

1 **Predicting olivine composition using Raman spectroscopy**
2 **through band shift and multivariate analyses**

3 **LAURA B. BREITENFELD,¹ M. DARBY DYAR,¹ CJ CAREY,² THOMAS J. TAGUE, JR.,³ PENG**
4 **WANG³, TERRY MULLEN⁴, AND MARIO PARENTE⁴**

5 ¹Department of Astronomy, Mount Holyoke College, South Hadley, MA 01075, U.S.A.

6 ²College of Information and Computer Sciences, University of Massachusetts Amherst, Amherst
7 MA 01003, U.S.A.

8 ³Bruker Optics, Inc., Billerica, MA 01821, U.S.A.

9 ⁴Department of Electrical and Computer Engineering, University of Massachusetts Amherst,
10 Amherst MA 01003, U.S.A.

11
12
13
14
15 Revision #2
16
17

18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40

ABSTRACT

Olivine group minerals are ubiquitous in extrusive igneous rocks, and play an important role in constraining equilibria for samples in the upper mantle and above. All Raman spectra of the olivine group minerals in the solid solution between forsterite (Fo, Mg_2SiO_4) and fayalite (Fa, Fe_2SiO_4) have a high intensity doublet between 800 and 880 cm^{-1} . Previous studies used small sample suites with limited compositional ranges and varying spectrometers to relate energy shifts of these two bands to Mg/Fe contents. In this work, Raman spectra of 93 olivine samples were acquired on either Bruker's 532 nm (laser wavelength) Senterra or BRAVO (785/852.3 nm) spectrometer. This paper compares the two-peak band shift univariate method with two multivariate methods: partial least squares (PLS) and the least absolute shrinkage operator (Lasso). Datasets from several instruments are also examined to assess the most accurate method for predicting olivine composition from a Raman spectrum.

Our 181-spectra PLS model is recommended for use when determining olivine composition from a Raman spectrum. For Raman spectra of mixed phases where only the olivine doublet can be identified, composition can best be determined using the position of the peak ca. 838-857 cm^{-1} through use of the equation $\% \text{Fo} = -0.179625x^2 + 310.077x - 133717$ regression equation (where x = DB2 centroid in units of cm^{-1}).

In situ methods for predicting mineral composition on planetary surfaces are critically important to extraterrestrial exploration going forward; of these, Raman spectroscopy is likely the best, as evidenced by the impending deployment of several Raman instruments to Mars (*ExoMars* and *Mars 2020*). More broadly, application of machine learning methods to spectral data processing have implications to multiple fields that use spectroscopic data.

Keywords: Raman spectroscopy, olivine, forsterite, fayalite, PLS, Lasso

41

INTRODUCTION

42

43

44

45

46

47

48

49

50

51

Olivine group minerals control many of the properties of Earth's upper mantle, affect rheology, and may be diagnostic of crystallization temperature in terrestrial and extraterrestrial rocks. Their solid solution commonly spans the range between forsterite (Fo, or Mg_2SiO_4) and fayalite (Fa, or Fe_2SiO_4) with minor substitutions of alternative cations such as Mn and Ni. Because olivine composition provides an important petrogenetic indicator, development of convenient methods to measure it without microanalytical techniques that require sample preparation are desirable. This paper focuses on determination of olivine composition using Raman spectroscopy for this purpose. It has the potential to enable compositions to be conveniently determined in the laboratory, in field identifications with portable units, and on other planets such as Mars.

52

53

54

55

56

57

58

59

60

61

62

63

This problem has been extensively studied using conventional regression-based approaches, but generalization of their results is arguable given their very small (<20) sample suites and coverage of olivine composition. These prior studies (Kuebler et al. 2006, Foster et al. 2007, Gaisler and Kolesov 2007, Mouri and Enami, 2008, Yasuzuka et al. 2009, Ishibashi et al. 2011) have related olivine composition to the peak positions of a high intensity doublet in the range of $800\text{-}880\text{ cm}^{-1}$ (Figure 1). Peak centroids are regressed against composition to obtain an equation suitable for prediction of composition in unseen samples. Previous workers have used R^2 values to characterize their prediction algorithms, making their accuracy difficult to quantify and comparisons across models impossible. Moreover, a variety of Raman instruments with varying resolution and laser wavelength, and a combinations of single crystal and powdered samples, were used. It is thus difficult to assess which model to use to estimate olivine composition on unseen data from a different spectrometer than the one used in each study.

64 In this study, all known publicly-accessible olivine Raman data are considered. In
65 addition, new Raman spectra were acquired from a suite of 93 well-characterized synthetic and
66 naturally-occurring olivines using Bruker BRAVO and Senterra spectrometers. The accuracies of
67 linear regression (univariate) methods for various datasets are quantified to enable useful
68 comparisons. Univariate methods are compared and contrasted with two multivariate analysis
69 approaches: partial least squares (Stone and Brooks 1990) and the least absolute shrinkage
70 operator (Hastie et al. 2009) to evaluate the best prediction models for use with Raman spectra of
71 olivines to determine composition with known accuracy.

72 BACKGROUND

73 Assignments of Raman active forsterite and fayalite modes have evolved over time
74 (Table 1). Generally, forsterite Raman bands above 500 cm^{-1} can be classified as internal
75 movement within the $(\text{SiO}_4)^{4-}$ tetrahedra. Below this threshold energy, peaks are caused by
76 rotation and translation of the tetrahedra as well as divalent cation motion. Forsterite and fayalite
77 have 84 vibrational modes; only 36 are Raman active ($11A_g + 11 B_{1g} + 7B_{2g} + 7B_{3g}$) (Mckeown
78 et al. 2010).

79 Key to this study are the two principle Raman bands that form a doublet composed of
80 five vibrational modes ($2A_g + 2B_{1g} + B_{2g}$) (Table 1). This doublet occurs between $\sim 815\text{-}825\text{ cm}^{-1}$
81 (DB1) and $\sim 838\text{-}857\text{ cm}^{-1}$ (DB2) (Kuebler et al. 2006) and is primarily attributed to A_g , though
82 B_{1g} and B_{2g} also affects the shape and intensity of the spectrum. The energy shift of the A_g
83 stretch from the SiO_4 tetrahedra is caused by changes in site geometry due to cation substitutions
84 in adjacent sites. Cation substitutions between forsterite and fayalite thus result in band shifts
85 (Figure 2) as well as changes in the shape and intensity of the peaks. Many previous workers
86 (Table 2) have used the peak centroids of the DB1 and DB2 doublet peaks to derive olivine

87 composition. However, this practice does not allow other information in the spectra to be
88 utilized, such as shifts arising from minor modes that affect the shape of the primary doublet and
89 give rise to other, more subtle features elsewhere in the wavenumber range.

90 Other peaks within olivine spectra have been utilized rarely for prediction of
91 composition, such as $\sim 200\text{-}230\text{ cm}^{-1}$, $\sim 290\text{-}310\text{ cm}^{-1}$, $\sim 410\text{-}440\text{ cm}^{-1}$, $\sim 540\text{-}553\text{ cm}^{-1}$, $\sim 881\text{-}883$
92 cm^{-1} , $\sim 914\text{-}920\text{ cm}^{-1}$, $950\text{-}966\text{ cm}^{-1}$ (Table 2). However, these features are relatively low in
93 intensity compared to those of the DB1 and DB2 doublet, making fitting of peak centroids
94 difficult and less accurate. Raman bands caused by different vibrational modes should not be
95 affected by octahedral substitutions. For example, features between 400 and 700 cm^{-1} have been
96 attributed to the internal bending modes of the anion, which have minimal centroid shifts
97 (Chopelas, 1991; Kuebler et al., 2006).

98 This study evaluates the relative usefulness of the most prominent bands in the Raman
99 spectra of olivine group minerals using a combination of conventional regression/peak fitting
100 methods and more recently-developed multivariate methods. The latter have the advantage of
101 weighing the relative importance of different spectral energies in determining Fe/Mg ratio,
102 enhancing our understanding of the underlying physical processes that give rise to the features.

103 **METHODS**

104 **Sample provenance**

105 Natural samples (Table 3) came from collections of the Mineral Spectroscopy Lab at
106 Mount Holyoke, the National Museum of Natural History (NMNH, Smithsonian), and from S.A.
107 Morse (University of Massachusetts Amherst) (Morse 2001). This is the largest suite of naturally
108 occurring olivine samples studied by Raman (or any other type of) spectroscopy. Roughly one-
109 third of the natural samples came from previous studies of olivines from mantle xenoliths

110 (McGuire et al. 1991, Dyar et al. 1989, Dyar et al. 1992) or Fe³⁺-bearing samples studied by
111 Schaefer (1983), Banfield et al. (1992), and Dyar et al. (1998). Another group of samples was
112 provided by S.A. Morse of the University of Massachusetts Amherst. They come from the
113 Kiglapait layered mafic body, a large 1.3 Ga layered intrusion on the coast of Labrador, Canada
114 (Morse 1996, Morse 2001). As the original melt crystallized, the Fe/Mg ratio of the remaining
115 liquid changed, so a range of olivine compositions were produced. Lower Mg and higher Fe
116 contents occurred successively higher within the intrusion. Finally, several samples came from
117 the NMNH (see Table 3).

118 Naturally occurring olivine typically has high Fo content of roughly 89.5%. For a solid
119 %Fo prediction model, wide representation of the Fo-Fa continuum is needed. Because samples
120 with intermediate Fo/Fa content are rare in nature (except at specific localities such as the
121 Kiglapait, as noted above), synthetic samples were added to our collection of naturally formed
122 olivines to represent %Fo from 0 to 100 (see Dyar et al. 2009 for sample descriptions). Synthetic
123 samples were synthesized by Donald Lindsley in his laboratory at SUNY Stony Brook. First, a
124 silicon and hematite mixture was ground for 1-2 hours with ethanol. Next, an iron sponge was
125 added and for less than one hour grinding continued. The product was enveloped in silver foil
126 and put in a glass silicon capsule. The center of the capsule was drawn out into a capillary while
127 one end of the capsule was sealed, leaving the sample by the sealed end. Near the open end of
128 the capsule, and Fe getter was placed. For 10-20 minutes, the capsule was placed into a ~800°C
129 vertical tube furnace (the Fe getter remained at ~600°C). Finally, the capsule was removed from
130 the furnace and sealed across the capillary. The completely sealed capsule section, which
131 contained the sample, was next cooked for 10 days in a horizontal tube furnace at ~920-940°C
132 (Sklute, 2006). This sample suite has been studied with a wide range of other spectroscopic

133 techniques (Dyar et al. 2009, Lane et al. 2011, Isaacson et al. 2014).

134 **Sample characterization**

135 Olivines examined were either a single crystal or powdered samples. To produce a
136 powdered sample, each sample was first visually inspected and handpicked for purity. Then each
137 grain was treated using oxalic acid (2 tsp. in 2 gal. of water) for one hour to remove surface
138 weathering, followed by three cycles of washing and rinsing with clean water. As needed,
139 samples were either crushed in a tungsten shatterbox or ground by hand in a diamonite mortar.
140 Because crystal orientation affects the Raman spectrum, we chose to study both single crystals
141 and powders (Price et al. 1987), affording the opportunity to compare those results. The spot
142 sizes of the Senterra and BRAVO spectrometers differ, so that single crystals were analyzed with
143 the Senterra while the BRAVO examined powders. However, repeated analyses of the same
144 sample on each instrument showed no evidence for heterogeneity, as expected given our careful
145 sample preparation and use of homogeneous starting material. The only difference was a
146 consistent offset (as discussed below) due to differences in calibration.

147 Many natural samples from Dyar's collections already had published compositions
148 (Table 3) that included Mössbauer studies to determine Fe³⁺ contents. Where needed, additional
149 samples were analyzed by Mössbauer spectroscopy using standard methods (Sklute, 2006). Rh
150 was used on a WEB Research Co. model W100 spectrometer equipped with a Janus closed-cycle
151 He refrigerator. Run times ranged from 2-12 hours; results were calibrated against α -Fe foil.
152 Typical count rates were between 500,000 and 900,000 non-resonant counts/hour. Most samples
153 contained no Fe³⁺.

154 As needed, new electron microprobe analyses of 10 spots on each sample were acquired
155 either by Molly McCanta at the University of Tennessee in Knoxville or at Brown University by

156 Joseph Boesenberg; in both cases, standard operation conditions were used. Figure 3 shows the
157 calculated %Fo for each sample that was determined by normalizing the contents to contain only
158 Mg and Fe, as commonly done with the formula $\%Fo = (100 \times Mg) / (Mg + Fe_{total})$, where
159 $Fe_{total} = \Sigma Fe^{2+} + Fe^{3+}$. This represented only a minor adjustment because only very minor
160 substitutions of other cations were observed, as seen in compositions of the natural samples as
161 given in Table S1. Synthetic samples are as-named in Dyar et al. (2012a).

162 **Raman measurements**

163 Spectra of powdered samples were acquired on a BRAVO dual laser (785 and 852.3 nm
164 simultaneous DuoLaser™) system (2.0 cm⁻¹/channel spectral resolution) with three sample scans
165 and an integration time of 10s. Because the BRAVO samples a large area (~2 mm diameter), it
166 required sample masses of >100 mg, which were only available for 25 samples. The remainder
167 of the sample suite (68 samples) was run on a Bruker Senterra spectrometer using the 532 nm
168 laser and a microscope attachment to probe single grains. The Senterra used 10 mW laser power
169 for two sample scans and integrated for 10s, analyzed through a 20× objective. The highest
170 Senterra resolution available of 0.5 cm⁻¹/channel was utilized.

171 The Senterra calibration was performed automatically and was anchored by the NIST
172 standards, acetaminophen and silicon, resulting in a wavelength accuracy of 0.2 cm⁻¹. The
173 photometric accuracy was verified using NIST traceable glass (Allen et al. 2000). The BRAVO
174 wavelength was similarly calibrated with the wavelength accuracy being 1 cm⁻¹ or better.
175 Multiple sample scans were acquired for each sample ensuring reproducibility of the spectral
176 data acquired. Pre-processed data that included dark subtractions and baseline removal were
177 converted from Bruker's Opus format into ascii files and uploaded to the lab web site, currently
178 at nemo.cs.umass.edu:54321.

179 **Data analysis**

180 Because many spectra showed residual features after the baseline was mostly removed by
181 the Opus algorithm, we applied additional baseline removal using the adaptive iteratively
182 reweighted penalized least squares (AirPLS) method (Zhang et al. 2010), which uses the sum of
183 differences between signal and baseline to adjust weights intelligently. Smoothness is the
184 adjustable baseline removal parameter, for which a value of 100 was used. Multiple types of
185 normalization were tested on these data including normalizing to the maximum value (L_{∞} norm),
186 the sum of absolute values (L1 norm), the sum of squared values (L2 norm), and scaling to
187 intensity at several specific energies. Normalization to the maximum peak intensity
188 outperformed all other methods and therefore it was used in subsequent analyses throughout.
189 Normalization was executed to account for arbitrary intensity differences between the two
190 spectrometers. DB1 and DB2 were peak fitted for each spectrum using Gaussian and Lorentzian
191 peak shapes and a method that simply sums all the counts in the region of interest (e.g., 800-880
192 cm^{-1}). Pre-processing of spectra used the superman website nemo.umass.cs.edu:54321 (Carey et
193 al. 2017).

194 Next, PLS and Lasso models were applied for multivariate analysis of %Fo. PLS
195 regresses one response variable (%Fo) against multiple explanatory variables (intensity at each
196 channel of the spectra). PLS predictions utilize every channel of the spectral range, assigning
197 coefficients to every single channel. Because PLS utilizes all available variables (channels) and
198 eliminates multicollinearity (peaks whose intensities are dependent, as is the case for the doublet
199 in the Raman spectra of olivine). This algorithm was created for the analysis of data with high
200 collinear explanatory (p) variables, which are significantly greater in number compared to the
201 observations (N). Therefore, $p \gg N$ (Butler and Denham 2000). PLS can predict multiple

202 dimensional datasets and has been utilized for the specific application of spectroscopy (Wold et
203 al. 1983). This paper utilizes PLS2 (hereafter referred to as PLS) rather than alternative versions.

204 Lasso is a continuous shrinkage, which allows for the production of coefficient values
205 to be reduced even to as small as zero (Tibshirani 1995). This shrinkage is in agreement with the
206 shrinkage parameter t , by shrinking the residual sum of squares based upon the sum of the
207 absolute value of the coefficients. In other words, this method selects a subset of predictors with
208 the strongest effect on the response variable. Unlike PLS, the Lasso produces a sparse models
209 with few coefficients (depending on the value of the α parameter), with most channel intensities
210 set to zero. The relative merits of PLS versus Lasso in spectroscopic methods (e.g., Dyar et al.
211 2012a) are just beginning to be explored and there is as yet no consensus for which method is
212 better; their usefulness appears to be highly variable for each dataset and application-dependent.

213 **Model comparisons**

214 Use of the R^2 parameter to describe the fit of a regression model (here %Fo is the
215 dependent variable and peak centroid is the independent variable) is not helpful for drawing
216 comparisons between different models because R^2 depends on the error associated with each
217 measurement. A more appropriate metric for cross-comparison is root mean square error
218 (RMSE), which calculates the square root of the average difference between predicted and true
219 %Fo. RMSE is useful in this application because it is expressed in the same unit as the
220 measurement – in our case, %Fo. This paper uses RMSE in three different ways. Internal RMSE
221 describes the prediction error of an expression that is created using all the data in the dataset. In
222 other words, if there are 25 samples, the regression expression utilizes all of them. Internal
223 RMSE is useful in comparing one model to another, but inappropriate for evaluating errors on
224 unseen data. In contrast, leave-one-out cross-validated RMSE (LOO RMSE-CV) removes one

225 sample at a time, uses a regression model based on the other $n-1$ samples to predict the n th
226 sample, and then repeats the process n times, where n is the number of samples in the dataset.
227 Thus LOO RMSE-CV gives the best estimate of how the model will perform on unseen data.
228 Finally, RMSE-test is used to describe the RMSE of comparisons between true and predicted
229 values in completely unseen data.

230 UNIVARIATE (PEAK CENTROID) ANALYSIS

231 Univariate methods focus entirely on the two principle Raman bands in the five-mode
232 doublet between $\sim 815-825\text{ cm}^{-1}$ (DB1) and $\sim 838-857\text{ cm}^{-1}$ (DB2), as discussed above and used
233 by prior workers. Peak centroid positions (Tables S2 and S3) of the Raman spectra of our 93
234 synthetic and naturally-occurring olivines (Figure 2) were utilized for univariate predictions,
235 along with data from the RRUFF database and other publications for which data were provided.
236 There were 25 olivines for which there was sufficient sample to make measurements on the
237 Bruker BRAVO instrument, and those 25 samples plus an additional 68 were also run on the
238 Bruker Senterra spectrometer, which has a microbeam to enable analysis of individual grains or
239 small clumps. Different Raman instruments can produce spectra with equivalent bands at slightly
240 different wavenumber positions due to varying calibration protocols. Therefore, all spectra were
241 analyzed as raw data as well as after the BRAVO dataset was aligned to the Senterra data. This
242 was accomplished by aligning corresponding bands within spectra of 25 samples acquired on
243 both the BRAVO and Senterra spectrometers using a method described in Mullen et al. (2018).
244 These raw and aligned data results are compared to fits made to data taken from the RRUFF
245 database (Table S4). Two different sets of data from RRUFF were tested: all 188 spectra of
246 olivine group minerals, and a subset of 156 spectra designated as RRUFF*. The former group
247 includes 32 spectra from the RRUFF site listed as “broad scan with spectral artifacts,” while the

248 other 156 spectra lack that designation. This annotation refers to spectra acquired over a broad
249 energy range versus one with higher resolution. It is important to note that these data do not
250 represent 188 different samples, but in many cases include spectra of the same samples acquired
251 on multiple instruments, with depolarized versus polarized lasers, and on single crystals with
252 varying orientations. Comparisons are also made to RMSE values calculated using peak
253 positions given in papers by Kuebler et al. (2006), Yasuzuka et al. (2009) and Ishibashi et al.
254 (2011). The equations of the first-, second-, and third-order polynomial fits are reported for the
255 DB1 and DB2 in Table S5.

256 Univariate results from second-order polynomial fits to peak position versus %Fo content
257 are summarized in Table 4, which also includes the resolution of the spectra from each dataset
258 along with values for R^2 (coefficient of determination) of the internally cross-validated data, the
259 internal RMSE values and LOO RMSE-CV. Linear, second-, and third-order polynomial fits
260 relating peak centroid position to composition were created for the BRAVO and Senterra data by
261 Breitenfeld (2017). In all cases, second-order polynomial fits to the data produced more accurate
262 RMSE values than linear ones. Third-order polynomials produced identical or slightly better fits
263 than second-order ones, but the improvement was negligible and not significant. Thus results in
264 Table 4 use second-order fits, following the precedent of Kuebler et al. (2006).

265 While the DB1 BRAVO data conspicuously lie on a polynomial curve (Figure 4), results
266 from the Senterra for DB1 and DB2 (both spectrometers) produce trends that are closer to linear.
267 These differences may result from experimental parameters such as variable resolutions,
268 excitation laser wavelength, and detector sensitivities, all of which impact the consistency of
269 these univariate predictions.

270 Figure 5 displays data from sources other than this study, including the RRUFF* data, all
271 RRUFF olivine data and results of Kuebler et al. (2006), Yasuzuka et al. (2009), and Ishibashi et
272 al. (2011). Previous Raman studies of olivine that did not report peak centroids could not be
273 included on this plot, such as those in Wang et al. (2004), Gaisler and Kolesov (2007), and
274 Mouri and Enami (2008). In Figure 5, the Ishibashi et al. (2011) model predicts %Fo most
275 accurately. It is quickly apparent that the models based on RRUFF* data perform poorly in
276 comparison to the other univariate models, perhaps because that database combines spectra
277 acquired on different instrument using several excitation lasers for the model.

278 The usefulness of applying Kuebler et al.'s (2006) models to other datasets was evaluated
279 through the prediction of two aggregate datasets. When the Senterra + BRAVO + RRUFF*
280 (>50%Fo) dataset is predicted using the models of Kuebler et al. (2006), the RMSE-test values
281 for DB1 and DB2 are 9.99 and 8.10 %Fo, respectively. For the aligned Senterra + BRAVO data,
282 the RMSE-test values for DB1 and DB2 are 19.42 and 14.64%Fo, respectively. These RMSE-
283 test values are larger than the internal or LOO RMSE-CV values of the Kuebler et al. (2006)
284 dataset alone (Table 4).

285 **MULTIVARIATE (MACHINE LEARNING) ANALYSIS**

286 In other types of spectroscopy, it has been shown that multivariate predictions using the
287 entire spectrum produce more accurate predictions of composition. This has been demonstrated
288 for laser-induced breakdown spectroscopy (Tucker et al. 2010, Dyar et al. 2016a) and x-ray
289 absorption spectroscopy (Dyar et al. 2012b, Dyar et al. 2016b). Given the similarities between
290 those data and our Raman spectra, it was expected that Raman predictions of olivine composition
291 might follow this trend. Accordingly, both PLS and Lasso regression methods were tested on our
292 datasets using Raman spectra acquired on the BRAVO and Senterra, along with RRUFF and

293 RRUFF* data for which the full spectra are available online. Because these predictions require
294 use of the entire spectrum rather than just peak centroids, no other publicly accessible olivine
295 Raman spectra could be included in the multivariate analyses. Model comparisons included R^2 ,
296 internal PLS or Lasso RMSE, and LOO RMSE-CV, keeping in mind the caveats just discussed
297 (Table 5).

298 Choice of adjustable parameter is important to the outcomes of both multivariate
299 techniques. For PLS, the number of components in each model were tested using components
300 ranging from 1 to 10. The value producing the lowest (best) prediction accuracy over this range
301 was chosen (first local minimum); generally this value was 4-7 components (Table 5). Numbers
302 of components greater than 10 might produce more accurate predictions of %Fo but they
303 dramatically reduce the generalizability of the model to unseen data, so they were not
304 considered. This is comparable to the concept of using very high-order polynomials to predict
305 data – they can be quite accurate but not applicable to any other datasets.

306 For Lasso models, a value of α was chosen for each prediction to train the model
307 depending on the desired “sparseness” of the model – i.e., how few channels needed to employ
308 to predict %Fo. Variations in the value of α change the number of channels used by the model.
309 As α increases, fewer channels are examined in the multivariate analysis prediction (Figure 6).
310 As the number of channels in a model increases, the value of LOO RMSE-CV also decreases,
311 showing the importance of models with a large number of channels (small α). Use of large
312 numbers of channels may overtune the model and reduce its generalizability.

313 Multivariate models can be set up to use all or any subset of the spectral
314 data/wavenumber range. As discussed above and shown in Table 2, useful wavelength ranges for
315 prediction of olivine composition have been proposed to occur at $\sim 200\text{-}230\text{ cm}^{-1}$, $\sim 290\text{-}310\text{ cm}^{-1}$,

316 $\sim 410\text{-}440\text{ cm}^{-1}$, $\sim 540\text{-}553\text{ cm}^{-1}$, $\sim 881\text{-}883\text{ cm}^{-1}$, $\sim 914\text{-}920\text{ cm}^{-1}$, and $950\text{-}966\text{ cm}^{-1}$. Each of these
317 ranges was tested individually along with models covering all five of those regions as well as the
318 range from $400\text{-}700\text{ cm}^{-1}$ and $300\text{-}1500\text{ cm}^{-1}$. Internal RMSE model accuracies are given in
319 Table 6 for BRAVO data only, Senterra data only, and then the combined datasets.

320 For the small wavelength regions taken individually, the best prediction accuracy is
321 results from the energy range between 800 and 880 cm^{-1} ; although this was known from practice,
322 our data show the advantage quantitatively. The range between 400 and 700 cm^{-1} , attributed to
323 the internal bending modes of the anion as mentioned above, should cause minimal centroid
324 shifts (Chopelas, 1991; Kuebler et al., 2006), and this is reflected in the observed high RMSE
325 values for that range.

326 Interestingly, the best prediction accuracy comes from models that cover energy ranges
327 that include multiple peaks. The $300\text{-}1500\text{ cm}^{-1}$ PLS model resulted in RMSE values of ± 3.87 ,
328 ± 4.52 , and $\pm 5.48\%$ Fo for the BRAVO, Senterra, and BRAVO + Senterra models, respectively.
329 Better or comparable performance was found using only the five “useful” regions noted above:
330 ± 1.67 , ± 4.40 , and $\pm 5.79\%$ Fo for PLS models and ± 7.13 , ± 3.93 , and $\pm 5.06\%$ Fo for Lasso
331 models. The superiority of the five-region models over the whole-spectrum models was most
332 dramatic for the Lasso models in larger datasets.

333 DISCUSSION

334 Understanding centroid variability for equivalent %Fo

335 This study deliberately chose to include both natural and synthetic olivine samples.
336 Minor cation substitutions within the natural samples cause variations within the band centroid
337 positions affecting the accuracy of the models. Therefore samples with the same %Fo can
338 produce bands at slightly different wavenumber positions. This could be mitigated by reducing

339 the size of the model to only include samples that fall along the prediction line or by acquiring
340 many spectra of the same sample to build the model. However, the goal of this study was to
341 create a broadly applicable model to predict olivine composition in natural samples, so we chose
342 not to remove samples that did not fall on the line of a perfect fit.

343 Additionally, this study intentionally acquired data on multiple spectrometers to
344 understand the implications of making composition predictions across different laboratories.
345 Variations observed between the two spectrometers used in this study are likely typical of
346 comparisons that would be encountered in comparisons to data from other spectrometers (see
347 Dyar et al. 2016c). This dilemma justifies using multivariate analysis rather than the univariate
348 band shift method and can be mitigated by spectrometer alignment.

349 **Restricting energy range to common compositions**

350 The vast majority of naturally-occurring olivines contain relatively high Mg contents
351 (Figure 3), commonly around Fo₈₀₋₉₀. So a test was developed to determine if improved accuracy
352 could be obtained by limiting samples in the training set to those with Fo contents greater than
353 50%. An aggregated model of Senterra, BRAVO and RRUFF* data with %Fo values greater
354 than 50% was constructed. The LOO RMSE-CV of this model (± 4.22 for the five-region PLS
355 model) was smaller than the original dataset (± 9.45 %Fo). Model accuracy likely improves
356 because the fayalite doublet is often poorly resolved. Therefore, in the vast majority of
357 applications, it may only be necessary to distinguish between the compositional variations of
358 forsterite rather than the entire olivine solid-solution. In these cases, the >50% Fo model would
359 be preferable because of its smaller error.

360 **Univariate analyses**

361 Because the DB2 peak centroid covers a much wider energy shift with changing
362 composition (compare *x*-axis limits in Figure 4 top and bottom), it might be expected that its use
363 would result in better prediction accuracy than using the DB1 peak. This is indeed observed for
364 data from nearly all the datasets studied (Table 4). The exceptions are in the combined BRAVO
365 + Senterra datasets and for the datasets extracted from published papers by Kuebler et al. (2006)
366 and Yasuzuka et al. (2009), but there is no apparent effect due to dataset size and resolution. It is
367 notable that for the largest and most diverse datasets, DB2 fits always produce the best prediction
368 accuracy. In particular, the model utilizing all available data, including those from the Senterra,
369 and BRAVO instruments plus RRUFF* and other data from Kuebler et al. (2006), Yasuzuka et
370 al. (2009), and Ishibashi et al. (2011) for which forsterite composition is >50%Fo produces the
371 best prediction accuracy for univariate analyses, and thus its DB2 regression equation (%Fo = -
372 $0.179625x^2 + 310.077x - 133717$) is recommended for use in applications where only peak
373 centroids can be resolved, as might be the case in a rock spectrum where the region of interest is
374 highly overlapped.

375 **Multivariate analyses**

376 It is apparent from Table 4 that both DB1 and DB2 peaks contain information about
377 composition, so using only one of them for predictions discards useful information. In contrast,
378 multivariate analyses offer the possibility to leverage information from anywhere in the spectra.
379 Table 5 shows the relative accuracies of the BRAVO, Senterra, and RRUFF predictions
380 individually and collectively in different combinations. The latter include all 188 of the RRUFF
381 olivine data and the RRUFF*, in combinations with the BRAVO and Senterra data acquired for
382 this project (Figure 7). As observed in the previous section, five-region models covering the

383 known olivine peaks show superior prediction performance over the whole-spectrum models.

384 Why?

385 The answer to this question can be seen in Figure 8, which shows the magnitudes of the
386 PLS coefficients for models covering the 300-1500 cm^{-1} range along with the channels chosen by
387 a Lasso model with $\alpha=0.001$ for the combined BRAVO and Senterra models. It is clear that the
388 entire range from 300-1500 cm^{-1} is rich with information about olivine composition that has been
389 previously unutilized in models that employ only restricted energy ranges. However, the five-
390 region models likely outperform the whole-spectrum models because the latter may inadvertently
391 include to unwanted features resulting from sample heating, fluorescence and, cosmic spikes.
392 These types of noise are not consistent for each spectrum and they do not relate to Raman
393 features resulting from compositional variations. So there is justification for using as many
394 regions that correspond to known olivine modes as possible, and thus five-region models are
395 used for full LOO RMSE-CV models given in Table 5.

396 It must be noted that use of this approach will make composition difficult to predict in
397 practice because pure olivine is rarely encountered in field applications. In practice (as when
398 deployed on a planetary surface), it is far more likely that the olivine will be mixed in with other
399 phases such as glass and other minerals. Thus it is desirable to have an alternate method for
400 predicting %Fo that isolates the olivine part of the spectra, for which the range from 800 to 880
401 cm^{-1} is recommended if there is no overlap from other non-olivine features.

402 In the largest datasets, PLS has comparable error to univariate band shift predictions,
403 while being significantly less time consuming. Thus the most accurate %Fo prediction based on
404 Raman spectra of pure olivine samples would be acquired at the highest resolution data possible,

405 pre-processed to remove baseline, normalized, and predicted using a PLS algorithm built from
406 the maximum number of samples (here, 281).

407 Using solely the peak centroid to model olivine composition does not utilize information
408 contained in other characteristics within the spectrum such as band shape, intensity, FWHM, area
409 and anomalies/noise within the spectra. Multivariate analyses appear to overcome these effects.

410 PLS and Lasso models examine multiple channels within the spectra to build a %Fo
411 prediction model. The number of coefficients per model is based on the assigned number of
412 components from the alpha value. As the number of channels in a Lasso model increases,
413 RMSE-CV decreases, showing the importance of models with a large number of channels.
414 However, there is a trade-off between the generalizability of the model that is optimized by
415 smaller numbers of channels versus improved accuracy from using larger numbers of channels.

416 **Factors influencing prediction accuracy**

417 Both dataset size and spectral resolution influence prediction accuracy for both univariate
418 and multivariate models. Results presented in Table 5 inform the effects of these factors.

419 To test the effect of resolution, spectral data and models were resampled to 3.0
420 cm^{-1} /channel resolution and compared against the native resolution of each instrument, which is
421 2.0, 0.5 and 0.48-2.0 cm^{-1} /channel for the BRAVO, Senterra (532 nm) and RRUFF datasets,
422 respectively. For comparison, the SuperCam instrument on Mars will have a pixel resolution of
423 2.5 cm^{-1} (Wiens et al. 2017), while the ExoMars RLS will use $<1 \text{ cm}^{-1}$ (Moral et al. 2018). Based
424 on the data in this study, there is no systematic effect of spatial resolution on prediction accuracy.

425 The effects of dataset size can also be roughly evaluated using the data collected here.
426 Individual datasets produce smaller (more accurate) LOO RMSE-CV values than aggregated
427 ones because the instrument and operating conditions are identical. For example, the combined

428 BRAVO and Senterra datasets had PLS LOO RMSE-CV values of 7.69 for the 800-880 cm^{-1}
429 range, while their individual, independently produced PLS LOO RMSE-CVs are 6.85 (BRAVO)
430 and 7.06 (Senterra) (Table 5). However, these single instrument models are less generalizable.
431 As the aggregated models increase in size and spectral diversity (i.e., instrument, laser
432 wavelength), the LOO RMSE-CV for the multivariate models decreases and the advantages of a
433 sole-source dataset diminish. Additionally, LOO RMSE-CV values can be reduced for aggregate
434 datasets by aligning the spectral data of the multiple instruments (Table 5).

435 A five-region PLS model of the collective Senterra, BRAVO, and RRUFF* datasets with
436 solely spectra corresponding to $>50\% \text{Fo}$ is recommended for future work. PLS (LOO RMSE-CV
437 is 4.22 $\% \text{Fo}$) is less time-consuming than univariate analyses (LOO RMSE-CV 4.58 $\% \text{Fo}$), gives
438 comparable accuracy, and is more generalizable. The PLS prediction model will aid workers
439 using different spectrometers, incident laser wavelengths and other operating conditions.
440 Additionally, a $\% \text{Fo}$ prediction for forsterite is more likely necessary than that of fayalite,
441 especially for applications to Mars.

442 It is hoped that future workers will add Raman spectral data to the recommended models
443 presented here, for which data are available on the lab website (nemo.cs.umass.edu:54321)
444 (Carey et al. 2017). Increasing the number of spectra within the models on additional instruments
445 will likely improve the accuracy of prediction results. Future workers with Raman spectra of
446 olivines with known compositions are encouraged to contact these authors so that a new,
447 expanded olivine prediction model can be created. Current and future improved PLS models will
448 be available from the authors.

449 **Parameters for model comparisons**

450 Tabulated R^2 values in Table 4 and S5 make it apparent that R^2 is a biased and potentially
451 misleading parameter when used to compare models of differing data. In this application, R^2 is a
452 measure of the proportion of the total variation of peak position from the average peak position
453 in that dataset that is explained by the regression line (McKillip and Dyar 2010). For example, if
454 all the samples in the dataset have comparable %Fo contents, then the deviation from that
455 average will be small, and R^2 may be misleadingly high. Moreover, R^2 cannot be used to
456 evaluate whether the calculated regression function is a correct description of the relationship
457 between peak position and %Fo. Finally, the R^2 statistic does not evaluate potential performance
458 of the linear or polynomial trend when applied to unseen (i.e., from a different dataset or
459 instrument) results.

460 The importance of this point is reinforced by the data in Table 4, in which only RMSE
461 can be used to make apples-to-apples comparisons among models. For example, Kuebler et al.
462 (2006) modeled %Fo using a second-order polynomial fit to the DB1 and DB2 centroids plotted
463 against %Fo, yielding regression lines with R^2 values of 0.98 (DB1) and 0.97 (DB2) %Fo,
464 respectively. When these data were used to calculate LOO RMSE-CV, a different story emerges.
465 Corresponding LOO RMSE-CV values are ± 4.33 and ± 4.57 %Fo units. The inaccuracy of this
466 model becomes even more apparent when it is used to predict a different dataset (RMSE-test
467 values). These results underscore the importance of evaluating model accuracies based on use of
468 leave-one-out cross-validation and/or “unseen” external data. Failure to do so invalidates any
469 claims of accuracy for application to other datasets, such as those on Mars.

470 **Future work**

471 In a study using 3,950 RRUFF spectra to test for matching accuracy, Carey et al. (2015)
472 tested the effects of common spectrum pre-processing steps, such as intensity normalization,

473 smoothing, squashing, and customized baseline removal (Giguere et al. 2017). Although
474 differences in sample crystal orientation, laser polarization, focus, and other instrumental
475 parameters can have major effects on spectra, even on the same samples and identical
476 instruments, Carey et al. (2015) showed that pre-processing techniques can effectively
477 ameliorate these differences and improve mineral identification. It is likely that optimizing pre-
478 processing of olivine spectra from disparate sources might also improve prediction accuracy for
479 obtaining %Fo from Raman spectra of olivine.

480 However, time did not permit testing of various pre-processing techniques on our own
481 datasets, though this is obviously an area ripe for research. In future work, effects of baseline
482 removal methods, alternate methods for normalization, and squashing and smoothing techniques
483 will be evaluated. Although olivine is an important rock-forming mineral group, there are many
484 other mineral groups that will need to be evaluated and quantitatively described through Raman
485 spectroscopy. Given the ubiquitous presence of species from the pyroxene and feldspar mineral
486 groups in igneous rocks and on planetary surfaces, creation of equivalent multivariate Raman
487 models for these groups should be a high priority for further research. Eventually, bringing our
488 results together with those analogous models for other phases will make it possible to identify
489 and quantify compositions of these phases in mineral mixtures.

490 **IMPLICATIONS**

491 As the technology for micro-Raman, reductions in the price of lasers, and
492 implementations for portable and remote instruments continue, Raman spectroscopy should play
493 an increasingly important role in geosciences and planetary exploration. Studies such as this one
494 that relate prominent Raman peaks to mineral identification and composition are needed to
495 enhance the capabilities of Raman instruments across those many applications. This paper lays a

496 foundation for future analogous studies of important rock-forming minerals by demonstrating
497 that specific features and energy ranges can be mined for accurate predictions of chemistry. This
498 project focuses on olivine, a common liquidus phase in magmatic systems that is present in most
499 basalts, in meteorites, and on terrestrial surfaces beyond Earth. Characterizing the ratio of Fe to
500 Mg in olivine constrains phase relations and crystallization conditions, and is feasible in both
501 pure spectra of olivine and in mixtures where the most prominent olivine doublet can be
502 resolved. Algorithms presented here will assist the *ExoMars* and *Mars 2020* mission teams in
503 recognizing olivine and determining its composition with known accuracy. We recommend use
504 of either our PLS model if pure olivine is encountered, or the recommended $\%Fo = -0.179625x^2$
505 $+310.077x -133717$ regression equation (where $x = DB2$ in units of cm^{-1}) for olivine that is
506 present in mixtures.

507 ACKNOWLEDGMENTS

508 We thank the Massachusetts Space Grant Consortium for student funding in support of
509 this project.

510 REFERENCES CITED

- 511 Abdu, Y.A., Varela, M.E., and Hawthorne, F.C. (2011) Raman, FTIR, and Mössbauer
512 spectroscopy of olivines from the D'Orbigny meteorite. Annual Meteoritical Society
513 Meeting, 74, 5112 (abstract).
- 514 Allen, F.S., Zhao, J., and Butterfield, D.S. (2000) Apparatus for Measuring and Applying
515 Instrumentation Correction to Produce a Standard Raman Spectrum. 6,141,095.
- 516 Banfield, J.M., Dyar, M.D., and McGuire, A.V. (1992) The defect microstructure of oxidized
517 mantle olivine from Dish Hill, California. *American Mineralogist*, 77, 959-975.

- 518 Breitenfeld, L.B. (2017) Predicting olivine composition using Raman spectroscopy through band
519 shift and multivariate analyses, 55-56 p. B.A. thesis, Mount Holyoke College, South
520 Hadley.
- 521 Butler, N.A., and Denham, M.C. (2000) The peculiar shrinkage properties of partial least squares
522 regression. *Journal of the Royal Statistical Society*, 62, 585-593.
- 523 Byrne, S.A., Dyar, M.D., Bessette, E.E., Breitenfeld, L.B., Crowley, M.C., Hoff, C.M.,
524 Marchand, G.J., Ketley, M.N., Roberts, A.L., Sklute, E.C., and Parente, M. (2015) Pure
525 mineral separates for mixing experiments to simulate planetary surfaces. *Lunar and*
526 *Planetary Science Conference*, 46, 1499 (abstract).
- 527 Carey, C., Dyar, M.D., Boucher, T., and Mahadevan, S. (2015) Machine learning tools for
528 mineral recognition and classification from Raman spectroscopy. *Journal of Raman*
529 *Spectroscopy*, 46, 894-903.
- 530 Carey, C., Dyar, M.D., Boucher, T., and Giguere, S. (2017) Web-based software for
531 preprocessing, matching, fitting, prediction, and visualization of spectroscopic data: the
532 data exploration, visualization, and analysis of spectra (DEVAS) website. *Lunar and*
533 *Planetary Science Conference*, 48, 1097 (abstract).
- 534 Chopelas, A. (1991) Single crystal Raman spectra of forsterite, fayalite, and monticellite.
535 *American Mineralogist*, 76, 1101-1109.
- 536 Dyar, M.D. (2003) Ferric iron in SNC meteorites as determined by Mössbauer spectroscopy:
537 Implications for Martian landers and Martian oxygen fugacity. *Meteoritics & Planetary*
538 *Science*, 38, 1733-1752.
- 539 Dyar, M.D., McGuire, A.V., and Ziegler, R.D. (1989) Redox equilibria and crystal chemistry of
540 coexisting minerals from spinel lherzolite mantle xenoliths. *American Mineralogist*, 74,

- 541 969-980.
- 542 Dyar, M.D., McGuire, A.V., and Harrell, M.D. (1992) Crystal chemistry of iron in two styles of
543 metasomatism in the upper mantle. *Geochimica et cosmochimica acta*, 56, 2579-2586.
- 544 Dyar, M.D., Delaney, J.S., Sutton, S.R., and Schaefer, M.W. (1998) Fe³⁺ distribution in oxidized
545 olivine: A synchrotron micro-XANES study. *American Mineralogist*, 83, 1361-1365.
- 546 Dyar, M.D., Carmosino, M.L., Speicher, EA., Ozanne, M.V., Clegg, S.M., and Wiens, R.C.
547 (2012a) Comparison of partial least squares and lasso regression techniques for laser-
548 induced breakdown spectroscopy of geological samples. *Spectrochimica Acta B*, 70, 51-
549 67.
- 550 Dyar, M.D., Fassett, C.I., Giguere, S., Lepore, K., Byrne, S., Boucher, T., Carey, C., and
551 Mahadevan, S. (2016a) Comparison of univariate and multivariate models for prediction
552 of major and minor elements from laser-induced breakdown spectra with and without
553 making. *Spectrochimica Acta B*, 123, 93-104.
- 554 Dyar, M.D., Breves, E.A., Gunter, M.E., Lanzirotti, A., Tucker, J.M., Carey, C., Peel, S.E.,
555 Brown, E.B., Oberti, R., Lerotic, M., and Delaney, J.S. (2016b) Use of multivariate
556 analysis for synchrotron micro-XANES analysis of iron valence state in amphiboles.
557 *American Mineralogist*, 101, 1171-1189.
- 558 Dyar, M.D., Breitenfeld, L.B., Carey, CJ, Bartholomew, P., Tague, T.J., Wang, P., Mertzman, S.,
559 Byrne, S.A., Crowley, M.C., Leight, C., Watts, E., Campbell, J.C., Celestian, A.,
560 McKeeby, B., Jaret, S., Glotch, T., Berlanga, G., and Misra, A. (2016c) Interlaboratory
561 and cross-instrument comparison of Raman spectra of 96 minerals. *Lunar and Planetary
562 Science* 47, Abstract #2240.
- 563 Dyar, M.D., Breves, E.A., Emerson, E., Bell, S.W., Nelms, M., Ozannes, M.V., Peel, S.E.,

- 564 Carmosino, M.L., Tucker, J.M., Gunter, M.E., Delaney, J.S., Lanzirrotti, A., and
565 Woodland, A.B. (2012b) Accurate determination of ferric iron in garnets by bulk
566 Mössbauer spectroscopy and synchrotron micro-XANES. *American Mineralogist*, 97,
567 1726-1740.
- 568 Dyar, M.D., Sklute, E.C., Menzies, O.N., Bland, P.A., Lindsley, D., Glotch, T., Lane, M.D.,
569 Wopenka, B., Klima, R., Bishop, J.L., Hiroi, T., Pieters, C.M., and Sunshine, J. (2009)
570 Spectroscopic characteristics of synthetic olivines, with an emphasis on fayalite: An
571 integrated multi-wavelength approach. *American Mineralogist*, 94, 883-898.
- 572 Floran, R.J., Prinz, M., Hlava, P.F., Keil, K., Nehru, C.E., and Hinthorne, J.R. (1978) The
573 Chassigny meteorite: A cumulate dunite with hydrous amphibole-bearing melt inclusions.
574 *Geochimica et Cosmochimica Acta*, 42, 1213-1229.
- 575 Foster, N.J., Wozniakiewicz, P.J., Burchell, M.J., Kearsley, A.T., Creighton, A.J., and Cole, M.J.
576 (2007) Identification by Raman spectroscopy of processing effects in forsterite- fayalite
577 samples during hypervelocity impacts on foils and capture in aerogel. *Annual*
578 *Meteoritical Society Meeting*, 70, 5186 (abstract).
- 579 Gaisler, S.V., and Kolesov, B.A. (2007) Raman spectra of olivine solid solution $(\text{Fe}_x\text{Mg}_{1-x})_2\text{SiO}_4$
580 and spin-vibration interaction. *Journal of Structural Chemistry*, 48, 61-65.
- 581 Giguere, S., Boucher, T., Carey, C.J., Mahadevan, S., and Dyar, M.D. (2017) A fully-customized
582 baseline removal framework for spectroscopic applications. *Applied Spectroscopy*, 71,
583 1457-1470.
- 584 Guyot, F., Boyer, H., Madon, M., Velde, B., and Poirier, J.P. (1986) Comparison of the Raman
585 microprobe spectra of $(\text{Mg, Fe})_2\text{SiO}_4$ and Mg_2GeO_4 with olivine and spinel structures.
586 *Physics and Chemistry of Minerals*, 13, 91-95.

- 587 Hastie, T., Tibshirani, R., and Friedman, J. (2009) *The Elements of Statistical Learning*, 745 p.
588 Springer Science, New York.
- 589 Iishi, K. (1978) Lattice dynamics of forsterite. *American Mineralogist*, 63, 1198-1208.
- 590 Isaacson, P.J., Klima, R.L., Sunshine, J.M., Pieters, C.M., Hiroi, T., and Dyar, M.D. (2014)
591 Visible to near-infrared reflectance spectroscopy of pure synthetic olivine across the
592 olivine solid solution. *American Mineralogist*, 99, 467-478.
- 593 Ishibashi, H., Arakawa, M., Yamamoto, J., and Kagi, H. (2011) Precise determination of Mg/Fe
594 ratios applicable to terrestrial olivine samples using Raman spectroscopy. *Journal of*
595 *Raman Spectroscopy*, 43, 331-337.
- 596 Ishibashi, H., Arakawa, M., Ohi, S., Yamamoto, J., Miyake, A., and Kagi, H. (2008)
597 Relationship between Raman spectral pattern and crystallographic orientation of a rock-
598 forming mineral: a case study of $\text{Fo}_{89}\text{Fa}_{11}$ olivine. *Journal of Raman Spectroscopy*, 39,
599 1653-1659.
- 600 Kolesov, B.A., and Geiger, C.A., (2004) A Raman spectroscopic study of Fe–Mg olivines.
601 *Physics and Chemistry of Minerals*, 31, 142–154.
- 602 Kolesov, B.A., and Tanskaya, J.V. (1996) Raman spectra and cation distribution in the lattice of
603 olivines. *Materials Research Bulletin*, 31, 1035-1044.
- 604 Kuebler, K.E., Jolliff, B.L., Wang, A., and Haskin, L.A. (2006) Extracting olivine (Fo-Fa)
605 compositions from Raman spectral peak positions. *Geochimica et Cosmochimica Acta*,
606 70, 6201-6222.
- 607 Lane, M.D., Glotch, T.D., Dyar, M. D., Pieters, C.M., Klima, R., Hiroi, T., Bishop, J.L., and
608 Sunshine, J. (2011) Midinfrared spectroscopy of synthetic olivines: Thermal emission,
609 specular and diffuse reflectance, and attenuated total reflectance studies of forsterite to

- 610 fayalite. *Journal of Geophysical Research*, 116, E08010, DOI: 10.1029/2010JE003588.
- 611 McCanta, M.C., Treiman, A.H., Dyar, M.D., Alexander, C.M.O'D., Rumble, D., and Essene,
612 E.J. (2008) The LaPaz Icefield 04840 meteorite: Mineralogy, metamorphism, and origin
613 of an amphibole-and biotite-bearing R chondrite. *Geochimica et Cosmochimica Acta*, 72,
614 5757-5780.
- 615 McGuire, A.V., Dyar, M.D., and Nielson, J.E. (1991) Metasomatic oxidation of upper mantle
616 periodotite. *Contributions to Mineralogy and Petrology*, 109, 252-264.
- 617 McGuire, A.V., Francis, C.A., and Dyar, M.D. (1992). Mineral standards for electron
618 microprobe analysis of oxygen. *American Mineralogist*, 77, 1087-1087.
- 619 McKeown, D.A., Bell, M.I., and Caracas, R. (2010) Theoretical determination of the Raman
620 spectra of single-crystal forsterite (Mg₂SiO₄). *American Mineralogist*, 95, 980-986.
- 621 McKillop, S., and Dyar, M.D. (2010) *Geostatistics explained: An introductory guide for earth*
622 *scientists*, 396 p. Cambridge University Press, Cambridge, UK.
- 623 McSween, H.Y., and Jarosewich, E. (1983) Petrogenesis of the Elephant Moraine A79001
624 meteorite: Multiple magma pulses on the shergottite parent body. *Geochimica et*
625 *Cosmochimica Acta*, 47, 1501-1513.
- 626 Mikouchi, T., and Kurihara, T. (2008) Mineralogy and petrology of paired lherzolitic shergottites
627 Yamato 000027, Yamato 000047, and Yamato 000097: Another fragment from a Martian
628 "lherzolite" block. *Polar Science*, 2, 175-194.
- 629 Mohanan, K., Sharma, S.K., and Bishop F.C. (1993) A Raman spectral study of forsterite-
630 monticellite solid solutions. *American Mineralogist*, 78, 42-48.
- 631 Moral, A.G., Rull, F., Maurice, S., Hutchinson, I., Canora, C.P., Seoane, L., Canchal, R.,
632 Gallego, P., Ramos, G., Prieto, J.A.R., Santiago, A., Santamaría, P., Colombo, M.,

- 633 Belenguer, T., López, G., Quintana, C., Zafra, J., Berrocal, A., Pintor, C., Cabrero, J., and
634 Saiz, J. (2018) Raman laser spectrometer for 2020 ExoMars Mission. Engineering and
635 qualification model capabilities and future activities. Lunar and Planetary Science
636 Conference, 49, 2449 (abstract).
- 637 Morse, S.A. (1996) Kiglapait mineralogy .3. Olivine compositions and Rayleigh fractionation
638 models. *Journal of Petrology*, 37, 1037-1061.
- 639 Morse, S.A. (2001) Augite-olivine equilibria in the Kiglapait intrusion, Labrador, Canada.
640 *Canadian Mineralogist*, 39, 267-274.
- 641 Mouri, T., and Enami, M. (2008) Raman spectroscopic study of olivine-group minerals. *Journal*
642 *of Mineralogical and Petrological Sciences*, 103, 100-104.
- 643 Mullen, T., Dyar, M.D., Parente, M., and Breitenfeld, L.B. (2018) Improving matching accuracy
644 in Raman spectroscopy by quantifying the wavenumber shift between instruments. *Lunar*
645 *and Planetary Science*, 49, Abstract #1185.
- 646 Parques-Ledent, M.T., and Tarte, P. (1973) Vibrational studies of olivine-type compound-I. The
647 i.r. and Raman spectra of the isotopic species of Mg_2SiO_4 . *Spectrochimica Acta A*, 29,
648 1007-1016.
- 649 Piriou, B. (1983) The high-frequency vibrational spectra of vitreous and crystalline
650 orthosilicates. *American Mineralogist*, 68, 426-443.
- 651 Price, G.D., Parker, S.C., and Leslie, M., (1987) The lattice dynamics of forsterite. *Mineralogical*
652 *Magazine* 51, 157-170.
- 653 Schaefer, M.W. (1983) Crystal chemistry of ferric-rich fayalites, 38-56 p. Ph.D. thesis, MIT,
654 Cambridge.
- 655 Servoin, J., Piriou, B., and Alain, P. (1972) Diffusion Raman de la forsterite sythetique pure,

- 656 Mg_2SiO_4 . *Comptes Rendus de l'Académie des Sciences Paris*, 274, 135-137.
- 657 Sklute, E.C. (2006) Mössbauer spectroscopy of synthetic olivine across the Fe-Mg solid solution,
658 1-38 p. B.A. thesis, Mount Holyoke College, South Hadley.
- 659 Stone, M., and Brooks, R.J. (1990) Continuum regression: cross-validated sequentially
660 constructed predictions embracing ordinary least-squares, partial least-squares, and
661 principal components regression. *Journal of the Royal Statistical Society B*, 52, 237–269.
- 662 Tibshirani, R. (1995) Regression shrinkage and selection via the lasso. *Journal of the Royal*
663 *Statistical Society*, 58, 267-288.
- 664 Tucker, J.M., Dyar, M.D., Schaefer, M.W., Clegg, S.M., and Wiens, R.C. (2010) Optimization of
665 laser-induced breakdown spectroscopy for rapid geochemical analysis. *Chemical*
666 *Geology*, 277, 137-148.
- 667 Wang, A., Kuebler, K., Jolliff, B., and Haskin, L.A. (2004) Mineralogy of a Martian meteorite as
668 determined by Raman spectroscopy. *Journal of Raman Spectroscopy*, 35, 504-514.
- 669 Wiens, R.C., Newell, R., Clegg, S., Sharma, S.K., Misra, A., Bernardi, P., Maurice, S., McCabe,
670 K., Cais, P., and the SuperCam Science Team (2017) The SuperCam remote Raman
671 spectrometer for Mars 2020. *Lunar and Planetary Science*, 48, 2600 (abstract).
- 672 Wold, S., Martens, H., and Wold, H. (1983) The multivariate calibration problem in chemistry
673 solved by the PLS method. *Lecture Notes in Mathematics*, 973, 286-293.
- 674 Yasuzuka, T., Ishibashi, H., Arakawa, M., Yamamoto, J., and Kagi, H. (2009) Simultaneous
675 determination of Mg# and residual pressure in olivine using micro-Raman spectroscopy.
676 *Journal of Mineralogical and Petrological Sciences*, 104, 395-400.
- 677 Zhang, Z.M., Chen, S., and Liang, Y.Z. (2010) Baseline correction using adaptive iteratively
678 reweighted penalized least squares. *Royal Society of Chemistry*, 135, 1138-1146.

680

Figure captions

681 Figure 1. High intensity Raman doublet ($815\text{-}857\text{ cm}^{-1}$) of forsterite (purple) and fayalite (green).

682 Figure 2. Unaligned Raman spectra of olivine doublet (DB1 and DB2) of 93 samples acquired on
683 Bruker's Senterra and BRAVO spectrometers. All spectra were baseline removed using Air-PLS
684 and normalized to a maximum intensity of 1. Spectra are color-coded based on Fo content, where
685 forsterite is represented with yellow, fayalite with purple, and intermediate compositions in
686 between.

687 Figure 3. Histogram of 93 synthetic (blue) and natural (red) samples on the Fo-Fa series. Natural
688 olivines typically form with a %Fo of ~ 89.5 resulting in an unbalanced distribution on the Fo-Fa
689 series.

690 Figure 4. Fo by EMPA versus peak centroids of DB1 and DB2. Second order polynomial fits
691 and RSME-CV values are included for the unaligned data acquired on Bruker's BRAVO ($n=25$)
692 and Senterra ($n=68$) spectrometers. Error bars are smaller than the symbols and are given in
693 Tables S3 and S4.

694 Figure 5. %Fo by EMPA versus peak centroids positions of (a) DB1 and (b) DB2 for data
695 acquired in other studies. Including results from the RRUFF* database (see text for explanation
696 of notation), and studies by Kuebler et al. (2006), Yasuzuka et al. (2009), and Ishibashi et al.
697 (2011).

698 Figure 6. Variations in Lasso model accuracy as a function of the number of coefficients. As α
699 increases, fewer channels are examined: (a) 38 channels for $\alpha = 0.001$, (b) 27 channels for $\alpha =$
700 0.01 , and (c) 10 channels for $\alpha = 0.1$. As the number of channels examined is decreased (fewer

701 coefficients within the model), the RMSE-CV of the models increases in value (gets worse). This
702 demonstrates the value of a models that examines a high number of channels, which is achieved
703 in a small α value Lasso model or PLS models.

704 Figure 7. Bar graph comparing results from Tables 4 and 5. When the dataset is small and/or all
705 the data are acquired on the identical instrument, then univariate methods produce better results
706 than those using multivariate analyses. However, as the number of samples and instruments used
707 increase, PLS methods generally produce more accurate results.

708 Figure 8. (top) Plot of BRAVO and Senterra unaligned data in Table 6, along with circles
709 indicating the magnitude of PLS coefficients (right-hand y axis units). Note that PLS
710 coefficients are proportional to spectral intensity at each wavenumbers, so absolute values cannot
711 be compared on this plot. However, the PLS coefficients do demonstrate that the entire
712 wavenumber range contains information useful to predicting composition. (bottom) Analogous
713 plot for the same data but showing Lasso coefficients for a model with $\alpha = 0.001$.

Table 1. Assignment of Raman Modes in Olivine

~Band (cm ⁻¹)	Symmetry	Assignment	Literature
815-825 (DB1)	A _g	n.a.	Servoin et al. (1972)
		v ₁ +v ₃	Paques-Ledent and Tarte (1973)
		v ₁	Iishi (1978)
		v ₁ +v ₃	Pirou (1983)
		v ₁ +v ₃	Chopelas et al. (1991)
		SiO ₄ ²⁻ stretching	Kolesov and Tanskaya (1996)
		v ₁ +v ₃	Kolesov and Geiger (2004)
		Si-O stretch, v ₃	McKeown et al. (2010)
838	B _{1g}	v ₁	Iishi (1978)
		v ₁ (+v ₃)	Chopelas et al. (1991)
		v ₁ (+v ₃)	Kolesov and Geiger (2004)
		v ₁	McKeown et al. (2010)
837-857 (DB2)	A _g	n.a.	Servoin et al. (1972)
		v ₁ +v ₃	Paques-Ledent and Tarte (1973)
		v ₃	Iishi (1978)
		v ₁ (+v ₃)	Pirou (1983)
		v ₃	Price et al. (1987)
		v ₁ +v ₃	Chopelas et al. (1991)
		SiO ₄ ²⁻ stretching	Kolesov and Tanskaya (1996)
		v ₁ +v ₃	Kolesov and Geiger (2004)
		Si-O stretch; SiO ₄ breathing v ₃	McKeown et al. (2010)
866	B _{1g}	v ₃	Iishi (1978)
		v ₃	Price et al. (1987)
		v ₃ (+v ₁)	Chopelas et al. (1991)
		v ₃ (+v ₁)	Kolesov and Geiger (2004)
882	B _{2g}	v ₃	McKeown et al. (2010)
		v ₃	Paques-Ledent and Tarte (1973)
		v ₃	Iishi (1978)
		v ₃	Chopelas et al. (1991)
		v ₃	Kolesov and Geiger (2004)
		v ₃	McKeown et al. (2010)

1

2

Table 2. Summary of %Fo prediction models using band shift method

Paper	# samples*	Bands (cm⁻¹)
Iishi (1978)	1	All bands
Guyot et al. (1986)	4	815-825, 837-857, 881-883, 914-920, 950-966
Chopelas et al. (1991)	1	All bands
Mohanan et al. (1993)	1	All bands from 200-1000
Kolesov and Tanskaya, (1996)	2	All bands from 200-1000
Wang et al. (2004)	0	815-825, 837-857
Kuebler et al. (2006)	10	815-825, 837-857
Foster et al. (2007)	2	815-825, 837-857
Gaisler and Kolesov (2007)	0	200-230, 290-310, 410-440, 815-825, 837-857
Mouri and Enami (2008)	0	815-825, 837-857
Ishibashi et al. (2008)	1	815-825, 837-857
Yasuzuka et al. (2009)	10	540-553, 815-825, 837-857
McKeown et al. (2010)	1	All bands
Abdu et al. (2011)	3	815-825, 837-857
Ishibashi et al. (2011)	15	815-825, 837-857, 881-883, 914-920, 950-966
Weber et al. (2014)	5	815-825, 837-857

Note: *Number of samples with reported centroid positions and compositions.

3

Table 3. Natural Samples Studied

Sample Name	Locality	Chemistry	Mössbauer	Raman Instrument	
Ba-1-61	Dish Hill, CA	UTK	[12]	Senterra	5
Ba-1-74	Dish Hill, CA	UTK	[12]	Senterra	
Ba-2-1 WR1	Dish Hill, CA	Brown	[6]	Senterra	
Ba-2-1 WR2	Dish Hill, CA	Brown	[6]	Senterra	
Ba-2-1 WR3	Dish Hill, CA	Brown	[6]	Senterra	
Ba-2-1 WR4	Dish Hill, CA	Brown	[6]	Senterra	
Ba-2-1 D-1	Dish Hill, CA	[1]	[6]	Senterra	
Ci-1-183	Dish Hill, CA	[12]	[12]	Senterra	
Ci-1-25	Dish Hill, CA	[12]	[12]	Senterra	
DH101-B	Dish Hill, CA	Brown	[6]	Senterra	
DH101-C	Dish Hill, CA	Brown	[6]	Senterra	
DH101-D	Dish Hill, CA	Brown	[6]	Senterra	
DH101-E	Dish Hill, CA	Brown	[6]	Senterra	
Dyar 89-190	unknown	Brown	[12]	BRAVO	
Dyar 89-12	unknown	Brown	[12]	BRAVO	
Dyar 89-187	unknown	Brown	[12]	BRAVO	
Dyar 89-194	unknown	Brown	[12]	BRAVO	
Ep-1-13	Potrillo maar, NM	[12]	[7]	BRAVO	
Ep-3-139-C	Kilbourne Hole, NM	Brown	[8]	Senterra	
Ep-3-139-D	Kilbourne Hole, NM	Brown	[8]	Senterra	
Ep-3-44	Kilbourne Hole, NM	UTK	[12]	Senterra	
Ep-3-46	Kilbourne Hole, NM	UTK	[12]	Senterra	
Ep-3-72	Kilbourne Hole, NM	UTK	[12]	Senterra	
Ep-3-7A	Kilbourne Hole, NM	Univ. Houston	[12]	Senterra	
KI-3003	Kiglapait Formation	Brown	[12]	Senterra	
KI-3373	Kiglapait Formation	Brown	[12]	Senterra	
NMNH 112085	Red Rock Ridge	UTK, Brown	[12]	Senterra	
NMNH 1210672	Germany Greifensteiner Kalk	UTK, Brown	[12]	Senterra	
NMNH 135841	Sweden Nykopig Tunaberg	Brown	[12]	Senterra	
NMNH 85539	unknown	UTK, Brown	[12]	Senterra	
Rockport	Rockport	Brown	[9]	Senterra	
Globe	Globe, AZ	[1]	[12]	BRAVO	
H279-12	Harrat al Kishb, Saudi Arabia	[12]	[12]	Senterra	
H30-82-8	Harrat al Kishb, Saudi Arabia	UTK	[12]	Senterra	
H30-B1	Harrat al Kishb, Saudi Arabia	Brown	[10]	BRAVO	
H30-B2	Harrat al Kishb, Saudi Arabia	[12]	[7]	BRAVO	
H30-B3	Harrat al Kishb, Saudi Arabia	UTK	[12]	Senterra	
H30-B4	Harrat al Kishb, Saudi Arabia	UTK, Brown	[12]	Senterra	
H30-B5	Harrat al Kishb, Saudi Arabia	UTK	[12]	Senterra	
H312-1	Harrat Uwayrid, Saudi Arabia	[12]	[12]	Senterra	
H366-28	Harrat Hutaymah, Saudi Arabia	[12]	[12]	BRAVO	
H366-30	Harrat Hutaymah, Saudi Arabia	[12]	[12]	Senterra	
NMNH 9140	Orange Co. NY	UTK, Brown	[12]	Senterra	
KBH-94-23-B	Kilbourne Hole, NM	UTK	[12]	Senterra	
KBH-94-23-E	Kilbourne Hole, NM	UTK	[12]	Senterra	

Notes: Sources are abbreviated as follows: [1] Byrne et al. (2015), [2] Floran et al. (1978), [3] McSween and Jarosewich (1983), [4] McCanta et al. (2008), [5] Mikouchi and Kurihara (2008), [6] McGuire et al. (1991), [7] Dyar et al. (1989), [8] Dyar et al. (1992), [9] Schaefer (1983), [10] McGuire et al. (1992), [11] Dyar (2003), [12] this study.

Table 3. (continued). **Natural Samples Studied**

Sample Name	Locality	Chemistry	Mössbauer	Raman Instrument
KBH-94-23-E	Kilbourne Hole, NM	UTK	[12]	Senterra
KI-3005	Kiglapait Formation	Brown	[12]	Senterra
KI-3289	Kiglapait Formation	Brown	[12]	Senterra
KI-3362	Kiglapait Formation	Brown	[12]	Senterra
KI-3648	Kiglapait Formation	UTK	[12]	Senterra
KI-4030	Kiglapait Formation	Brown	[12]	Senterra
Ki-5-16	Cima volcanic field, CA	[12]	[12]	Senterra
Ki-5-235	Cima volcanic field, CA	UTK	[12]	Senterra
Ki-5-35	Cima volcanic field, CA	UTK	[12]	Senterra
Ki-5-62	Cima volcanic field, CA	UTK	[12]	BRAVO
Pakistan	Pakistan Sapatime Kohistan District	Brown	[12]	BRAVO
San Carlos AZ	San Carlos AZ	[1]	[12]	BRAVO
ALHA 77005	Mars	UTK	[11]	Senterra
ALHA-77005-193	Mars	UTK	[11]	Senterra
Chassigny USNM E24	Mars	[2]	[11]	Senterra
EETA-79001 60B	Mars	[3]	[11]	Senterra
EETA-79001-A	Mars	[3]	[11]	Senterra
LAP-0484016	Mars	[4]	[12]	Senterra
NWA2737	Mars	Brown	[12]	Senterra
Y000097 86	Mars	[5]	[12]	Senterra

Notes: Sources are abbreviated as follows: [11] Dyar (2003), [12] this study.

6

7

Table 4. Model accuracy by LOO RMSE-CV for univariate analyses using 2nd order polynomial fits

Data	#Spectra	Resolution (cm ⁻¹ /channel)	Model	Internal R ²	Internal RMSE	LOO RMSE-CV
BRAVO	25	2.0	DB1	0.88	9.87	7.69
			DB2	0.97	5.26	4.35
Senterra	68	0.5	DB1	0.94	7.38	4.94
			DB2	0.97	5.20	3.35
Senterra + BRAVO	93	0.5 / 2.0	DB1	0.89	10.23	7.43
			DB2	0.89	10.20	7.91
Senterra + BRAVO aligned	93	0.5 / 2.0	DB1	0.92	8.72	5.64
			DB2	0.94	7.36	5.05
Kuebler et al. (2006)	13	6.2	DB1	0.98	8.63	4.33
			DB2	0.97	5.05	4.57
Yasuzuka et al. (2009)	10	0.05	DB1	0.97	1.49	1.68
			DB2	0.97	1.57	1.84
Ishibashi et al. (2011)	15	1.5	DB1	0.98	1.30	1.32
			DB2	0.99	0.92	0.84
RRUFF*	156	0.48 - 1.4	DB1	0.92	11.58	8.75
			DB2	0.94	9.93	8.28
RRUFF (all)	188	2.0	DB1	0.86	15.50	11.29
			DB2	0.93	10.70	8.86
Senterra + BRAVO+RRUFF*	249	0.5/2.0 /0.48 – 1.4	DB1	0.90	12.04	9.36
			DB2	0.92	10.58	8.91
Senterra + BRAVO + RRUFF (all)	281	0.5/2.0 /0.48 – 2.0	DB1	0.86	14.35	10.55
			DB2	0.92	11.02	9.20
Senterra + BRAVO+RRUFF* (>50%Fo)	181	0.5/2.0 /0.48 – 1.4	DB1	0.56	6.60	5.01
			DB2	0.61	6.20	4.58
All data (only RRUFF*)	287	see above	DB1	0.90	11.49	8.70
			DB2	0.92	10.26	8.56
All data for samples >50%Fo (only RRUFF*)	213	see above	DB1	0.64	6.30	4.86
			DB2	0.69	5.85	4.27

*Includes all olivine group spectra on RRUFF except those designated as “Broad Scan with Spectral Artifacts.”

8

Table 5. Model accuracy by LOO RMSE-CV for Multivariate Analyses

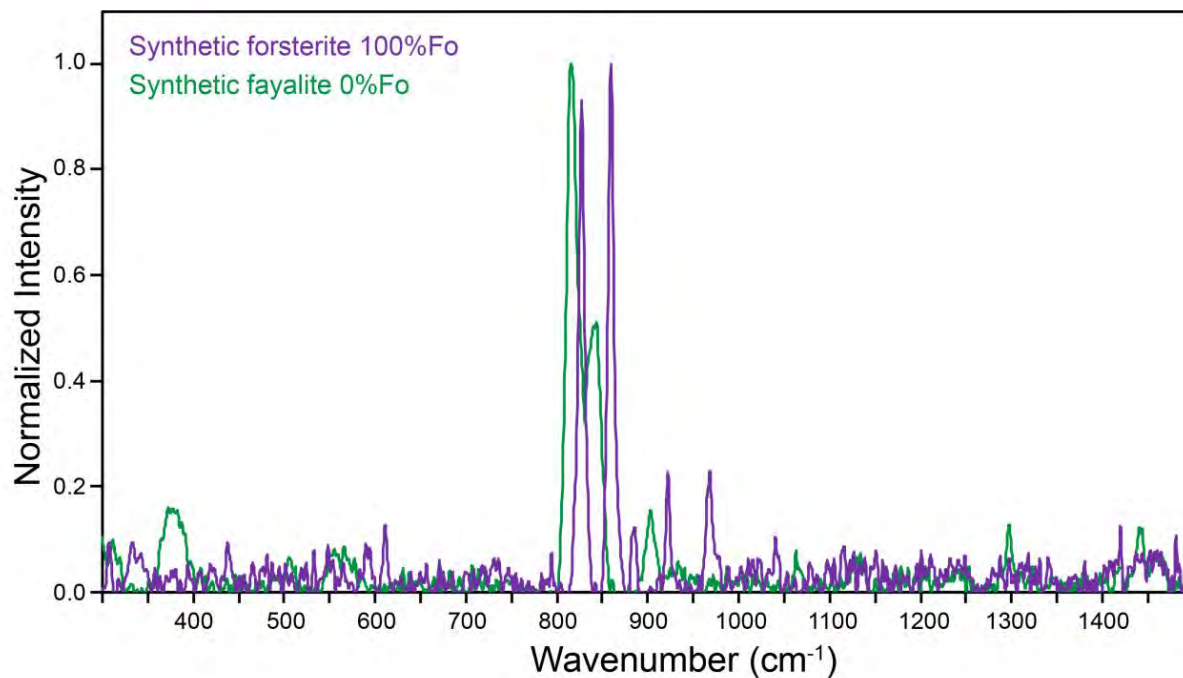
Data	#	Res.	Model	σ	C/a	Internal R ²	Internal RMSE	LOO RMSE-CV
BRAVO	25	3.0	PLS	800-880	7	0.99	2.56	6.85
		3.0	Lasso	800-880	0.001	0.99	1.13	6.87
		2.0	PLS	800-880	2	0.96	5.75	7.43
		2.0	Lasso	800-880	0.005	0.99	1.88	9.24
		2.0	PLS	5 regions	9	0.99	1.67	8.17
		2.0	Lasso	5 regions	0.3	0.93	7.13	7.91
Senterra	68	3.0	PLS	800-880	6	0.97	5.57	7.06
		3.0	Lasso	800-880	0.01	0.93	8.52	10.62
		0.5	PLS	800-880	7	0.96	5.85	7.31
		0.5	Lasso	800-880	0.008	0.95	7.37	11.75
		2.0	PLS	5 regions	9	0.98	4.40	6.45
		2.0	Lasso	5 regions	0.005	0.98	3.93	7.50
Senterra + BRAVO	93	2.0	PLS	800-880	5	0.95	6.63	7.79
		2.0	Lasso	800-880	0.002	0.95	6.95	9.38
		2.0	PLS	5 regions	9	0.96	5.79	7.91
		2.0	Lasso	5 regions	0.006	0.97	5.06	8.72
Senterra + BRAVO aligned	93	2.0	PLS	800-880	7	0.97	5.69	6.93
		2.0	Lasso	800-880	0.003	0.95	7.19	9.42
		2.0	PLS	5 regions	8	0.96	5.87	7.90
		2.0	Lasso	5 regions	0.02	0.95	7.07	9.74
	118	2.0	PLS	800-880	7	0.97	5.56	6.63
		2.0	Lasso	800-880	0.002	0.94	7.34	9.52
2.0		PLS	5 regions	9	0.96	5.78	7.59	
RRUFF*	156	1.5	PLS	800-880	5	0.98	6.10	6.65
		1.5	Lasso	800-880	0.01	0.94	9.97	12.25
		1.5	PLS	5 regions	9	0.99	4.56	5.97
		1.5	Lasso	5 regions	0.008	0.98	6.46	7.02
RRUFF (all)	188	2.0	PLS	800-880	10	0.97	7.13	8.34
		2.0	Lasso	800-880	0.001	0.93	11.00	13.15
		2.0	PLS	5 regions	10	0.97	7.30	8.53
		2.0	Lasso	5 regions	0.001	0.96	7.85	10.99
Senterra + BRAVO + RRUFF*	249	2.0	PLS	800-880	10	0.97	6.97	7.83
		2.0	Lasso	800-880	0.002	0.97	6.97	11.94
		2.0	PLS	5 regions	9	0.96	7.12	8.34
		2.0	Lasso	5 regions	0.002	0.97	7.06	11.29
Senterra + BRAVO + RRUFF (all)	281	2.0	PLS	800-880	9	0.95	8.87	9.64
		2.0	Lasso	800-880	0.001	0.91	11.77	13.41
		2.0	PLS	5 regions	10	0.95	8.57	9.45
		2.0	Lasso	5 regions	0.015	0.93	10.08	12.14
Senterra + BRAVO + RRUFF* (>50%Fo)	181	2.0	PLS	800-880	6	0.82	4.23	4.86
		2.0	Lasso	800-880	0.004	-0.28	11.25	13.74
		2.0	PLS	5 regions	9	0.89	3.30	4.22
		2.0	Lasso	5 regions	0.02	0.24	8.65	13.28

Notes: # = number of spectra, Res. = Resolution (cm⁻¹ /channel), σ = wavenumber range (cm⁻¹), C/a = number of components for PLS models and a for Lasso models. Five regions = 410-440 cm⁻¹, 538-556 cm⁻¹, 800-880 cm⁻¹, 908-926 cm⁻¹ and 950-968 cm⁻¹ combined. *Includes all olivine group spectra on RRUFF except those designated as "Broad Scan with Spectral Artifacts."

10

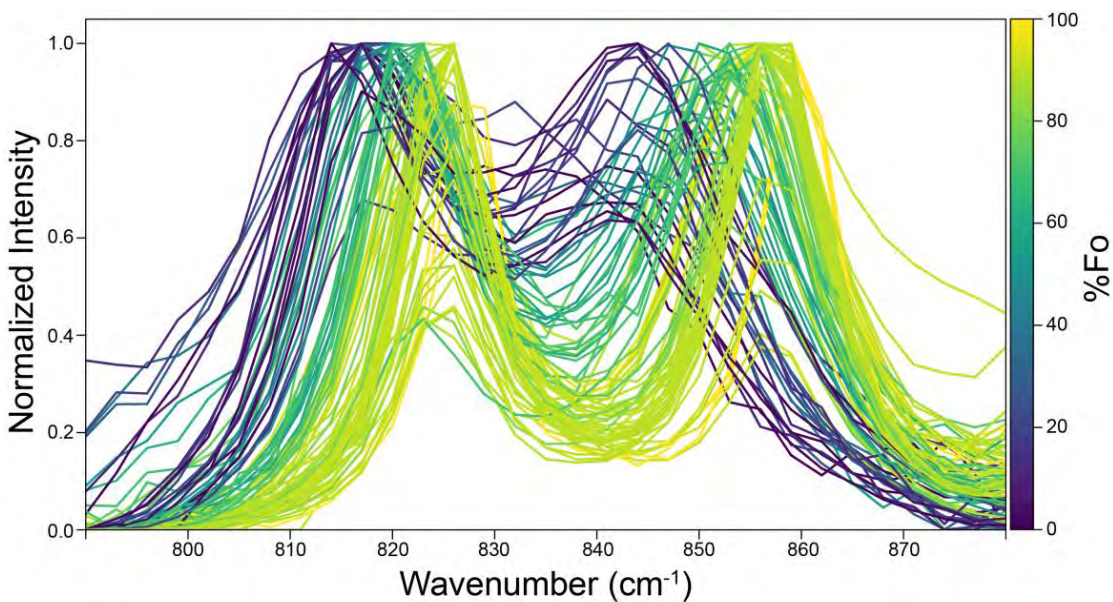
Table 6. Comparison of Internal RMSE Results by Energy Region with 2.0 cm⁻¹ /channel resolution

Number of spectra Model used	BRAVO 25		Senterra 68		BRAVO + Senterra 93	
	PLS	Lasso	PLS	Lasso	PLS	Lasso
410-440 cm ⁻¹	18.96	12.70	15.30	16.46	20.87	20.54
538-556 cm ⁻¹	20.82	19.43	22.60	34.09	23.11	31.77
800-880 cm ⁻¹	5.75	1.88	5.44	6.04	6.63	6.95
908-926 cm ⁻¹	16.04	14.46	12.91	14.74	14.45	16.14
950-968 cm ⁻¹	18.24	15.38	19.11	20.97	19.99	24.61
5 regions above	1.67	7.13	4.40	3.93	5.79	5.06
400-700 cm ⁻¹	17.34	16.99	8.01	6.37	13.07	12.74
300-1500 cm ⁻¹	3.87	6.76	4.52	9.42	5.48	3.74



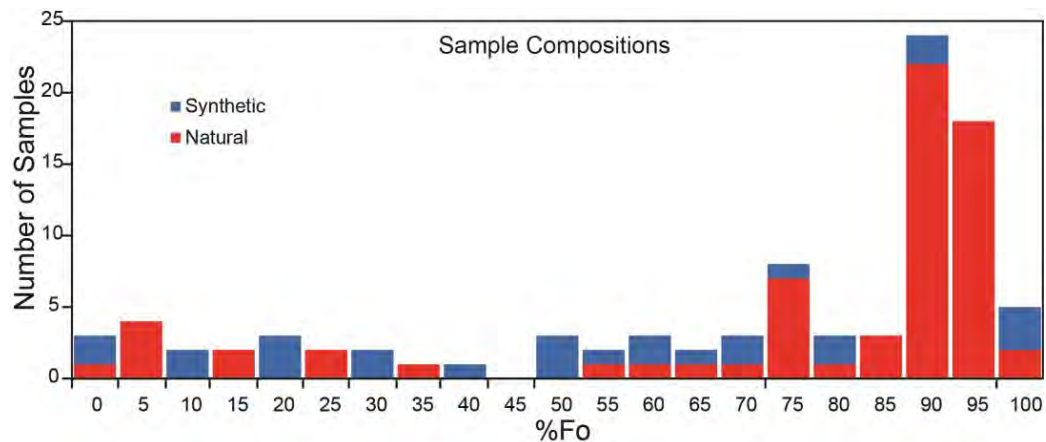
11

12 Figure 1. High intensity Raman doublet (815-857 cm^{-1}) of forsterite (purple) and fayalite (green).



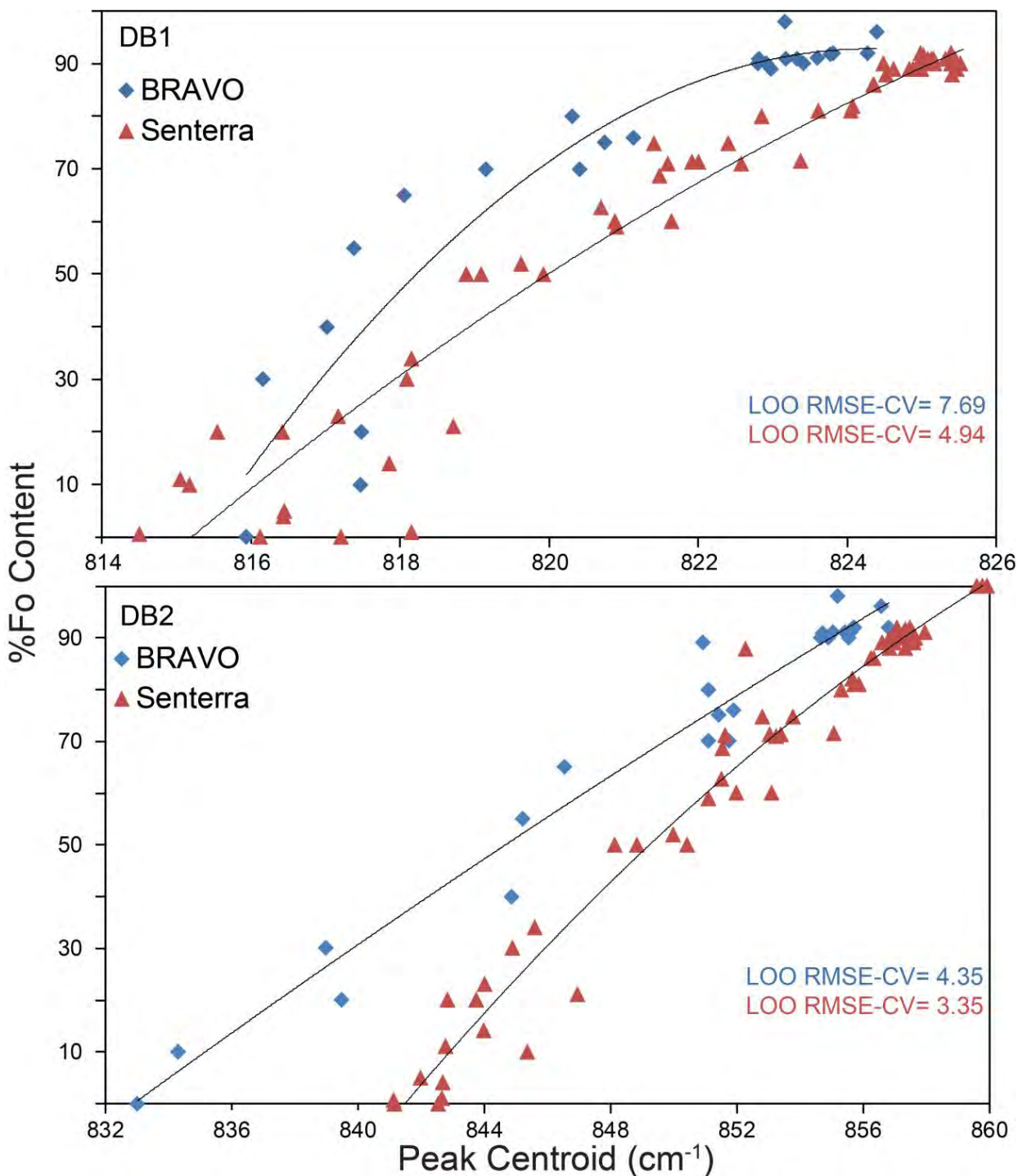
13

14 Figure 2. Unaligned Raman spectra of olivine doublet (DB1 and DB2) of 93 samples acquired on
15 Bruker's Senterra and BRAVO spectrometers. All spectra were baseline removed using Air-PLS
16 and normalized to a maximum intensity of 1. Spectra are color-coded based on Fo content, where
17 forsterite is represented with yellow, fayalite with purple, and intermediate compositions in
18 between.



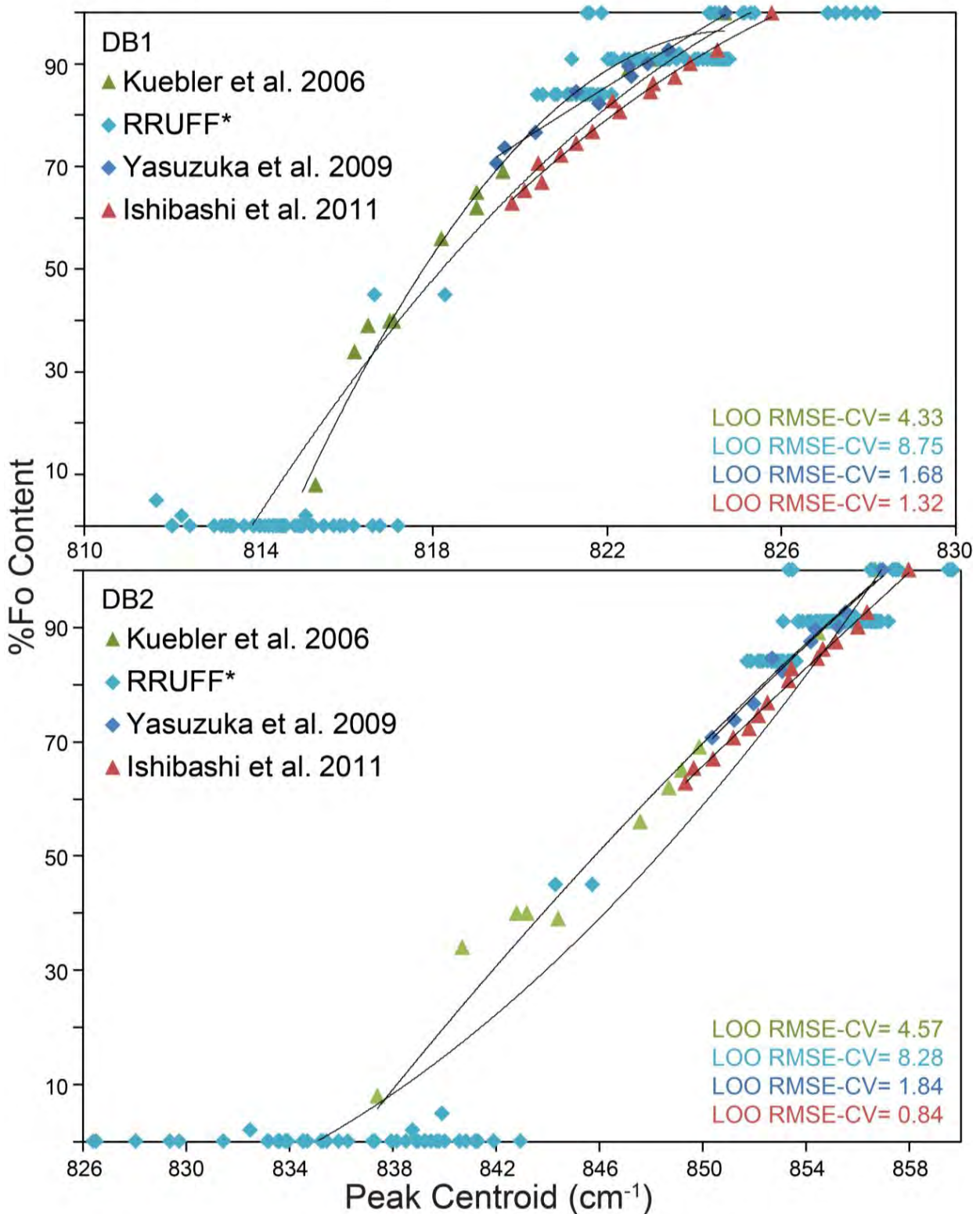
19
20 Figure 3. Histogram of 93 synthetic (blue) and natural (red) samples on the Fo-Fa series. Natural
21 olivines typically form with a %Fo of ~89.5 resulting in an unbalanced distribution on the Fo-Fa
22 series.

23



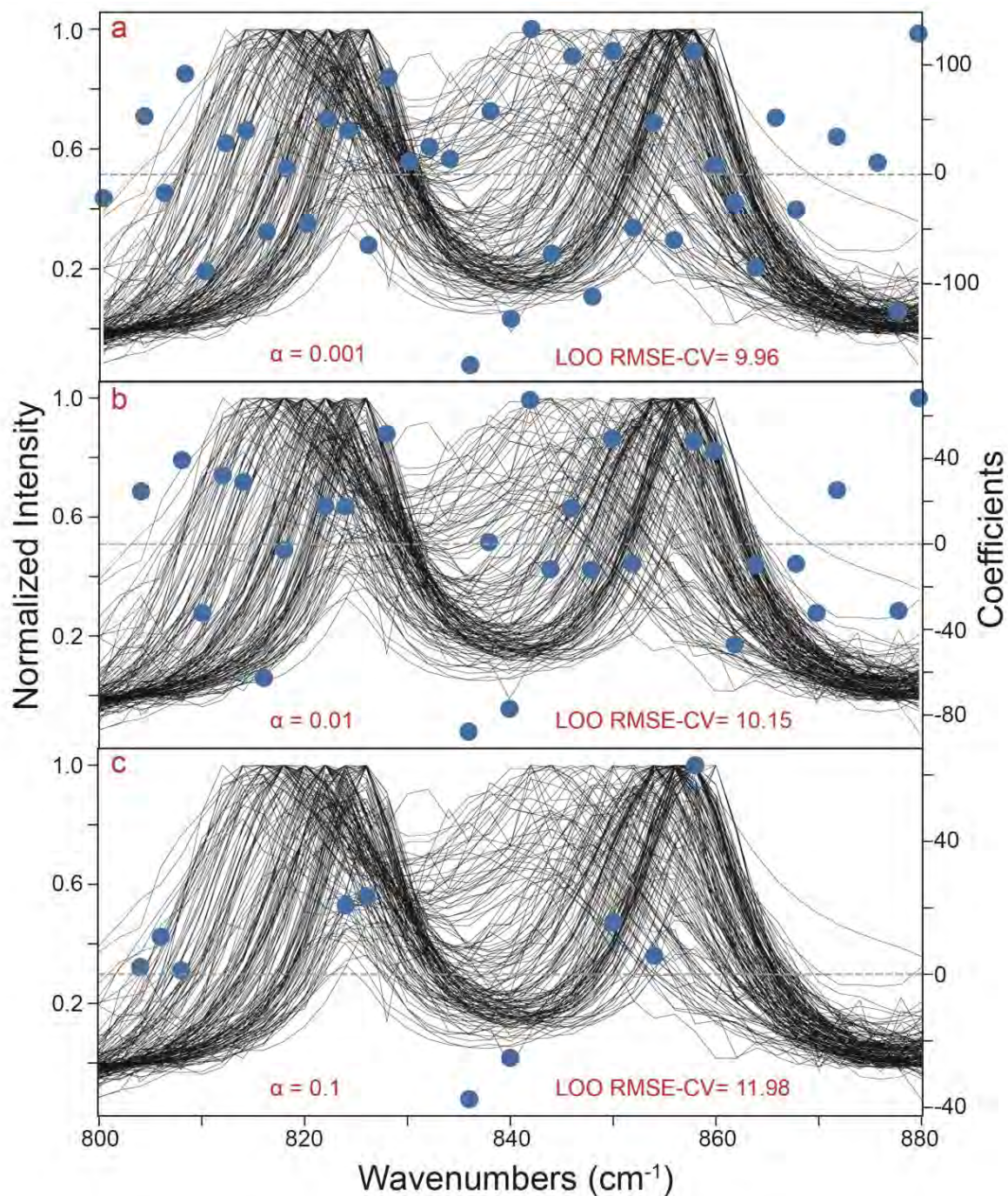
24

25 Figure 4. Fo by EMPA versus peak centroids of DB1 and DB2. Second order polynomial fits
26 and RSME-CV values are included for the unaligned data acquired on Bruker's BRAVO ($n=25$)
27 and Senterra ($n=68$) spectrometers. Error bars are smaller than the symbols and are given in
28 Tables S3 and S4.

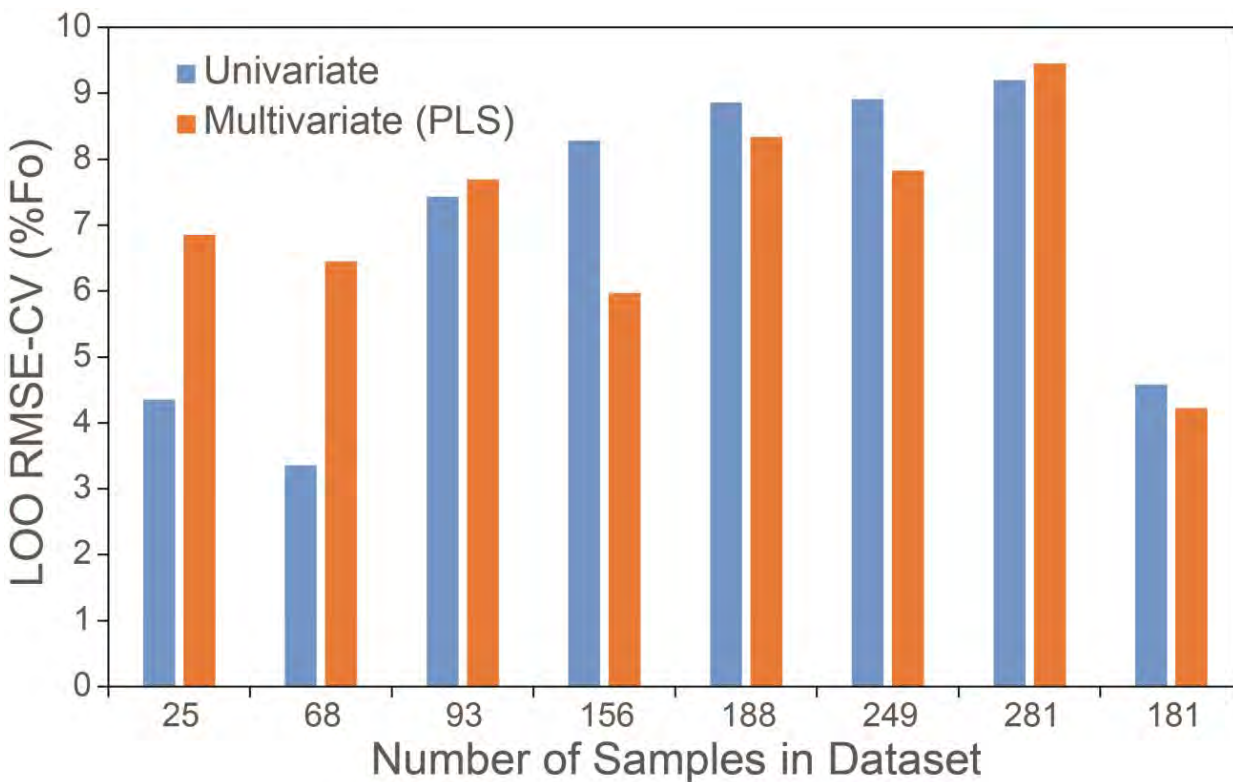


29
30
31
32
33

Figure 5. %Fo by EMPA versus peak centroids positions of (a) DB1 and (b) DB2 for data acquired in other studies. Including results from the RRUFF* database (see text for explanation of notation), and studies by Kuebler et al. (2006), Yasuzuka et al. (2009), and Ishibashi et al. (2011).



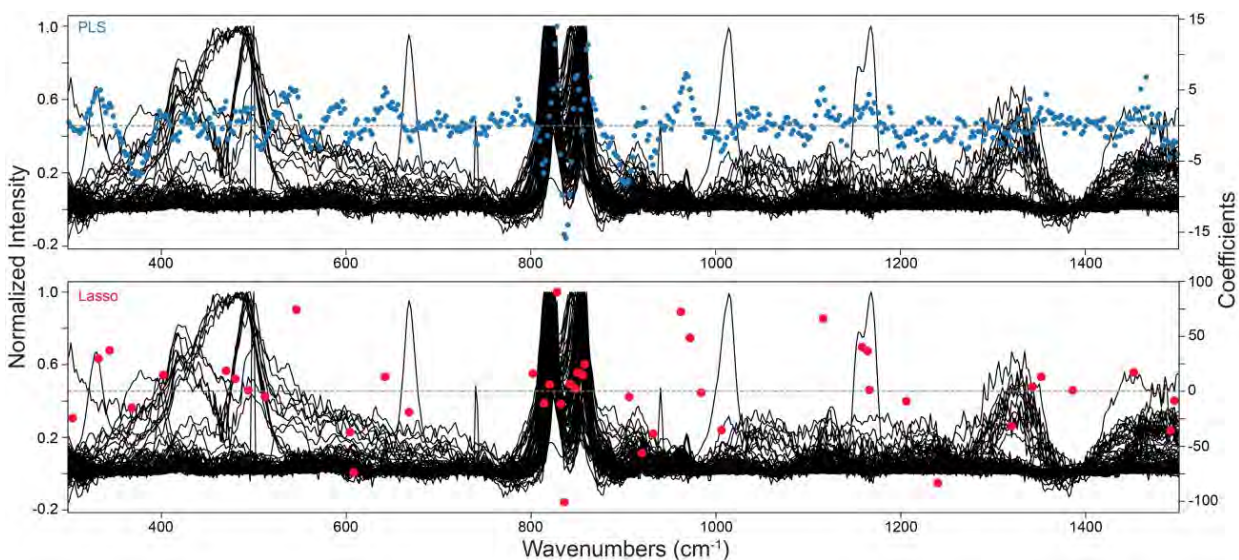
34
35 Figure 6. Variations in Lasso model accuracy as a function of the number of coefficients. As α
36 increases, fewer channels are examined: (a) 38 channels for $\alpha = 0.001$, (b) 27 channels for $\alpha =$
37 0.01, and (c) 10 channels for $\alpha = 0.1$. As the number of channels examined is decreased (fewer
38 coefficients within the model), the RMSE-CV of the models increases in value (gets worse). This
39 demonstrates the value of a models that examines a high number of channels, which is achieved
40 in a small α value Lasso model or PLS models.
41



42

43 Figure 7. Bar graph comparing results from Tables 4 and 5. When the dataset is small and/or all
44 the data are acquired on the identical instrument, then univariate methods produce better results
45 than those using multivariate analyses. However, as the number of samples and instruments used
46 increase, PLS methods generally produce more accurate results.

47



48

49 Figure 8. (top) Plot of BRAVO and Senterra unaligned data in Table 6, along with circles
50 indicating the magnitude of PLS coefficients (right-hand y axis units). Note that PLS
51 coefficients are proportional to spectral intensity at each wavenumbers, so absolute values cannot
52 be compared on this plot. However, the PLS coefficients do demonstrate that the entire
53 wavenumber range contains information useful to predicting composition. (bottom) Analogous
54 plot for the same data but showing Lasso coefficients for a model with $\alpha = 0.001$.