traditional methods using
419 only two or four elements. Furthermore, *k*-means clustering showed that the Application and
420 Labeled datasets had similar data distributions, indicating that the models trained on the
421 Labeled data can confidently be used to predict whether clinopyroxenes experienced
422 metasomatism based on unlabeled data. Finally, our results show that many locations are likely
423 to have undergone metasomatism and that metasomatism is heterogeneously distributed
424 worldwide.

425

426 **Data Availability Statement**

427 All data and code used in this study are available at <https://doi.org/10.5281/zenodo.6466993>.

428

429 **Acknowledgments**

430 The authors appreciate the assistance of Yifan Zhang, Yunzhu He, and Jinfeng Sun for data
431 compilation. We thank Can He, Shengxin Wang, and Anzhou Li in for their contribution in
432 running the models. The work benefits from discussions with Prof. Hongfu Zhang, Xinmiao
433 Zhao, and Yan Xiao. We thank Robert Dennen for polishing the language of the paper. J
434 ZhangZhou acknowledges support from NSFC grant No. 42072066 and startup funds from
435 Zhejiang University. Andre Python was funded by Zhejiang University under grant No.
436 2021QN81029 (fundamental research funds for the central universities).

437

438 **References**

439 Achterbergh, E. V., Griffin, W. L., & Stienhofer, J. (2001). Metasomatism in mantle
440 xenoliths from the Letlhakane kimberlites: estimation of element fluxes. *Contributions to
441 Mineralogy and Petrology*, 141(4), 397-414. <https://doi.org/10.1007/s004100000236>

442 Aiuppa, A., Casetta, F., Coltorti, M., Stagno, V., & Tamburello, G. (2021). Carbon
443 concentration increases with depth of melting in Earth's upper mantle. *Nature
444 Geoscience*, 14(9), 697-703. <https://doi.org/10.1038/s41561-021-00797-y>

445 Aitchison, J. (1982). The statistical analysis of compositional data. *Journal of the Royal
446 Statistical Society: Series B (Methodological)*, 44(2), 139-160.
447 <https://doi.org/10.1111/j.2517-6161.1982.tb01195.x>

448 Araújo, D. P., Griffin, W. L., & O'Reilly, S. Y. (2009). Mantle melts, metasomatism and
449 diamond formation: Insights from melt inclusions in xenoliths from Diavik, Slave
450 Craton. *Lithos*, 112, 675-682. <https://doi.org/10.1016/j.lithos.2008.09.015>

451 Artemieva, I. M. (2009). The continental lithosphere: reconciling thermal, seismic, and
452 petrologic data. *Lithos*, 109(1-2), 23-46. <https://doi.org/10.1016/j.lithos.2008.09.015>

453 Azodi, C. B., Tang, J., & Shiu, S. H. (2020). Opening the Black Box: Interpretable machine
454 learning for geneticists. *Trends in Genetics*, 36(6), 442-455.
455 <https://doi.org/10.1016/j.tig.2020.03.005>

456 Boser, B. E., Guyon, I. M. & Vapnik, V. N. (1992). A training algorithm for optimal margin
457 classifiers. In Haussler, D. (Ed.), *5th Annual ACM Workshop on COLT* (pp. 144–152).
458 Pittsburgh, PA: ACM Press.

459 Breiman, L. (2001). Random Forests. *Machine Learning*, 45, 5–32.
460 <https://doi.org/10.1023/A:1010933404324>

461 Brey, G. P., Bulatov, V. K., Girnis, A. V., & Lahaye, Y. (2008). Experimental melting of
462 carbonated peridotite at 6–10 GPa. *Journal of Petrology*, 49(4), 797-821.
463 <https://doi.org/10.1093/petrology/egn002>

464 Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In *KDD '16:*
465 *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery*
466 *and Data Mining* (pp. 785-794). New York, NY: Association for Computing Machinery.
467 <https://doi.org/10.1145/2939672.2939785>

468 Coltorti, M., Bonadiman, C., Hinton, R. W., Siena, F., & Upton, B. G. J. (1999). Carbonatite
469 metasomatism of the oceanic upper mantle: evidence from clinopyroxenes and glasses in
470 ultramafic xenoliths of Grande Comore, Indian Ocean. *Journal of Petrology*, 40(1), 133-165.
471 <https://doi.org/10.1093/petroj/40.1.133>

472 Cook-Patton, S. C., Leavitt, S. M., Gibbs, D., Harris, N. L., Lister, K., Anderson-Teixeira, K.
473 J., et al. (2020). Mapping carbon accumulation potential from global natural forest regrowth.
474 *Nature* 585, 545–550. <https://doi.org/10.1038/s41586-020-2686-x>

475 Dalou, C., Koga, K. T., Hammouda, T., & Poitrasson, F. (2009). Trace element partitioning
476 between carbonatitic melts and mantle transition zone minerals: Implications for the source
477 of carbonatites. *Geochimica et Cosmochimica Acta*, 73(1), 239-255.
478 <https://doi.org/10.1016/j.gca.2008.09.02>

479 Dasgupta, R., & Hirschmann, M. M. (2006). Melting in the Earth's deep upper mantle caused
480 by carbon dioxide. *Nature*, 440(7084), 659-662. <https://doi.org/10.1038/nature04612>

481 Dawson, J. B. (1984). Contrasting types of upper-mantle metasomatism? *Developments in*
482 *Petrology*, 11(2), 289-294. <https://doi.org/10.1016/B978-0-444-42274-3.50030-5>

483 Dice, L. R. (1945). Measures of the amount of ecologic association between species. *Ecology*,
484 26(3), 297-302. <https://doi.org/10.2307/1932409>

485 Frey, F. A., & Green, D. H. (1974). The mineralogy, geochemistry and origin of Iherzolite
486 inclusions in Victorian basanites. *Geochimica et Cosmochimica Acta*, 38(7), 1023-1059.
487 [https://doi.org/10.1016/0016-7037\(74\)90003-9](https://doi.org/10.1016/0016-7037(74)90003-9)

488 Frey, F. A., & Prinz, M. (1978). Ultramafic inclusions from San Carlos, Arizona: petrologic
489 and geochemical data bearing on their petrogenesis. *Earth and Planetary Science
490 Letters*, 38(1), 129-176. [https://doi.org/10.1016/0012-821X\(78\)90130-9](https://doi.org/10.1016/0012-821X(78)90130-9)

491 Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals
492 of Statistics*, 29(5), 1189-1232. <https://www.jstor.org/stable/2699986>

493 Green, T. H., Adam, J., & Siel, S. H. (1992). Trace element partitioning between silicate
494 minerals and carbonatite at 25 kbar and application to mantle metasomatism. *Mineralogy
495 and Petrology*, 46(3), 179-184. <https://doi.org/10.1007/BF01164645>

496 Griffin, W. L., Kobussen, A. F., Babu, E. V. S. S. K., O'Reilly, S. Y., Norris, R., & Sengupta,
497 P. (2009). A translithospheric suture in the vanished 1-Ga lithospheric root of South India:
498 evidence from contrasting lithosphere sections in the Dharwar Craton. *Lithos*, 112, 1109-
499 1119. <https://doi.org/10.1016/j.lithos.2009.05.015>

500 Hirschmann, M. M. (2000). Mantle solidus: Experimental constraints and the effects of
501 peridotite composition. *Geochemistry, Geophysics, Geosystems*, 1(10).
502 <https://doi.org/10.1029/2000GC000070>

503 Klemme, S., Van der Laan, S. R., Foley, S. F., & Günther, D. (1995). Experimentally
504 determined trace and minor element partitioning between clinopyroxene and carbonatite
505 melt under upper mantle conditions. *Earth and Planetary Science Letters*, 133(3-4), 439-
506 448. [https://doi.org/10.1016/0012-821X\(95\)00098-W](https://doi.org/10.1016/0012-821X(95)00098-W)

507 Kohavi, R. (1995). A study of cross-validation and bootstrap for accuracy estimation and model
508 selection. In *Proceedings of the Fourteenth International Joint Conference on Artificial
509 Intelligence* (Vol. 2, pp. 1137-1145).

510 Liu, J., Pearson, D. G., Wang, L. H., Mather, K. A., Kjarsgaard, B. A., Schaeffer, A. J., et al.
511 (2021). Plume-driven retronization of deep continental lithospheric
512 mantle. *Nature*, 592(7856), 732-736. <https://doi.org/10.1038/s41586-021-03395-5>

513 Lloyd, F. E., & Bailey, D. (1975). Light element metasomatism of the continental mantle: the
514 evidence and the consequences. *Physics and Chemistry of the Earth*, 9, 389-416.
515 [https://doi.org/10.1016/0079-1946\(75\)90030-0](https://doi.org/10.1016/0079-1946(75)90030-0)

516 McDonough, W. F., & Sun, S. S. (1995). The composition of the Earth. *Chemical*
517 *Geology*, 120(3-4), 223-253. [https://doi.org/10.1016/0009-2541\(94\)00140-4](https://doi.org/10.1016/0009-2541(94)00140-4)

518 Menzies, M., & Murthy, V. R. (1980). Nd and Sr isotope geochemistry of hydrous mantle
519 nodules and their host alkali basalts: implications for local heterogeneities in
520 metasomatically veined mantle. *Earth and Planetary Science Letters*, 46(3), 323-334.
521 [https://doi.org/10.1016/0012-821X\(80\)90048-5](https://doi.org/10.1016/0012-821X(80)90048-5)

522 Norman, M. D. (1998). Melting and metasomatism in the continental lithosphere: laser ablation
523 ICPMS analysis of minerals in spinel lherzolites from eastern Australia. *Contributions to*
524 *Mineralogy and Petrology*, 130(3), 240-255. <https://doi.org/10.1007/s004100050363>

525 O'Reilly, S. Y., & Griffin, W. L. (2013). Mantle metasomatism. In Harlov, D. E., & Austrheim,
526 H. (Eds.), *Metasomatism and the Chemical Transformation of Rock, Lecture Notes in Earth*
527 *System Sciences* (pp. 471-533). Springer, Berlin, Heidelberg.
528 <https://link.springer.com/book/10.1007%2F978-3-642-28394-9>

529 Pearson, D. G., Scott, J. M., Liu, J., Schaeffer, A., Wang, L. H., van Hunen, J., et al. (2021).
530 Deep continental roots and cratons. *Nature*, 596(7871), 199-210.
531 <https://doi.org/10.1038/s41586-021-03600-5>

532 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., et al. (2011).
533 Scikit-learn: Machine learning in Python. *The Journal of Machine Learning Research*, 12,
534 2825-2830. <https://arxiv.org/abs/1201.0490v4>

535 Peng, Y., Manthilake, G., & Mookherjee, M. (2021). Electrical conductivity of metasomatized
536 lithology in subcontinental lithosphere. *American Mineralogist*.
537 <https://doi.org/10.2138/am-2021-7942>

538 Petrelli, M., & Perugini, D. (2016). Solving petrological problems through machine learning:
539 the study case of tectonic discrimination using geochemical and isotopic
540 data. *Contributions to Mineralogy and Petrology*, 171(10), 1-15.
541 <https://doi.org/10.1007/s00410-016-1292-2>

542 Petrelli, M., Caricchi, L., & Perugini, D. (2020). Machine Learning Thermo-Barometry:
543 Application to Clinopyroxene-Bearing Magmas. *Journal of Geophysical Research: Solid*
544 *Earth*, 125(9), e2020JB020130. <https://doi.org/10.1029/2020JB020130>

545 Python, A., Bender, A., Nandi, A. K., Hancock, P. A., Arambepola, R., Brandsch, J., & Lucas,
546 T. C. (2021). Predicting non-state terrorism worldwide. *Science Advances*, 7(31), eabg4778.
547 <https://doi.org/10.1126/sciadv.abg4778>

548 Ringwood, A. E. (1962). A model for the upper mantle. *Journal of Geophysical*
549 *Research*, 67(2), 857-867. <https://doi.org/10.1029/JZ067i002p00857>

550 Roden, M. F., & Murthy, V. R. (1985). Mantle metasomatism. *Annual Review of Earth and*
551 *Planetary Sciences*, 13(1), 269-296. <https://doi.org/10.1146/annurev.ea.13.050185.001413>

552 Rousseeuw, P. J. (1987). Silhouettes: a graphical aid to the interpretation and validation of
553 cluster analysis. *Journal of Computational and Applied Mathematics*, 20, 53-65.
554 [https://doi.org/10.1016/0377-0427\(87\)90125-7](https://doi.org/10.1016/0377-0427(87)90125-7)

555 Rudnick, R. L., McDonough, W. F., & Chappell, B. W. (1993). Carbonatite metasomatism in
556 the northern Tanzanian mantle: petrographic and geochemical characteristics. *Earth and*
557 *Planetary Science Letters*, 114(4), 463-475. [https://doi.org/10.1016/0012-821X\(93\)90076-](https://doi.org/10.1016/0012-821X(93)90076-)
558 L

559 Sarafian, E., Gaetani, G. A., Hauri, E. H., & Sarafian, A. R. (2017). Experimental constraints
560 on the damp peridotite solidus and oceanic mantle potential
561 temperature. *Science*, 355(6328), 942-945.

562 <https://www.science.org/doi/10.1126/science.aaj2165>

563 Sarbas, B. (2008). The GEOROC database as part of a growing geoinformatics network. In
564 Brady, S. R., Sinha, A. K., Gundersen, L. C. (Eds.), *Geoinformatics 2008—Data to
565 Knowledge, Proceedings, Geoinformatics 2008—Data to Knowledge* (Scientific
566 Investigations Report 2008-5172, pp. 42-43). Reston, VA: USGS. https://gfzpublic.gfz-potsdam.de/pubman/item/item_10062

568 Smith, L. I. (2002). *A tutorial on Principal Components Analysis* (Technical report). University
569 of Otago. <http://hdl.handle.net/10523/7534>

570 Sørensen, T. J. (1948). *A method of establishing groups of equal amplitude in plant sociology
571 based on similarity of species content and its application to analyses of the vegetation on
572 Danish commons*. Copenhagen: I kommission hos E. Munksgaard.

573 Stehman, S. V. (1997). Selecting and interpreting measures of thematic classification accuracy.
574 *Remote Sensing of Environment*, 62(1), 77-89.
575 <https://www.sciencedirect.com/science/article/abs/pii/S0034425797000837>

576 Sun, C., & Dasgupta, R. (2019). Slab–mantle interaction, carbon transport, and kimberlite
577 generation in the deep upper mantle. *Earth and Planetary Science Letters*, 506, 38-52.
578 <https://doi.org/10.1016/j.epsl.2018.10.028>

579 Sweeney, R. J., Prozesky, V., & Przybylowicz, W. (1995). Selected trace and minor element
580 partitioning between peridotite minerals and carbonatite melts at 18–46 kb
581 pressure. *Geochimica et Cosmochimica Acta*, 59(18), 3671-3683.
582 [https://doi.org/10.1016/0016-7037\(95\)00270-A](https://doi.org/10.1016/0016-7037(95)00270-A)

583 Thomson, A. R., Walter, M. J., Kohn, S. C., & Brooker, R. A. (2016). Slab melting as a barrier
584 to deep carbon subduction. *Nature*, 529(7584), 76-79.
585 <https://www.nature.com/articles/nature16174>

586 Thomson, A. R., Kohn, S. C., Prabhu, A., & Walter, M. J. (2021). Evaluating the formation
587 pressure of diamond-hosted majoritic garnets: A machine learning majorite
588 barometer. *Journal of Geophysical Research: Solid Earth*, 126(3), e2020JB020604.
589 <https://doi.org/10.1029/2020JB020604>

590 Ueki, K., Hino, H., & Kuwatani, T. (2018). Geochemical discrimination and characteristics of
591 magmatic tectonic settings: A machine-learning-based approach. *Geochemistry,
592 Geophysics, Geosystems*, 19(4), 1327-1347. <https://doi.org/10.1029/2017GC007401>

593 Valetich, M. J., Le Losq, C., Arculus, R. J., Umino, S., & Mavrogenes, J. (2021). Compositions
594 and Classification of Fractionated Boninite Series Melts from the Izu–Bonin–Mariana Arc:
595 A Machine Learning Approach. *Journal of Petrology*, 62(2), egab013.
596 <https://doi.org/10.1093/petrology/egab013>

597 Wang, C., Jin, Z., Gao, S., Zhang, J., & Zheng, S. (2010). Eclogite-melt/peridotite reaction:
598 Experimental constrains on the destruction mechanism of the North China Craton. *Science
599 China Earth Sciences*, 53(6), 797-809. <https://doi.org/10.1007/s11430-010-3084-2>

600 Wang, X., Wang, Z., Cheng, H., Zong, K., Wang, C. Y., Ma, L., et al. (2022). Gold endowment
601 of the metasomatized lithospheric mantle for giant gold deposits: Insights from
602 lamprophyre dykes. *Geochimica et Cosmochimica Acta*, 316, 21-40.
603 <https://doi.org/10.1016/j.gca.2021.10.006>

604 Wang, Y., Qiu, K. F., Müller, A., Hou, Z. L., Zhu, Z. H., & Yu, H. C. (2021). Machine Learning
605 Prediction of Quartz Forming-Environments. *Journal of Geophysical Research: Solid
606 Earth*, 126(8), e2021JB021925. <https://doi.org/10.1029/2021JB021925>

607 Yaxley, G. M., & Green, D. H. (1998). Reactions between eclogite and peridotite: mantle
608 refertilisation by subduction of oceanic crust. *Schweizerische mineralogische und*
609 *petrographische Mitteilungen*, 78(2), 243-255.

610 Zhang, H. F. (2009). Peridotite-melt interaction: a key point for the destruction of cratonic
611 lithospheric mantle. *Chinese Science Bulletin*, 54, 3417-3437.
612 <https://doi.org/10.1007/s11434-009-0307-z>

613 Zhang, H. F., Goldstein, S. L., Zhou, X. H., Sun, M., & Cai, Y. (2009). Comprehensive
614 refertilization of lithospheric mantle beneath the North China Craton: further Os–Sr–Nd
615 isotopic constraints. *Journal of the Geological Society*, 166(2), 249-259.
616 <https://doi.org/10.1144/0016-76492007-152>

617 Zhao, Y., Zhang, Y., Geng, M., Jiang, J., & Zou, X. (2019). Involvement of slab-derived fluid
618 in the generation of Cenozoic basalts in Northeast China inferred from machine
619 learning. *Geophysical Research Letters*, 46(10), 5234-5242.
620 <https://doi.org/10.1029/2019GL082322>

621 Zheng, Z., Zhao, L. & Oleson, K.W. Large model structural uncertainty in global projections
622 of urban heat waves. *Nature Communications* 12, 3736 (2021).
623 <https://doi.org/10.1038/s41467-021-24113-9>

624 Zong, K., & Liu, Y. (2018). Carbonate metasomatism in the lithospheric mantle: Implications
625 for cratonic destruction in North China. *Science China Earth Sciences*, 61(6), 711-729.
626 <https://doi.org/10.1007/s11430-017-9185-2>

627

628 **Figure Captions**

629 **Figure 1.** Locations of sample analyses used in this study. The color of each sampling point
630 represents the number of analyses performed on clinopyroxenes in mantle xenoliths from that
631 location.

632 **Figure 2.** The application of elemental ratios proposed in previous studies to attempt to identify
633 metasomatism in the global dataset. Symbols indicate whether each sample was
634 petrographically identified as metasomatized ('positive') or not affected by metasomatism
635 ('negative'). (a) The accuracy (Eq. 2) of CaO versus Al₂O₃ is 59.5%; (b) that of MgO versus
636 FeO^T is 43%; (c) that of Ti versus Eu is 57.5%; and (d) that of La versus Yb is 76.5%.

637 **Figure 3.** Operational flow chart of the methods used in this study. Step I: the Labeled dataset
638 was used to train the XGBoost models and evaluate model performance. Step II: the
639 preprocessed Parent dataset was used to train *k*-means clustering models to verify that the data
640 distributions of the Labeled and Application datasets were similar. Step III: the best model was
641 used to predict the probability of metasomatism worldwide within 1° × 1° grid cells based on
642 the Application dataset.

643 **Figure 4.** Results of the XGBoost model trained on the Labeled dataset to classify
644 clinopyroxenes as affected or unaffected by metasomatism. (a, c) Confusion matrices of
645 classification results based on major and trace element compositions from the testing subset,
646 respectively. (b, d) The relative feature importances of major and trace elements, respectively.

647 **Figure 5.** Heat-map matrix of linear correlations (Pearson coefficients) between major
648 elements concentrations in clinopyroxenes of the Parent dataset (21,605 observations).

649 **Figure 6.** Unsupervised learning results illustrating the similarity of the (a) major and (b) trace
650 element data distributions in the Labeled (training, orange diamonds) and Application datasets
651 (gray circles).

652 **Figure 7.** Probability map of mantle metasomatism at 972 unique sampling locations. Symbol
653 color indicates the predicted probability of metasomatism from 0 (blue) to 1 (red).

654 **Figure 8.** Probability distributions of metasomatism at four selected sampling locations: (a)
655 Hannuoba, North China Craton; (b) Zealandia, South Pacific Ocean; (c) Pulpwood Harbour,
656 South Canadian Shield; and (d) Finsch, Kaapvaal Craton.

657 **Table 1.** Summary of XGBoost model performance.