



**HAL**  
open science

# The SoilExp software: An open-source Graphical User Interface (GUI) for post-processing spatial and temporal soil surveys

G. Boudoire, M. Liuzzo, S. Cappuzzo, G. Giuffrida, P. Cosenza, A. Derrien,  
E.E. Falcone

## ► To cite this version:

G. Boudoire, M. Liuzzo, S. Cappuzzo, G. Giuffrida, P. Cosenza, et al.. The SoilExp software: An open-source Graphical User Interface (GUI) for post-processing spatial and temporal soil surveys. *Computers & Geosciences*, 2020, 142, pp.104553. 10.1016/j.cageo.2020.104553 . hal-02901361

**HAL Id: hal-02901361**

**<https://uca.hal.science/hal-02901361>**

Submitted on 4 Sep 2020

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

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

1 **The SoilExp software: an open-source Graphical User Interface**  
2 **(GUI) for post-processing spatial and temporal soil surveys**

3

4 **G. BOUDOIRE**<sup>1,2\*</sup>, M. LIUZZO<sup>1</sup>, S. CAPPUZZO<sup>1</sup>, G. GIUFFRIDA<sup>1</sup>, P. COSENZA<sup>1</sup>,  
5 A. DERRIEN<sup>3</sup>, E.E. FALCONE<sup>1</sup>

6

7 <sup>1</sup>Istituto Nazionale di Geofisica e Vulcanologia, Sezione di Palermo, Via Ugo La Malfa  
8 153, 90146 Palermo, Italy

9

10 <sup>2</sup>Université Clermont Auvergne, CNRS, IRD, OPGC, Laboratoire Magmas et Volcans,  
11 6 avenue Blaise Pascal, 63178 Aubière, France

12

13 <sup>3</sup>Observatoire Volcanologique du Piton de la Fournaise (OVPF), Institut de Physique  
14 du Globe de Paris (IPGP), Sorbonne Paris-Cité, UMR 7154 CNRS, Université Paris  
15 Diderot, Bourg Murat, France

16

17 \* Corresponding author. Present address: Laboratoire Magmas et Volcans – 6 avenue  
18 Blaise Pascal – 63170 Aubière (France). Telephone: +33 07 85 22 19 11. E-mail:  
19 guillaume.boudoire@uca.fr

20

21

22

23

24

G. Boudoire has developed the SoilExp software, and performed field and laboratory tests. M. Liuzzo has conceived and designed the MEGA instrument and performed field and laboratory tests. S. Cappuzzo has conceived and designed the MEGA instrument. G. Giuffrida has performed field and laboratory tests. P. Cosenza has contributed to the design of MEGA instrument. A. Derrien has performed field tests and E.E. Falcone contributed to the writing of the manuscript.

## 25 **Abstract**

26 Preliminary interpretation of geological processes during field measurement campaigns  
27 require fast data analysis to adapt ongoing target strategies. It is the case of soil  
28 investigations where coupling geochemical and geophysical records favor a better  
29 understanding of subsurface processes. This task requires (i) statistical analysis ~~is~~  
30 ~~needed~~ to identify areas of interest during spatial surveys and (ii) signal processing ~~is~~  
31 ~~required~~ to analyze temporal series.

32 Here we present SoilExp, an open-source Python-based Graphical User  
33 Interface (GUI) that permits to process spatial and temporal surveys of soil gases (e.g.  
34 soil CO<sub>2</sub> flux) combined with common physical parameters (e.g. self-potential,  
35 temperature) that are synchronously recorded on the field. SoilExp mixes innovative  
36 algorithms with the more common tools used for the analysis of both spatial surveys or  
37 temporal series. It offers the possibility to display distribution plots, maps, comparative  
38 plots, spectra and spectrograms, as well as data statistical analysis, in order to deal  
39 efficiently with datasets acquired on the field. Field measurements performed at  
40 Stromboli (Italy) supports that such software solution facilitates a quick visualization  
41 of the data output and is a powerful tool on the geochemical and geophysical analysis.

42

## 43 **Keywords**

44 Geo-spatial survey, time series, soil CO<sub>2</sub> degassing, self-potential, Python

45

## 46 **1. Introduction**

47 Identifying hidden geologic structures and studying gas and hydrothermal fluid  
48 circulation within the ground is of first interest in many disciplines as agriculture  
49 (Kucera & Kirkham, 1971), mineral resources (Hinkle & Dilbert, 1984; Lovell *et al.*,

50 1983), geothermy (Chiodini *et al.*, 2001, 2005), geological storage (Sandig *et al.*, 2014)  
51 and natural hazards (Allard *et al.*, 1991; Finizola *et al.*, 2002; Hernandez *et al.*, 2001;  
52 Irwin & Barnes, 1980). Coupling geochemical and geophysical records has  
53 demonstrated a real complementarity to characterize soil heterogeneities and related  
54 fluid circulations (Aubert *et al.*, 1984; Boudoire *et al.*, 2018; Elskens *et al.*, 1964;  
55 Finizola *et al.*, 2003; Gaudin *et al.*, 2015; Giammanco *et al.*, 1997). In particular,  
56 diffusive CO<sub>2</sub> degassing (CO<sub>2</sub>), self-potential (SP) and temperature (T) measurements  
57 are among the most common methods used by the scientific and industrial community  
58 to perform both spatial surveys or temporal records for monitoring purposes (Boudoire  
59 *et al.*, 2018; Byrdina *et al.*, 2012; Finizola *et al.*, 2003; Gresse *et al.*, 2016; Pearson *et*  
60 *al.*, 2008).

61         These measurements often require (i) the use of self-alone instruments and (ii)  
62 a preliminary data treatment to be used reliably. For instance, (i) measurements of  
63 diffusive CO<sub>2</sub> degassing (CO<sub>2</sub>) may require the use of a stainless steel probe (active  
64 method) or an accumulation chamber (passive method) connected to infrared  
65 spectrometers, self-potential may require the use of non-polarizable Cu/CuSO<sub>4</sub>  
66 electrodes coupled with a high impedance voltmeter and temperature (T) measurements  
67 may be performed with K-type thermal probes and a digital thermometer or with a  
68 pyrometer (Finizola *et al.*, 2010). Additionally (ii), spatial surveys often need a quick  
69 first idea of results during the daily performed acquisitions in order to identify the main  
70 areas of interests and eventually adapt or correct the ongoing fieldwork strategy  
71 (Chatterjee *et al.*, 2019). Meanwhile, temporal series are often subjected to an  
72 environmental influence that needs to be corrected before an accurate use of the signals  
73 as regression, moving average or cut-band filter for the most common ones (Boudoire  
74 *et al.*, 2017a; Liuzzo *et al.*, 2013; Padron *et al.*, 2008; Viveiros *et al.*, 2008). Many  
75 industrial software packages or homemade codes are able to deal efficiently with this

76 kind of data but often required to be used additionally to cover the whole range of  
77 expected common data treatment tools (with various file formatting). It is often time  
78 consuming and limit a fast and efficient evaluation of the datasets.

79 Here we present a new user-friendly Python-based GUI (Graphical User  
80 Interface) software: Soil Exploration (SoilExp). SoilExp is able to analyze both spatial  
81 and temporal datasets obtained on the field and respecting some file formatting rules.  
82 The final aim of SoilExp is to provide to the geologic-environmental researchers  
83 community both innovative and classical tools for a first data processing: (i) data  
84 correction (linear regression, moving average, cut-band filter), (ii) data analysis  
85 (statistical analysis, populations identification), (iii) data comparison (correlations,  
86 cross-correlations) and, (iv) graphical representation (distribution plots, comparative  
87 plots, spectra, spectrograms, maps). To illustrate the potentiality of SoilExp to sustain  
88 field surveys and to address scientific issues, both soil CO<sub>2</sub> flux and self-potential  
89 measurements were performed at Stromboli (Italy). Results are presented in a final  
90 section and discussed with respect to those obtained from previous field surveys.

91

## 92 **2. Overview on the SoilExp software**

93 The SoilExp 1.0 software distribution is written in Python 2.7 (Fig. 1). The Graphical  
94 User Interface (GUI) is based on the Tkinter library. It requires the following libraries:  
95 Pandas, Numpy, SciPy, Matplotlib, Scikit-Learn, PySerial is required. Thus, as  
96 processed during SoilExp 1.0 development, we recommend to the user to install the  
97 Anaconda distribution on their machines in order to benefit of the Spyder open source  
98 cross-platform integrated development environment (IDE) with scientific libraries. Full  
99 details are provided in the user guide. Information related to the installation of the  
100 software distribution (SoilExp 1.0), to its step-by-step use and to potential script  
101 modifications are reported in the associated user manual.

102           Indeed, in this study, we focus on the main functionalities provided by the  
103 SoilExp 1.0 distribution. These main functionalities are exposed through three  
104 independent scripts described in the following parts (Fig. 2). The first script is dedicated  
105 to save/reset field data from the MEGA (Multisensors Electrical and Gas Analyzer)  
106 instrument and calibrate its sensors using an USB-Serial connection (so-called “Serial”  
107 option in the following parts) (Fig. 3a). The second script is dedicated to the analysis  
108 of spatial surveys (so-called “Space” option in the following parts) (Fig. 3b). The third  
109 script is dedicated to time series processing (so-called “Time” option in the following  
110 parts) (Fig. 3c).

111           The “Serial” option is dedicated for applications based on the use of the MEGA  
112 instrument that has been entirely conceived and designed at the INGV of Palermo by  
113 two of the current authors (Liuzzo & Cappuzzo). The MEGA instrument is not the focus  
114 of this paper and will be better presented to the community in future specific  
115 contributions. The other two options are composed of four panels (Fig. 3): (i) the first  
116 one (e.g. ‘1. File Treatment’) is used to treat raw data files and create intermediate  
117 formatted files; (ii) the second one (e.g. ‘2. Data Processing’) is used to select a  
118 parameter of interest from an intermediate formatted file, modify its corresponding  
119 series (correction, filtering, average) and, display resulted plots (for the “Time” option);  
120 (iii) the third one (e.g. ‘3. Data Analysis’) is dedicated to show the results of the data  
121 analysis as correlations, cross-correlations, statistics, populations identification and  
122 plots (for the “Space” option); (iv) the last one (e.g. ‘4. Save’) is used to save final  
123 processed datasets and related information (as populations) in .csv files.

124

## 125 **3. Main functionalities**

### 126 **3.1. Initialization**

127 The initialization step (first panel of the “Time” and “Space” options) aims at  
128 converting raw data files to intermediate files that may be manipulated in the other  
129 panels (Fig. 3b, c). Raw data files are “.csv” files downloaded from the MEGA  
130 instrument or created by the user in a compatible format to be correctly processed (see  
131 the user manual).

132 In the “Space” and “Time” options (Fig. 3b, c) the user can choose either the  
133 instrumental calibration by default or new calibration parameters to recalculate data  
134 series. Data reduction then performed in order to identify internal errors (typographical  
135 errors) or unreliable measurements, i.e. out of the range of values defined for the battery  
136 voltage, the pump flux and the horizontal dilution of precision (HDOP) of the global  
137 positioning system (GPS). These limit values used to define outliers are set by default  
138 but may be changed by the user directly in the GUI. Rows containing bad values are  
139 either linearly interpolated with the “Time” option in order to correctly apply further  
140 time series analysis (in this case interpolated rows are kept in memory in order to be  
141 removed in final .csv files) or left as empty rows with the “Space” option. Additionally,  
142 as soil CO<sub>2</sub> measurements may be acquired with the accumulation chamber method in  
143 the “Space” option, an additional filter is applied on the r-squared value of computed  
144 soil CO<sub>2</sub> flux. In this case, values of soil CO<sub>2</sub> flux out of the r-squared range will be  
145 considered as outliers values and thus set at 0 gm<sup>-2</sup>d<sup>-1</sup> (no flux). In the case where soil  
146 CO<sub>2</sub> flux has to be calculated from the dynamic concentration method, the software  
147 integrates the possibility in both options to convert the CO<sub>2</sub> %molar contents in flux  
148 (Camarda *et al.*, 2006; Gurrieri and Valenza, 1988; Liuzzo *et al.*, 2015). The conversion  
149 is made by the use of the equation of Camarda *et al.* (2006) that takes into account the  
150 soil permeability value defined by the user.

151 Once the raw data file is cleaned, the last step will generate intermediate files  
152 that are of first interests to keep processing the data thanks to the other panels of the

153 GUI. The definition of the time-lag period (in seconds) is here fundamental. We define  
154 the time-lag as the period in seconds separating two independent series of  
155 measurements (i.e. the period during which the used instrument will be in stand-by).  
156 With the “Time” option, a new intermediate file will be created each time that two  
157 consecutive rows are separated by a duration greater than the time-lag. With the  
158 “Space” option, an unique intermediate file will be created with the median values of  
159 each independent series (the median value being here considered as representative of  
160 the acquisition to avoid the effect of potential spikes on the average). The time-lag must  
161 not be confused with the sampling rate that is necessary shorter and defines the period  
162 (in seconds) between two measurements within the same series. The sampling rate is  
163 used to correctly adapt the scale and the legend of plots (spectra, spectrogram) by  
164 converting a range of measurements (number of records) into a time range (number of  
165 seconds).

166

## 167 **3.2. The “Space” option: dealing with spatial surveys**

168 The “Space” option aims to propose innovative and classical tools to deal with spatial  
169 surveys, i.e. with datasets where each point of measurements is defined by distinct  
170 geographical coordinates (Fig. 3b).

171

### 172 ***3.2.1. Correlations and preliminary data correction***

173 The “Space” option allows coefficients of linear correlations to be identified (slope,  
174 offset, r-squared) between the parameter of interest (e.g. CO<sub>2</sub>) and other records (e.g.  
175 temperature, pressure, wind speed) obtained by the user (Fig. 3b). In particular, these  
176 coefficients are often useful during soil surveys, where records can be affected by  
177 external parameters. For instance, soil CO<sub>2</sub> flux may be slightly dependent on pressure



178 (Barde-Cabusson *et al.*, 2009; Liuzzo *et al.*, 2013; Viveiros *et al.*, 2008). The equation  
179 used for the linear regression is the following:

180

$$181 \quad \text{Signal}_{Clr}(x,y) = \text{Signal}_{Raw}(x,y) - (a \times (\text{Param}_{Raw}(x,y) - \text{Param}_{Avg}))$$

182

183 where  $\text{Signal}_{Clr}(x,y)$  is the value of the parameter of interest at the geographical position  
184  $(x,y)$  after the correction by linear regression,  $\text{Signal}_{Raw}(x,y)$  is the value before the  
185 correction,  $a$  is the slope of the linear correlation,  $\text{Param}_{Raw}(x,y)$  is the parameter used  
186 to performed the regression at the geographical position  $(x,y)$  and,  $\text{Param}_{Avg}$  is the  
187 average of this parameter for the whole dataset in order to correct the offset linked to  
188 the correction. After the correction, coefficients of linear correlations are reprocessed  
189 and automatically updated in order to verify the efficiency of the correction and identify  
190 potential needs of further steps of correction.

191

### 192 **3.2.2. Statistical analysis**

193 In order to better constrain the data distribution of the parameter of interest, the “Space”  
194 option allows the user to display also a probability histogram together with the best fit  
195 line of the potential normal distribution (see Supplementary Material). The normality  
196 of the data series is tested via the Anderson-Darling normality test (SciPy library;  
197 Anderson and Darling, 1952; Stephens, 1974, 1976). In our case, i.e. making the  
198 assumption that both mean and variance are initially unknown, the Anderson-Darling  
199 normality test rejects the hypothesis of normality with a 95% significance level if  $A^2$   
200 (the squared of the test statistic A) exceeds 0.752 for data series owning more than 8  
201 samples (D’Agostino, 1986).

202 Together with the Anderson-Darling normality test, the “Space” option gives  
203 the opportunity to calculate some classical statistical values: mean, standard deviation,

204 median, minimum, maximum, kurtosis, skewness (Fig. 3b). Here we focus on the last  
205 two indicators less common for non-regular users of statistical tools. The kurtosis is a  
206 measure of the tailedness of the probability distribution of a random variable, i.e.  
207 describing the shape of the probability distribution (Zwillinger and Kokoska, 2000).  
208 Using the Fisher's definition, normally distributed data should provide a result of 0.  
209 Skewness is the measure of the asymmetry of the probability distribution of a random  
210 variable with respect to its mean. The skewness, which could be either positive or  
211 negative, should be about 0 for normally distributed data (Zwillinger and Kokoska,  
212 2000).

213         If the indicators described above (skewness, kurtosis, Anderson-Darling  
214 normality test) do not argue in favor of a normal distribution, it may be due to the  
215 presence of more than one population in the data series. Indeed, during spatial survey,  
216 one subject of major interest is often to discriminate the different populations that  
217 contribute to the data series. This is crucial in order to recognize the existence of distinct  
218 sources as, for example, biogenic or magmatic ones for soil CO<sub>2</sub> flux in volcanic context  
219 (Boudoire *et al.*, 2017b; Liuzzo *et al.*, 2015; Viveiros *et al.*, 2008). In order to address  
220 this specific issue, we have developed a new algorithm in SoilExp able to combine the  
221 two statistical methods more used in environmental scientific research, which can  
222 distinguish various populations from log-normally distributed data (Fig. 4a, b). The first  
223 is the graphical method based on probability plots known as Sinclair method (Chiodini  
224 *et al.*, 1998; Giammanco *et al.*, 2010; Sinclair, 1974); the second is the maximum-  
225 likelihood numerical method based on the use of Gaussian Mixture Model (GMM)  
226 implementing an expectation-maximization (EM) algorithm (Benaglia *et al.*, 2009;  
227 Boudoire *et al.*, 2018; Elio *et al.*, 2016). The Sinclair method provides an user-friendly  
228 view of the populations and mixed values, however, it has two main shortcomings. One  
229 is related to the low accuracy for datasets counting less than 100 values (Sinclair, 1974).

230 The second limitation is related to the difficulty to precisely estimate the confidence  
231 intervals. These problems are solved using the maximum-likelihood (ML) method that  
232 fits finite mixtures of normal distributions: we have implemented a Scikit-Learn-based  
233 algorithm that simulates such fitting with 1 to 10 populations with 1000 iterations for  
234 each simulation. The best simulation is then selected based on the value of the Bayesian  
235 Information Criterion (BIC) developed for model selection among finite set of  
236 simulations (Ghosh *et al.*, 2006) and displayed on the GUI (Fig. 3b) . Finally, for each  
237 value of the data serie, the algorithm predicts the probability that the value belongs to  
238 one of the defined populations. In our case, we have considered that if one value shows  
239 a probability to be defined by a single population greater than 95% thus it will be  
240 considered as part of this population. If not, this value is considered as an intermediate  
241 value (or mixed value) between the two neighboring populations. A the end, the  
242 algorithm allows the user to automatically see the result of this ML-based partitioning  
243 of the values on probability plots (Fig. 4a, b). Furthermore, the users can simulate  
244 different partitioning by modifying the number of inferred populations directly on the  
245 GUI (Fig. 3b), if the first step of differentiation is not satisfying.

246

### 247 **3.2.3. Mapping**

248 After having performed data correction and statistical analysis, it is possible to obtain  
249 a first idea of the two-dimensions (2D) distribution of the data (Fig. 4c, d). Our aim is  
250 not to develop complex interpolating algorithms for which many software are already  
251 built. Here we propose a simple graphical representation of the data through two distinct  
252 maps. The first one uses a simple color gradient to show the 2D evolution of the values.  
253 The second one is more innovative (presented on Fig. 4 for SP and CO<sub>2</sub> data obtained  
254 at Stromboli), meaning that the map-builder takes into consideration the results of the  
255 population analysis described above, generating and displaying a repartition of the

256 values between the different populations (and related mixing values). If an internet  
257 connection and an API key are available  
258 (<https://developers.google.com/maps/documentation/javascript/get-api-key>), a  
259 background satellite map will be automatically downloaded and georeferenced from  
260 the Google Maps Platform. If not, the background will remain neutral. However, the  
261 upper left box (Fig. 4c, d) highlights the coordinates of the corners to facilitate the  
262 extraction of an adequate background map from other sources.

263

### 264 **3.3. The “Time” option: processing time series**

265 The “Time” option aims to propose classical tools to deal with time series, i.e. with  
266 datasets where the measurements have specific frequency (here defined as the sampling  
267 rate) (Fig. 3c).

268

#### 269 ***3.3.1. Correlations and cross-correlations***

270 “Time” option allows to identify coefficients of linear correlations (slope, offset, r-  
271 squared) between the parameter of interest and other records, where the control panel  
272 (Fig. 3c) is similar to the one in "Space" option. Sometimes some signals may have a  
273 time delay between them, which can be attributed either to an instrumental lag or to an  
274 effect caused by a natural phenomenon. To take into account these effects, we have  
275 implemented a SciPy-based algorithm to calculate the cross-correlations between each  
276 parameters. The algorithm couples complex-valued functions with conjugates and Fast  
277 Fourier Transform (FFT) to numerically determine both lags and r-squared values  
278 between time series. Best results are shown in the table of the “Time” option GUI (Fig.  
279 3c).

280

#### 281 ***3.3.2. Signal processing***

282 The “Time” option gives the possibility to the user to apply three of the most common  
 283 signal processing tools used in the geo-scientific community: (i) linear regression, (ii)  
 284 moving average and, (iii) cut band filter.

285 The linear regression method is the same than in the “Space” option and only  
 286 require to select the parameter used for the regression and to compute the corresponding  
 287 coefficients. This method is used to remove short-term environmental influence on  
 288 geochemical and geophysical signals (Boudoire *et al.*, 2017a; Liuzzo *et al.*, 2013).

289 The moving average method is a type of finite impulse response filter used to  
 290 smooth out short-term signal variations. This method performs an average on a defined  
 291 subset of the data series, then shifts forward to repeat the calculations, excluding the  
 292 first value of the previous subset and including the next one. Using the convolution  
 293 operator of the Numpy library, we have implemented a simple moving average method,  
 294 i.e. giving the same weight to each value  $a_j$ :

295

$$296 \quad movav(a_i) = \left(\frac{1}{k}\right) \times \sum_{j=i-k/2}^{i+k/2} a_j \text{ for } i \in \left] \frac{k}{2}, \frac{n-k}{2} \right[$$

297

298 where  $i$  is the position of the value  $a_i$  in the data series on which the moving average is  
 299 applied,  $n$  the length of the data series and  $k$  the size of the subset. To deal with border  
 300 effects (i.e. when the number of available values to perform the moving average is  
 301 lower than the size of the defined subset), we have adapted the convolution to the  
 302 number of available values:

303

$$304 \quad movav(a_i) = \left(\frac{1}{i}\right) \times \sum_{j=0}^i a_j \text{ for } i \in \left[ 0, \frac{k}{2} \right]$$

305

306 
$$movav(a_i) = \left(\frac{1}{n-i}\right) \times \sum_{j=i}^n a_j \text{ for } i \in \left[\frac{n-k}{2}, n\right]$$

307

308 To enhance the reliability of the calculations linked to correlations and cross-  
309 correlations, the moving average method is applied to all data series when computed.

310 Finally, to treat long-term signal variations, we have used the Fast Fourier  
311 Transform (FFT) package of the SciPy library to develop a cut (or block) band filter.  
312 This filter removes from the signal spectra (cf. *fft*) the frequencies belonging to an  
313 interval defined by the user before making the inverse operation to rebuild the signal  
314 (cf. *ifft*).

315

### 316 **3.3.3. Graphical representation**

317 When pressing the plot-related buttons of the “Time” option, the user automatically  
318 applies the correction and filtering methods that has been defined previously (Fig. 2).

319 Consequently, the user may decide to perform several combination between the  
320 signals which is intended to compare:

- 321 (i) Compare the raw signal with the new corrected and filtered signal, and  
322 eventually reinitialize the signal to apply a distinct protocol (Fig. 5a). Based  
323 on the same statistical algorithm used with the “Space” option to  
324 characterize populations, we have implemented an option allowing the user  
325 to directly show on the plot the values belonging to the “highest” population  
326 (often considered as representative of anomalous values with respect to the  
327 background; Boudoire *et al.*, 2017a; Liuzzo *et al.*, 2013, Liuzzo *et al.*, 2015);
- 328 (ii) Compare the treated signal with another signal of interest (Fig. 5b). This  
329 plot may be particularly useful to investigate well correlated or cross-  
330 correlated signals;

- 331 (iii) See the FFT spectrum on which are displayed the three greatest frequencies  
332 (Fig. 5c). Thanks to the labels indicating the corresponding number of  
333 measurements, the user may define the frequency interval on which  
334 applying the cut-band filter;
- 335 (iv) See the corresponding spectrogram that is a different visual representation  
336 of the FFT spectrum, extensively used in geophysical signal processing (Fig.  
337 5d). It is particularly useful to detect periodic components and signal  
338 perturbations that may affect all frequencies. Here we use the ‘*specgram*’  
339 function of the Matplotlib library with a linear detrend and a magnitude  
340 mode of 256 NFFT of default (Nonequispaced Fast Fourier Transform: the  
341 number of points in each processed block) and a 128 noverlap (the number  
342 of points of overlap between processed blocks). The user is free to modify  
343 these parameters directly in the Python 2.7 script (see user manual).

344 The signal analysis depends on the sampling rate, therefore we cannot use an unique  
345 legend for spectrum and spectrogram axes. Consequently, we have adapted the  
346 algorithms to show both the results of the raw signal analysis (in term of number of  
347 measurements) and their meaning using more classical units. For the last one, we have  
348 coupled the number of measurements and the sampling rate to have a real temporal  
349 scale (i) in seconds (between parenthesis) on the spectrum and (ii) in hertz on the  
350 spectrogram.

351

### 352 **3.4. Saving and exporting results**

353 The SoilExp software gives the opportunity to save every graphical object with  
354 different extensions (.png, .eps ...), which can be easily further modified later.

355 Additionally both “Space” and “Time” GUI options have dedicated buttons to  
356 save .csv files. In the “Space” option, the final .csv file is similar to the intermediate

357 file but takes into account the results of the linear regressions that could be applied to  
358 correct the dataset. Additionally, it is possible to save a .csv file recording the data  
359 repartition between the defined populations and mixing groups. Both files aim to be  
360 eventually further processed through software dedicated to complementary and more  
361 specific tools as e.g. data interpolation, kriging, sequential Gaussian simulation (SGS).  
362 In the “Time” option, the final .csv file is also similar to the intermediate file but (i) has  
363 one supplemental column for the corrected and filtered data series and (ii) shows empty  
364 rows for missing values, which have been interpolated for the needs of signal  
365 processing. Such final file may be then processed through other complementary  
366 software for measurements of volcanic gas in plume or other environmental  
367 applications in atmospheric measurements (Fig. 1e, f).

368

#### 369 **4. SoilExp application: an example at Stromboli (Italy)**

370 In volcanic environment, two of the main goals of soil surveys are (i) the identification  
371 of volcano-tectonic structures (Giammanco et al., 1997; Finizola et al., 2002, 2010) and  
372 (ii) the characterization of hydrothermal fluid circulation (Revil et al., 2011; ; Boudoire  
373 et al., 2018). Once, because these low permeable structures may favor the ascent of  
374 magmatic fluids leading to fissural eruptions (Boudoire et al., 2017b). Moreover, such  
375 structural interfaces may raise important issues concerning soil stability and thus  
376 landslide outbreak (Neri et al., 2004). To test the ability of SoilExp to deal with such  
377 goals, we have performed a spatial soil survey at Stromboli (Sicily, Italy) by the mean  
378 of the MEGA instrumental kit (Fig. 1). Three transects were performed with a 20 m-  
379 spacing for a total of 45 measurements of soil CO<sub>2</sub> flux and self-potential (dataset  
380 available with our distribution as “intermediate” test file). Here, we focused on the first  
381 transect (14 measurements), the one on the northern flank of the volcano which is the  
382 closest to populated areas (Fig. 1c).



383 Data analysis performed with the “Space” option of SoilExp reveals (i) the  
384 absence of correlation between soil CO<sub>2</sub> ‘dynamic’ concentration (‘CO<sub>2</sub>\_10’) and the  
385 environmental parameter (pressure ‘P\_atm’, temperature ‘T\_atm’, humidity ‘Rh’)  
386 during the transect and (ii) an important correlation ( $R^2 = 0.79$ ) between soil CO<sub>2</sub>  
387 ‘dynamic’ concentration (‘CO<sub>2</sub>\_10’) and self-potential measurements (‘SP’).  
388 Consequently, no correction from the environmental influence was applied (Viveiros  
389 et al., 2008) and we focus on both soil CO<sub>2</sub> flux and self-potential measurements in the  
390 following parts. The analysis performed by SoilExp shows that soil CO<sub>2</sub> ‘dynamic’  
391 concentration (‘CO<sub>2</sub>\_10’) varies from 0.07 to 0.95 %. Self-potential (‘SP’) varies from  
392 -155 to +77 mV. The Anderson-Darling normally test gives  $A^2$  equal to 14.3 and 3.3  
393 for ‘CO<sub>2</sub>\_10’ and ‘SP’, respectively. These values are well above 0.752, and testify  
394 that both datasets do not present a normal distribution (at 95% of significance level). It  
395 means that these datasets are better explained by the presence of two or more  
396 populations. Actually, the new statistical algorithm developed in SoilExp highlights the  
397 presence of two populations of values for both parameters (Fig. 4a, b). Soil CO<sub>2</sub>  
398 ‘dynamic’ concentration shows the presence of two populations: one with high values  
399 (>0.20 % for 7.1% of the dataset; Fig. 4a) and the other with low values (<0.20 % for  
400 92.9% of the dataset). We applied the equation of Camarda et al. (2006) to convert soil  
401 CO<sub>2</sub> ‘dynamic’ concentration in soil CO<sub>2</sub> flux for a range of soil permeability between  
402 15 and 50, i.e. the most common values for volcanic soils (Camarda et al., 2006). The  
403 calculated upper limit of the population of low soil CO<sub>2</sub> flux does not exceed 42 gm<sup>-2</sup>d<sup>-1</sup>  
404 <sup>1</sup>. This value is in accordance with the definition of a “background” population  
405 characterized by low soil CO<sub>2</sub> flux and generally ascribed to the biological soil activity  
406 (Liuzzo et al., 2015; Boudoire et al., 2017b). Conversely, the population of higher soil  
407 CO<sub>2</sub> flux (up to 233 gm<sup>-2</sup>d<sup>-1</sup>) is consistent with a magmatic-hydrothermal origin of the  
408 released fluids (Giammanco et al., 1997; Liuzzo et al., 2015). Self-potential shows also

409 the presence of one population of high values (from -9 up to +77 mV for 14.3% of the  
410 dataset; Fig. 4b) whereas most of the dataset is defined by a population of more negative  
411 values (from -169 up to -100 mV for 85.7% of the dataset).

412 Interestingly, the map-building of the soil CO<sub>2</sub> ‘dynamic’ concentration (Fig.  
413 4c) and self-potential (Fig. 4d), based on this population analysis, shows that the  
414 population of high soil CO<sub>2</sub> ‘dynamic’ concentration spatially correlates with the high  
415 self-potential measurements. This positive correlation between the two parameters is  
416 consistent with an upward migration of hydrothermal fluids in a restricted part of the  
417 transect (<40 m-wide) as documented for other volcanic systems (Barde-Cabusson et  
418 al., 2009; Bennati et al., 2011). Actually, this restricted part of the transect is cut by the  
419 Nel Cannestrà eruptive fissure that is known representing a low permeability structure,  
420 in relation with N41° inferred regional faults, (Finizola et al., 2002, 2010; Carapezza et  
421 al., 2009). The identification and characterization of such structure that favors the  
422 ascent of magmatic fluids raise important civil protection issues (Boudoire et al.,  
423 2017b). Current monitoring is performed in this area by the Istituto Nazionale di  
424 Geofisica e Vulcanologia (INGV) (Carapezza et al., 2009).

425

## 426 **5. Conclusion**

427 In this work we presented an open-source Graphical User Interface (GUI) software,  
428 SoilExp, which is written in Python language and is able to provide statistical and  
429 spectrum analysis as well as several options on filtering and correcting analysis on  
430 records acquired during spatial/temporal surveys. The software is based on two main  
431 options. Firstly, the “Space” option, aims to display the main statistical indicators used  
432 to study spatial surveys, to test the normality of data series, to identify and define the  
433 populations constituting the dataset through an innovative algorithm, and to show  
434 results on satellite maps. The second one, the “Time” option, aims to process time series

435 through classical tools used in signal processing (linear regression, moving average,  
436 cut-band filter, cross-correlations) and in signal representations (scatter plots, spectra,  
437 spectrogram). Beyond facilitating the fast outcome from field surveys by offering  
438 filtering tools, graphical results and statistical analyses, SoilExp gives to the users the  
439 possibility to integrate all the results in a unique tool of elaboration, improving the  
440 research potential of the scientific community dealing with spatial and temporal soil  
441 surveys.

442

### 443 **Acknowledgments**

444 S. Gurrieri is acknowledged for constructive discussions and advices on the script. P.  
445 Boudoire, A. Finizola and T. Ricci are acknowledged for providing technical support  
446 crucial for instrumental tests. People proposing or commenting open source Python  
447 codes and solutions on forums are gratefully acknowledged. We thank the Associate  
448 Editor and Jean Vandemeulebrouck for their suggestions that greatly improved the  
449 clarity and quality of the paper. This work has been funded by the Fondo Sociale  
450 Europeo (PO FSE 2014-2020) in the frame of the project “Metodi di controllo  
451 geochimico e geofisico dei fenomeni naturali sul campo ed in laboratorio”. We also  
452 acknowledge the French government IDEX-ISITE initiative 16-IDEX-0001 (CAP 20-25).

453

### 454 **Computer code availability**

455 The SoilExp software including full scripts, user manual and example files may be  
456 freely downloaded from <https://github.com/FreeMindsObservatory/SoilExp> (main  
457 developer: Dr. Guillaume Boudoire; corresponding author). The first distribution  
458 started to be developed in 2017 (SoilExp 1.0 based on Python 2.7; 16.6 Mo), made  
459 available for MacOSX and Windows platforms. It may be modified by any contributor

460 according his needs thanks to a Creative Commons Attribution-NonCommercial-  
461 ShareAlike 4.0 International License (CC BY-NC-SA 4.0)  
462 (<https://creativecommons.org/licenses/by-nc-sa/4.0/>). The software requires the  
463 installation of the Anaconda distribution to use the Spyder open source cross-platform  
464 integrated development environment (IDE) integrated all required scientific libraries.  
465 Please contact the corresponding author for any support regarding the SoilExp software.  
466

## 467 **References**

468 Allard, P., Carbonnelle, J., Dajlevic, D., Le Bronec, J., Morel, P., Robe, M.C.,  
469 Maurenas, J.M., Faivre-Pierret, R., Martin, D., Sabroux, J.C., Zettwoog, P., 1991.  
470 Eruptive and diffuse emissions of CO<sub>2</sub> from Mount Etna. *Nature* 351(6325), 387-391.

471

472 Anderson, T.W., Darling, D. A., 1952. Asymptotic theory of certain "goodness-of-fit"  
473 criteria based on stochastic processes. *Annals of Mathematical Statistics* 23, 193–212,  
474 doi:10.1214/aoms/1177729437.

475

476 Aubert, M., Camus, G., Fournier, C. 1984. Resistivity and magnetic surveys in  
477 groundwater prospecting in volcanic areas—case history maar of Beaunit, Puy de  
478 Dome, France. *Geophysical prospecting* 32(4), 554-563.

479

480 Barde-Cabusson, S., Finizola, A., Revil, A., Ricci, T., Piscitelli, S., Rizzo, E., Angeletti,  
481 B., Balasco, M., Bennati, L., Byrdina, S., Carzaniga, N., Crespy, A., Di Gangi, F.,  
482 Morin, J., Perrone, A., Rossi, M., Roulleau, E., Suski, B., Villeneuve, N., 2009. New  
483 geological insights and structural control on fluid circulation in La Fossa cone  
484 (Vulcano, Aeolian Islands, Italy). *J. Volcanol. Geotherm. Res.* 185, 231–245, doi:  
485 10.1016/j.jvolgeores.2009.06.002.

486

487 Benaglia, T., Chauveau, D., Hunter, D., Young, D., 2009. mixtools : An R package for  
488 analyzing finite mixture models. *Journal of Statistical Software* 32(6), 1-29.

489

490 Bennati, L., Finizola, A., Walker, J.A., Lopez, D.L., Higuera-Diaz, I.C., Schütze, C.,  
491 Barahona, F., Cartagena, R., Conde, V., Funes, R., 2011. Fluid circulation in a complex  
492 volcano-tectonic setting, inferred from self-potential and soil CO<sub>2</sub> flux surveys: the  
493 Santa María–Cerro Quemado–Zunil volcanoes and Xela caldera (Northwestern  
494 Guatemala). *J. Volcanol. Geotherm. Res.* 199 (3–4), 216–229.  
495 <http://dx.doi.org/10.1016/j.jvolgeores.2010.11.008>.

496

497 Boudoire, G., Di Muro, A., Liuzzo, M., Ferrazzini, V., Peltier, A., Gurrieri, S., Michon,  
498 L., Giudice, G., Kowalski, P., Boissier, P., 2017a. New perspectives on volcano  
499 monitoring in a tropical environment: continuous measurements of soil CO<sub>2</sub> flux at  
500 Piton de la Fournaise (La Réunion Island, France). *Geophysical Research Letters*  
501 44(16), 8244-8253.

502

503 Boudoire, G., Finizola, A., Di Muro, A., Peltier, A., Liuzzo, M., Grassa, F., Delcher,  
504 E., Brunet, C., Boissier, P., Chaput, M., Ferrazzini, V., Gurrieri, S., 2018. Small-scale  
505 spatial variability of soil CO<sub>2</sub> flux: Implication for monitoring strategy. *Journal of*  
506 *Volcanology and Geothermal Research* 366, 13-26.

507

508 Boudoire, G., Liuzzo, M., Di Muro, A., Ferrazzini, V., Michon, L., Grassa, F., Derrien,  
509 A., Villeneuve, N., Bourdeu, A., Brunet, C., Giudice, G., Gurrieri, S., 2017b.  
510 Investigating the deepest part of a volcano plumbing system: evidence for an active

511 magma path below the western flank of Piton de la Fournaise (La Réunion Island).  
512 Journal of Volcanology and Geothermal Research 341, 193-207.  
513  
514 Byrdina, S., Rücker, C., Zimmer, M., Friedel, S., Serfling, U., 2012. Self potential  
515 signals preceding variations of fumarole activity at Merapi volcano, Central Java.  
516 Journal of Volcanology and Geothermal Research 215, 40-47.  
517  
518 Camarda, M., Gurrieri, S., Valenza, M., 2006. In situ permeability measurements based  
519 on a radial gas advection model : Relationships between soil permeability and diffuse  
520 CO<sub>2</sub> degassing in volcanic areas. Pure and applied geophysics 163(4), 897-914.  
521  
522 Carapezza, M.L., Ricci, T., Ranaldi, M., Tarchini, L., 2009. Active degassing structures  
523 of Stromboli and variations in diffuse CO<sub>2</sub> output related to the volcanic activity.  
524 Journal of Volcanology and Geothermal Research 182(3-4), 231-245.  
525  
526 Chatterjee, S., Deering, C.D., Waite, G.P., Prandi, C., Lin, P., 2019. An adaptive  
527 sampling strategy developed for studies of diffuse volcanic soil gas emissions. Journal  
528 of Volcanology and Geothermal Research.  
529  
530 Chiodini, G., Cioni, R., Guidi, M., Raco, B., Marini, L., 1998. Soil CO<sub>2</sub> flux  
531 measurements in volcanic and geothermal areas. Applied Geochemistry 13(5), 543-  
532 552.  
533  
534 Chiodini, G., Frondini, F., Cardellini, C., Granieri, D., Marini, L., Ventura, G., 2001.  
535 CO<sub>2</sub> degassing and energy release at Solfatara volcano, Campi Flegrei, Italy. Journal of  
536 Geophysical Research : Solid Earth 106(B8), 16213-16221.

537

538 Chiodini, G., Granieri, D., Avino, R., Caliro, S., Costa, A., Werner, C., 2005. Carbon  
539 dioxide diffuse degassing and estimation of heat release from volcanic and  
540 hydrothermal systems. *Journal of Geophysical Research : Solid Earth* 110(B8).

541

542 D'Agostino, R.B., 1986. Tests for the Normal Distribution. In *Goodness-of-Fit*  
543 *Techniques*. New York, Marcel Dekker, ISBN 0-8247-7487-6.

544

545 Elio, J., Ortega, M.F., Nisi, B., Mazadiego, L.F., Vaselli, O., Caballero, J., Chacon, E.,  
546 2016. A multi-statistical approach for estimating the total output of CO<sub>2</sub> from diffusive  
547 soil degassing by the accumulation chamber method. *International Journal of*  
548 *Greenhouse Gas Control* 47, 351-363.

549

550 Elskens, I., Tazieff, H., Tonani, F., 1964. A new method for volcanic gas analyses in  
551 the field. *Bulletin Volcanologique* 27(1), 347-350.

552

553 Finizola, A., Ricci, T., Deiana, R., Cabusson, S.B., Rossi, M., Praticelli, N., Giocoli,  
554 A., Romano, G., Delcher, E., Suski, B., Revil, A., Menny, P., Di Gangi, F., Letort, J.,  
555 Peltier, A., Villasante-Marcos, V., Douillet, G., Avard, G., Lelli, M., 2010. Adventive  
556 hydrothermal circulation on Stromboli volcano (Aeolian Islands, Italy) revealed by  
557 geophysical and geochemical approaches: implications for general fluid flow models  
558 on volcanoes. *Journal of Volcanology and Geothermal Research* 196(1-2), 111-119.

559

560 Finizola, A., Sortino, F., Lénat, J.F., Valenza, M., 2002. Fluid circulation at Stromboli  
561 volcano (Aeolian Islands, Italy) from self-potential and CO<sub>2</sub> surveys. *Journal of*  
562 *Volcanology and Geothermal Research* 116(1), 1-18.

563

564 Gaudin, D., Finizola, A., Delcher, E., Beauducel, F., Allemand, P., Delacourt, C.,  
565 Brothelande, E., Peltier, A., Di Gangi, F., 2015. Influence of rainfalls on heat and steam  
566 fluxes of fumarolic zones: Six months records along the Ty fault (Soufrière of  
567 Guadeloupe, Lesser Antilles). *Journal of Volcanology and Geothermal Research* 302,  
568 273-285.

569

570 Ghosh, J.K., Delampady, M., Samanta, T., 2006. *An Introduction to Bayesian Analysis:*  
571 *Theory and Methods*. Springer-Verlag, New York.

572

573 Giammanco, S., Gurrieri, S., Valenza, M., 1997. Soil CO<sub>2</sub> degassing along tectonic  
574 structures of Mount Etna (Sicily) : the Pernicana fault. *Applied Geochemistry* 12(4),  
575 429-436.

576

577 Giammanco, S., Bellotti, F., Gropelli, G., Pinton, A., 2010. Statistical analysis reveals  
578 spatial and temporal anomalies of soil CO<sub>2</sub> efflux on Mount Etna volcano (Italy).  
579 *Journal of Volcanology and Geothermal Research* 194(1), 1-14.

580

581 Gresse, M., Vandemeulebrouck, J., Byrdina, S., Chiodini, G., Bruno, P. P., 2016.  
582 Changes in CO<sub>2</sub> diffuse degassing induced by the passing of seismic waves. *Journal Of*  
583 *Volcanology And Geothermal Research* 320, 12–18.

584

585 Gurrieri, S., Valenza, M., 1988. Gas transport in natural porous mediums : a method  
586 for measuring CO<sub>2</sub> flows from the ground in volcanic and geothermal areas. *Rend. Soc.*  
587 *Ital. Mineral. Petrol* 43, 1151-1158.

588



589 Hernández, P.A., Salazar, J.M., Shimoike, Y., Mori, T., Notsu, K., Pérez, N., 2001.  
590 Diffuse emission of CO<sub>2</sub> from Miyakejima volcano, Japan. *Chemical Geology* 177(1),  
591 175-185.  
592

593 Hinkle, M.E., Dilbert, C.A., 1984. Gases and trace elements in soils at the North Silver  
594 Bell deposit, Pima County, Arizona. *Journal of Geochemical Exploration* 20(3), 323-  
595 336.  
596

597 Irwin, W.P., Barnes, I., 1980. Tectonic relations of carbon dioxide discharges and  
598 earthquakes. *Journal of Geophysical Research : Solid Earth* 85(B6), 3115-3121.  
599

600 Kucera, C., Kirkham, D.R., 1971. Soil respiration studies in tallgrass prairie in  
601 Missouri. *Ecology* 52(5), 912-915.  
602

603 Liuzzo, M., Di Muro, A., Giudice, G., Michon, L., Ferrazzini, V., Gurrieri, S., 2015.  
604 New evidence of CO<sub>2</sub> soil degassing anomalies on Piton de la Fournaise volcano and  
605 the link with volcano tectonic structures. *Geochemistry, Geophysics, Geosystems*  
606 16(12), 4388-4404.  
607

608 Liuzzo, M., Gurrieri, S., Giudice, G., Giuffrida, G., 2013. Ten years of soil CO<sub>2</sub>  
609 continuous monitoring on Mt. Etna : Exploring the relationship between processes of  
610 soil degassing and volcanic activity. *Geochemistry, Geophysics, Geosystems* 14(8),  
611 2886-2899.  
612

613 Lovell, J.S., Hale, M., Webb, J.S., 1983. Soil air carbon dioxide and oxygen  
614 measurements as a guide to concealed mineralization in semi-arid and arid regions.  
615 *Journal of Geochemical Exploration* 19(1-3), 305-317.  
616

617 Neri, M., Acocella, V., Behncke, B., 2004. The role of the Pernicana Fault System in  
618 the spreading of Mt. Etna (Italy) during the 2002–2003 eruption. *Bulletin of*  
619 *Volcanology* 66(5), 417-430.  
620

621 Padrón, E., Melián, G., Marrero, R., Nolasco, D., Barrancos, J., Padilla, G., Hernández,  
622 P.A., Pérez, N.M., 2008. Changes in the diffuse CO<sub>2</sub> emission and relation to seismic  
623 activity in and around El Hierro, Canary Islands. In *Terrestrial Fluids, Earthquakes and*  
624 *Volcanoes : The Hiroshi Wakita Volume III*, 95-114, Birkhäuser Basel.  
625

626 Pearson, S.C.P., Connor, C.B., Sanford, W.E., 2008. Rapid response of a hydrologic  
627 system to volcanic activity: Masaya volcano, Nicaragua. *Geology* 36(12), 951-954.  
628

629 Revil, A., Finizola, A., Ricci, T., Delcher, E., Peltier, A., Barde-Cabusson, S., Avard,  
630 G., Bailly, T., Bennati, L., Byrdina, S., Colonge, J., Di Gangi, F., Douillet, G., Lupi,  
631 M., Letort, J., Tsang Hin Sun, E., 2011. Hydrogeology of Stromboli volcano, Aeolian  
632 Islands (Italy) from the interpretation of resistivity tomograms, self-potential, soil  
633 temperature and soil CO<sub>2</sub> concentration measurements. *Geophysical Journal*  
634 *International* 186, 1078–1094, <https://doi.org/10.1111/j.1365-246X.2011.05112.x>.  
635

636 Sandig, C., Sauer, U., Bräuer, K., Serfling, U., Schütze, C., 2014. Comparative study  
637 of geophysical and soil–gas investigations at the Hartoušov (Czech Republic) natural  
638 CO<sub>2</sub> degassing site. *Environmental earth sciences* 72(5), 1421-1434.

639

640 Sinclair, A.J., 1974. Selection of threshold values in geochemical data using probability  
641 graphs. *Journal of Geochemical Exploration* 3(2), 129-149.

642

643 Stephens, M.A., 1974. EDF Statistics for Goodness of Fit and Some Comparisons.  
644 *Journal of the American Statistical Association* 69, 730–737, doi:10.2307/2286009.

645

646 Stephens, M.A., 1976. Asymptotic Results for Goodness-of-Fit Statistics with  
647 Unknown Parameters. *Annals of Statistics* 4, 357-369.

648

649 Viveiros, F., Ferreira, T., Vieira, J.C., Silva, C., Gaspar, J.L., 2008. Environmental  
650 influences on soil CO<sub>2</sub> degassing at Furnas and Fogo volcanoes (São Miguel Island,  
651 Azores archipelago). *Journal of Volcanology and Geothermal Research* 177(4), 883-  
652 893.

653

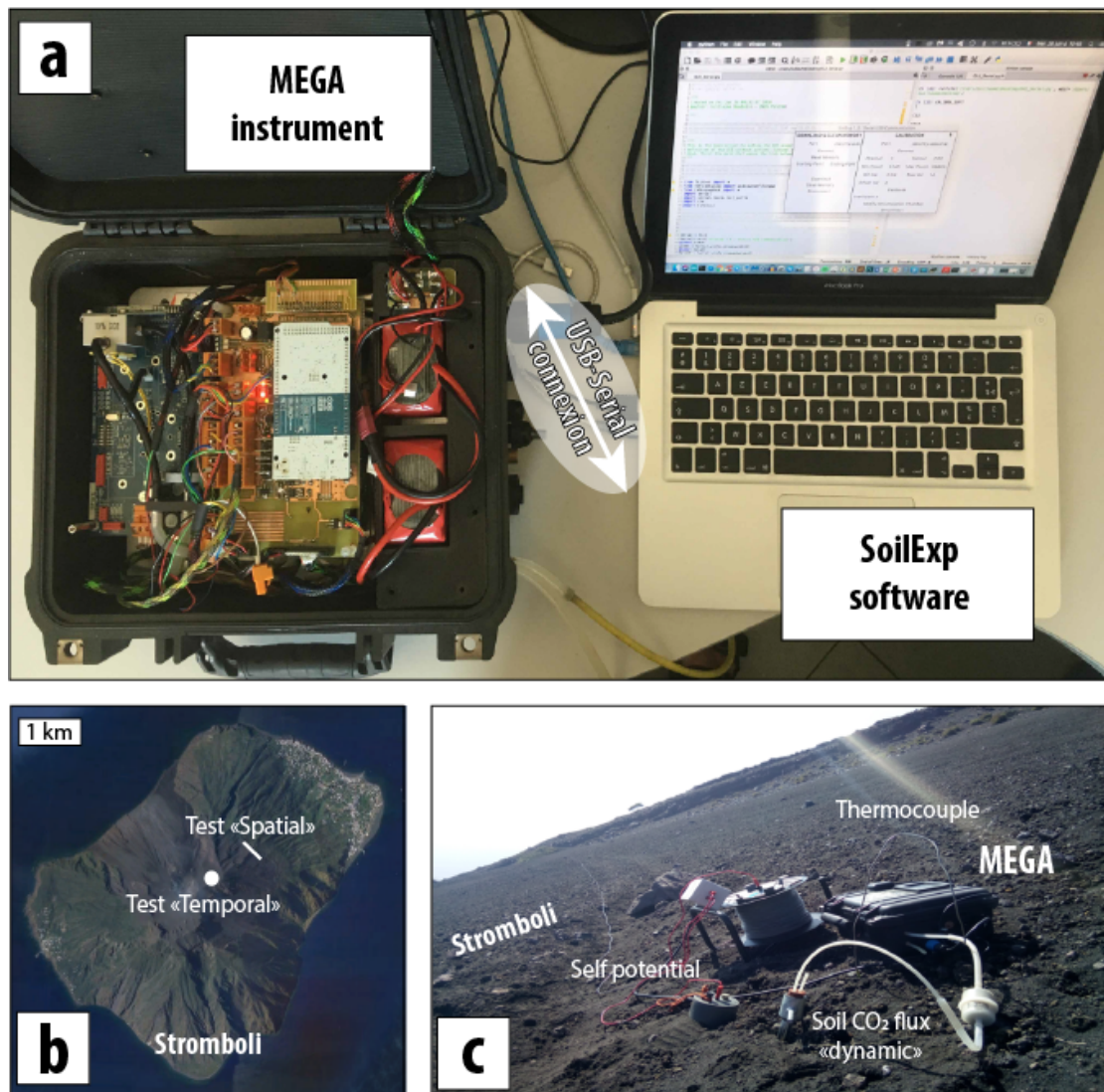
654 Shinohara, H., 2005. A new technique to estimate volcanic gas composition: Plume  
655 measurements with a portable multi-sensor system. *J. Volcanol. Geotherm. Res.* 143,  
656 319-333.

657

658 Zwillinger, D., Kokoska, S., 2000. *CRC Standard Probability and Statistics Tables and*  
659 *Formulae*. Chapman & Hall, New York.

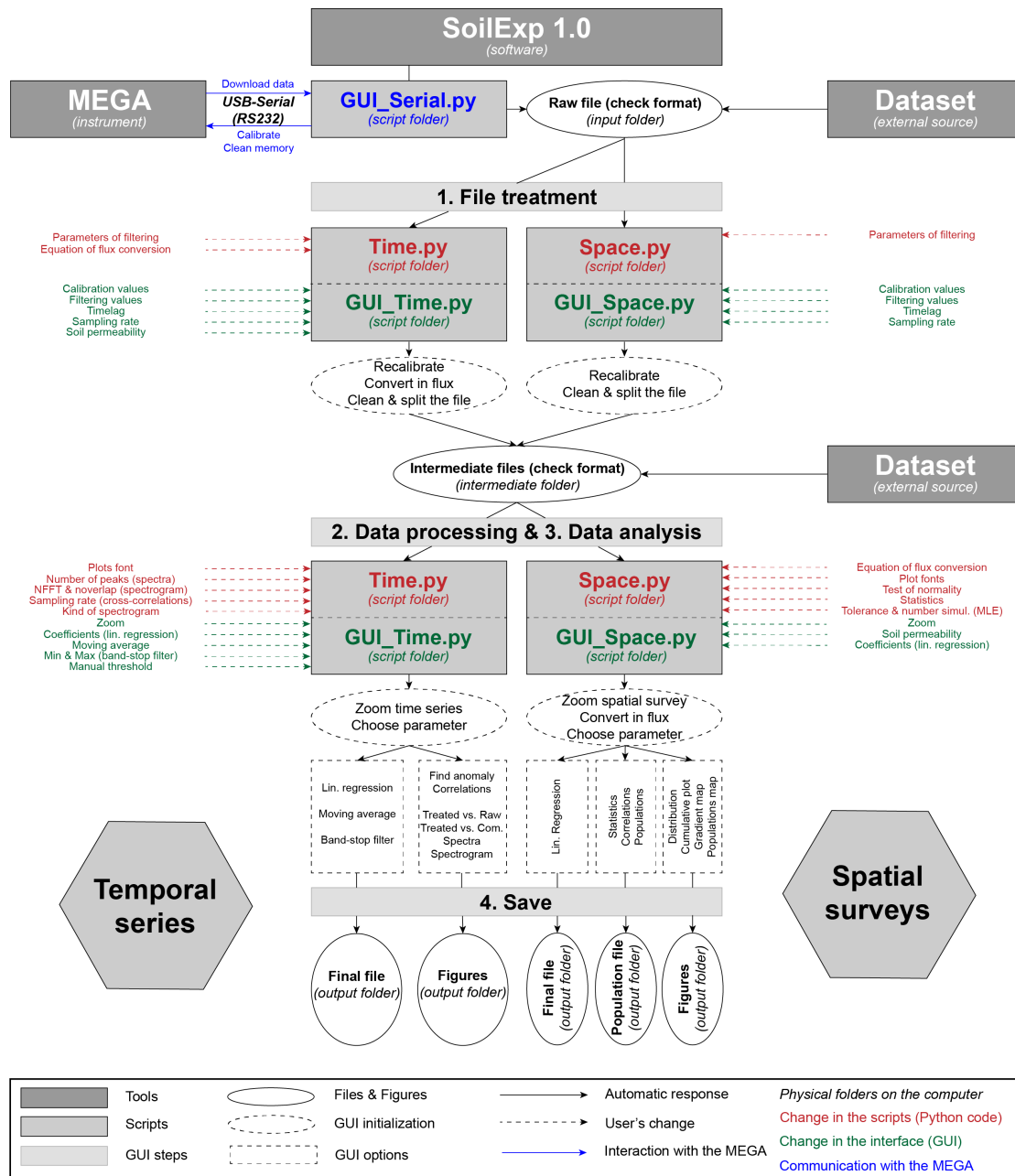
660

661 **Figures**



662

663 **Fig. 1. (a) The MEGA instrument and the SoilExp software, in evidence the USB-Serial**  
 664 **communication between the instrument and the software. (b) Tests of soil surveys at**  
 665 **Stromboli (Sicily, Italy) (c) based on soil CO<sub>2</sub> flux, self-potential and ground temperature.**  
 666 **These tests aim to illustrate the use of the SoilExp software in this study.**



667

668 Fig. 2. General scheme of use of the SoilExp software applied either to dataset acquired  
 669 with the MEGA instrument or though external sources.



SoilExp 1.0 : Serial USB Communication

DOWNLOAD & CLEAN MEMORY		CALIBRATION	
Port	/dev/tty.usbs	Port	/dev/tty.usbserial
Connect		Connect	
Read Memory		Channel	1
Starting Point	Ending Point	Sensor	CO2
		Min Count	5240
		Max Count	26650
		Min Val	0.04
		Max Val	10
		Offset Val	0
Download		Calibrate	
Clean Memory		Coefficient K	
Disconnect		Modify Accumulation Chamber	
		Disconnect	

**a**

SoilExp 1.0 : Space

1. FILE TREATMENT	2. DATA PROCESSING	3. DATA ANALYSIS																																																																								
<b>1a. Raw file</b> Load	<b>2a. Formatted file</b> Load & Init	<b>3a. Correlations</b> <table border="1"> <thead> <tr> <th></th> <th>Slope</th> <th>Offset</th> <th>R2</th> <th></th> <th>Slope</th> <th>Offset</th> <th>R2</th> </tr> </thead> <tbody> <tr> <td>CO2_100</td> <td>0.01</td> <td>-122.1</td> <td>0.18</td> <td>V_FLUX</td> <td>-0.16</td> <td>207.48</td> <td>0.02</td> </tr> <tr> <td>P_atm</td> <td>-0.00</td> <td>1027.9</td> <td>0.17</td> <td>P_in</td> <td>0.25</td> <td>-559.0</td> <td>0.05</td> </tr> <tr> <td>CO2_10</td> <td>0.03</td> <td>-128.7</td> <td>0.18</td> <td>Rh</td> <td>219.22</td> <td>-741.6</td> <td>0.09</td> </tr> <tr> <td>T_atm</td> <td>247.56</td> <td>-4247.</td> <td>0.21</td> <td>T_in</td> <td>-1.04</td> <td>-53.94</td> <td>0.01</td> </tr> <tr> <td>SP</td> <td>1.00</td> <td>0.00</td> <td>1.00</td> <td>V_BAT</td> <td>-22.90</td> <td>189.72</td> <td>0.00</td> </tr> <tr> <td>SO2</td> <td>147.28</td> <td>8165.2</td> <td>0.05</td> <td>PUMP_FLUX</td> <td>230.61</td> <td>-369.4</td> <td>0.09</td> </tr> <tr> <td>H2S</td> <td>394.12</td> <td>7128.4</td> <td>0.05</td> <td>Thermocouple</td> <td>-0.39</td> <td>-131.7</td> <td>0.04</td> </tr> <tr> <td>PYR</td> <td>-41557</td> <td>59273.</td> <td>0.06</td> <td>CO2_FLUX</td> <td>nan</td> <td>nan</td> <td>nan</td> </tr> </tbody> </table>		Slope	Offset	R2		Slope	Offset	R2	CO2_100	0.01	-122.1	0.18	V_FLUX	-0.16	207.48	0.02	P_atm	-0.00	1027.9	0.17	P_in	0.25	-559.0	0.05	CO2_10	0.03	-128.7	0.18	Rh	219.22	-741.6	0.09	T_atm	247.56	-4247.	0.21	T_in	-1.04	-53.94	0.01	SP	1.00	0.00	1.00	V_BAT	-22.90	189.72	0.00	SO2	147.28	8165.2	0.05	PUMP_FLUX	230.61	-369.4	0.09	H2S	394.12	7128.4	0.05	Thermocouple	-0.39	-131.7	0.04	PYR	-41557	59273.	0.06	CO2_FLUX	nan	nan	nan
	Slope		Offset	R2		Slope	Offset	R2																																																																		
CO2_100	0.01	-122.1	0.18	V_FLUX	-0.16	207.48	0.02																																																																			
P_atm	-0.00	1027.9	0.17	P_in	0.25	-559.0	0.05																																																																			
CO2_10	0.03	-128.7	0.18	Rh	219.22	-741.6	0.09																																																																			
T_atm	247.56	-4247.	0.21	T_in	-1.04	-53.94	0.01																																																																			
SP	1.00	0.00	1.00	V_BAT	-22.90	189.72	0.00																																																																			
SO2	147.28	8165.2	0.05	PUMP_FLUX	230.61	-369.4	0.09																																																																			
H2S	394.12	7128.4	0.05	Thermocouple	-0.39	-131.7	0.04																																																																			
PYR	-41557	59273.	0.06	CO2_FLUX	nan	nan	nan																																																																			
<b>1b. Linear Calibration</b> Parameter Slope Offset Recalculate	<b>2b. Subset</b> Start End 0 Subset -1 Zoom Point Permeability All 35 Convert CO2																																																																									
<b>1c. Treatment</b> TimeLag 60 Sampling Rate 1 Battery 12 14 Pump -15 15 HDOP 1 200 R-squared 0.9 1.0 Clean & Create	<b>2c. Parameter</b> Choose the parameter to study SP Statistics <b>2d. Regression</b> Slope Param Offset 0 0 Correct	<b>3b. Statistics</b> Mean -100.6 Median -119.1 Min -247.9 Max 141.4 Stdev 90.616 Kurtosis 0.5697 Skewness -0.107 N° Populations 1.0 Distribution Cumulative <b>3c. Maps</b> Google API API_KEY Zoom (5-16) 15 Scatter Map (gradient) Scatter Map (populations)																																																																								
<b>4. SAVE</b> Save Spatial Survey Save Populations																																																																										

**b**

SoilExp 1.0 : Time

1. FILE TREATMENT	2. DATA PROCESSING	3. DATA ANALYSIS																																																																																																
<b>1a. Raw file</b> Load	<b>2a. Formatted file</b> Load	<table border="1"> <thead> <tr> <th></th> <th>Slope</th> <th>Offset</th> <th>R2</th> <th>Delay</th> <th>R2</th> </tr> </thead> <tbody> <tr> <td>CO2_100</td> <td>0.83</td> <td>196.74</td> <td>0.74</td> <td>0.00</td> <td>0.74</td> </tr> <tr> <td>P_atm</td> <td>-15.52</td> <td>15849.</td> <td>0.05</td> <td>323.00</td> <td>0.19</td> </tr> <tr> <td>CO2_10</td> <td>1.00</td> <td>0.00</td> <td>1.00</td> <td>0.00</td> <td>1.00</td> </tr> <tr> <td>T_atm</td> <td>-72.34</td> <td>2706.6</td> <td>0.02</td> <td>255.00</td> <td>0.15</td> </tr> <tr> <td>SP</td> <td>-1.57</td> <td>1435.1</td> <td>0.00</td> <td>60.00</td> <td>0.14</td> </tr> <tr> <td>SO2</td> <td>-26.72</td> <td>-0.71</td> <td>0.00</td> <td>422.00</td> <td>0.22</td> </tr> <tr> <td>H2S</td> <td>-70.18</td> <td>211.44</td> <td>0.00</td> <td>422.00</td> <td>0.22</td> </tr> <tr> <td>PYR</td> <td>-1435E</td> <td>22025.</td> <td>0.02</td> <td>444.00</td> <td>0.13</td> </tr> <tr> <td>V_FLUX</td> <td>nan</td> <td>nan</td> <td>nan</td> <td>0.00</td> <td>nan</td> </tr> <tr> <td>P_in</td> <td>-0.02</td> <td>1536.9</td> <td>0.00</td> <td>665.00</td> <td>0.15</td> </tr> <tr> <td>Rh</td> <td>6.75</td> <td>1484.3</td> <td>0.00</td> <td>666.00</td> <td>0.03</td> </tr> <tr> <td>T_in</td> <td>0.49</td> <td>1478.8</td> <td>0.00</td> <td>409.00</td> <td>0.15</td> </tr> <tr> <td>V_BAT</td> <td>-62.42</td> <td>2302.6</td> <td>0.00</td> <td>668.00</td> <td>0.04</td> </tr> <tr> <td>PUMP_FLUX</td> <td>6.69</td> <td>1496.1</td> <td>0.00</td> <td>666.00</td> <td>0.03</td> </tr> <tr> <td>Thermocouple</td> <td>nan</td> <td>nan</td> <td>0.00</td> <td>0.00</td> <td>0.00</td> </tr> </tbody> </table>		Slope	Offset	R2	Delay	R2	CO2_100	0.83	196.74	0.74	0.00	0.74	P_atm	-15.52	15849.	0.05	323.00	0.19	CO2_10	1.00	0.00	1.00	0.00	1.00	T_atm	-72.34	2706.6	0.02	255.00	0.15	SP	-1.57	1435.1	0.00	60.00	0.14	SO2	-26.72	-0.71	0.00	422.00	0.22	H2S	-70.18	211.44	0.00	422.00	0.22	PYR	-1435E	22025.	0.02	444.00	0.13	V_FLUX	nan	nan	nan	0.00	nan	P_in	-0.02	1536.9	0.00	665.00	0.15	Rh	6.75	1484.3	0.00	666.00	0.03	T_in	0.49	1478.8	0.00	409.00	0.15	V_BAT	-62.42	2302.6	0.00	668.00	0.04	PUMP_FLUX	6.69	1496.1	0.00	666.00	0.03	Thermocouple	nan	nan	0.00	0.00	0.00
	Slope		Offset	R2	Delay	R2																																																																																												
CO2_100	0.83	196.74	0.74	0.00	0.74																																																																																													
P_atm	-15.52	15849.	0.05	323.00	0.19																																																																																													
CO2_10	1.00	0.00	1.00	0.00	1.00																																																																																													
T_atm	-72.34	2706.6	0.02	255.00	0.15																																																																																													
SP	-1.57	1435.1	0.00	60.00	0.14																																																																																													
SO2	-26.72	-0.71	0.00	422.00	0.22																																																																																													
H2S	-70.18	211.44	0.00	422.00	0.22																																																																																													
PYR	-1435E	22025.	0.02	444.00	0.13																																																																																													
V_FLUX	nan	nan	nan	0.00	nan																																																																																													
P_in	-0.02	1536.9	0.00	665.00	0.15																																																																																													
Rh	6.75	1484.3	0.00	666.00	0.03																																																																																													
T_in	0.49	1478.8	0.00	409.00	0.15																																																																																													
V_BAT	-62.42	2302.6	0.00	668.00	0.04																																																																																													
PUMP_FLUX	6.69	1496.1	0.00	666.00	0.03																																																																																													
Thermocouple	nan	nan	0.00	0.00	0.00																																																																																													
<b>1b. Linear Calibration</b> Parameter Slope Offset Recalculate	<b>2b. Initialize</b> Time Serie CO2_10 Zoom 0 start / stop -1 Initialize																																																																																																	
<b>1c. Treatment</b> TimeLag 60 Sampling Rate 1 Battery 12 14 Pump -15 15 HDOP 1 200 Soil Permeability 35 Clean & Split	<b>2c. Process</b> Regression 0 MovAverage 0 time window CutBand 1 min / max 1 Comparison Treated vs. Raw Treated vs. Comparative Spectra Spectrogram Correlations Find Anomaly Threshold 0																																																																																																	
<b>4. SAVE</b> Save Time Serie																																																																																																		

**c**

671 **Fig. 3. Graphical User Interface (GUI) of the “Serial” (a), “Space” (b) and “Time” (c)**  
672 **option of the SoilExp software. The GUI is divided in 4 panels. Panel (1) is dedicated to**  
673 **format the raw file in intermediate formatted files after applying potential distinct**  
674 **calibrations and conversions, and cleaning the dataset. Panel (2) aimed to process the data**  
675 **obtained from the intermediate formatted files either from the previous step or formatted**  
676 **independently by the user (conversion, moving average, linear regression, cut band filter).**  
677 **Panel (3) shows the result of the datasets processing and analysis (correlations, cross-**  
678 **correlations, statistics, analysis of populations, distribution, maps). Panel (4) allows to**  
679 **save the dataset transformed with the above operations in final .csv file.**

680

681

682

683

684

685

686

687

688

689

690

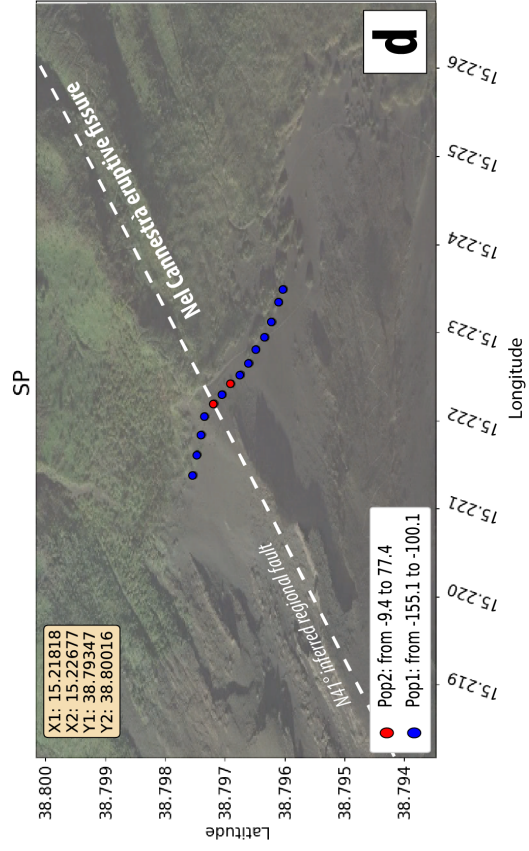
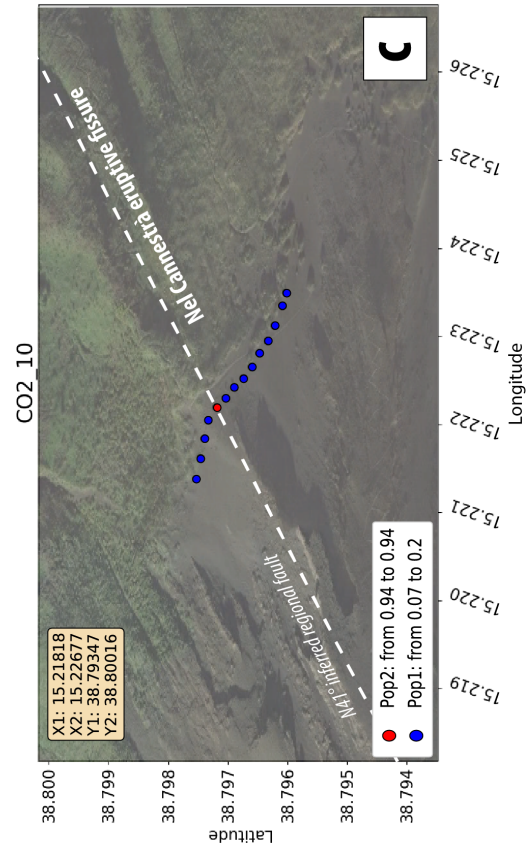
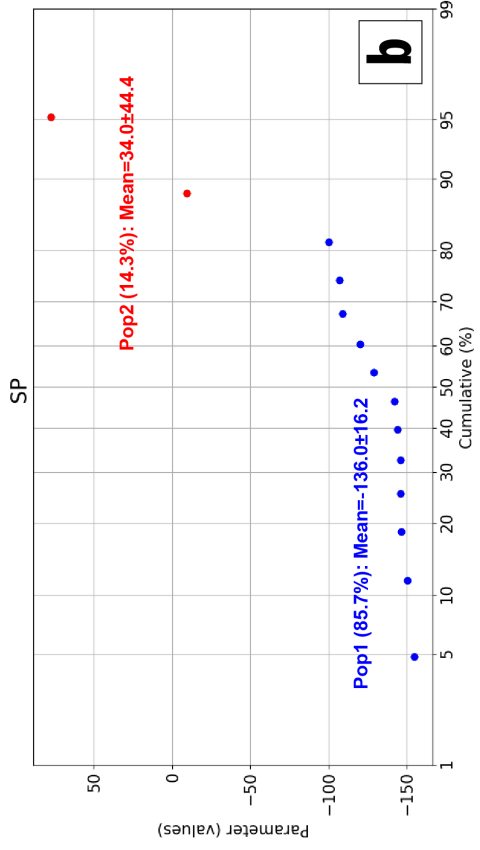
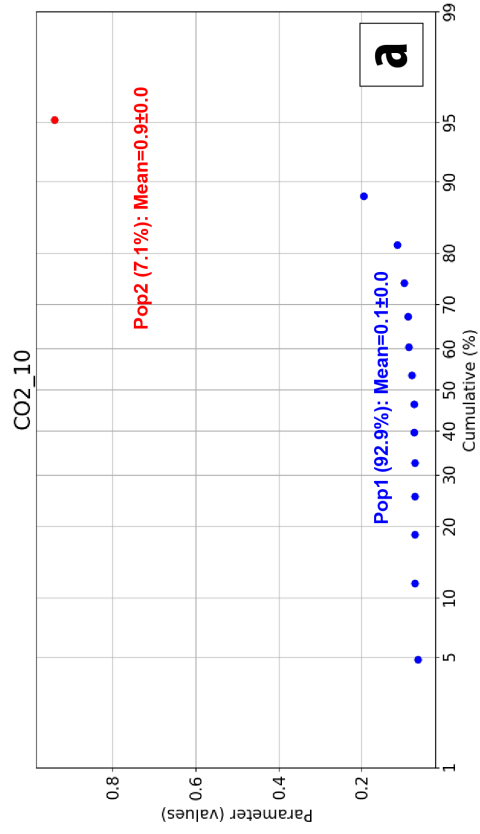
691

692

693

694

695





697 Fig. 4. Example of data analysis obtained by using the “Space” option at Stromboli (soil  
698 CO<sub>2</sub> flux and self-potential measurements along a transect with a 20 m-spacing; cf. Test  
699 “Spatial” on Fig. 1c, d). Probability plot of (a) soil CO<sub>2</sub> flux measurements obtained using  
700 a 0-10 %molar IR spectrometer (e.g. CO<sub>2</sub>\_10 by “dynamic” concentration; Gurrieri &  
701 Valenza, 1988; Camarda *et al.*, 2006) and (b) self-potential measurements carried out with  
702 a pair of non-polarizable Cu/CuSO<sub>4</sub> electrodes (e.g. SP; Finizola *et al.*, 2010). The  
703 identification of distinct populations is based on the maximum-likelihood numerical  
704 method (see text). Map highlighting the corresponding (c) soil CO<sub>2</sub> flux and (d) self-  
705 potential transect performed at Stromboli (cf. Fig. 1). The satellite map is obtained from  
706 Google Map. In case of absence of API key  
707 (<https://developers.google.com/maps/documentation/javascript/get-api-key>), the  
708 background will stay white. However, the (decimal) coordinates of the corners are  
709 reported in the upper left box in order to let the user free to download a map from distinct  
710 sources. In this example, the “Space” option allows to identify a soil CO<sub>2</sub> anomaly coupled  
711 with a positive SP anomaly that highlight an upward migration of hydrothermal fluids  
712 along the Nel Cannestrà eruptive fissure. This result is in accordance with previous study  
713 (Finizola *et al.*, 2002, 2010; Carapezza *et al.*, 2009).

714

715

716

717

718

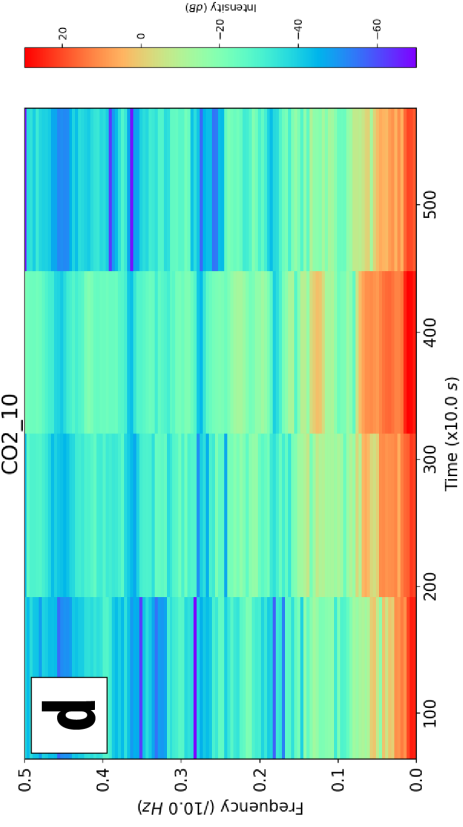
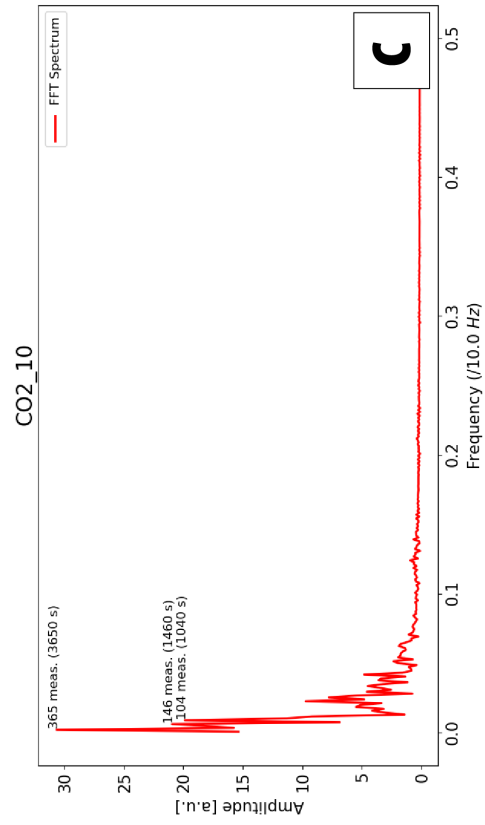
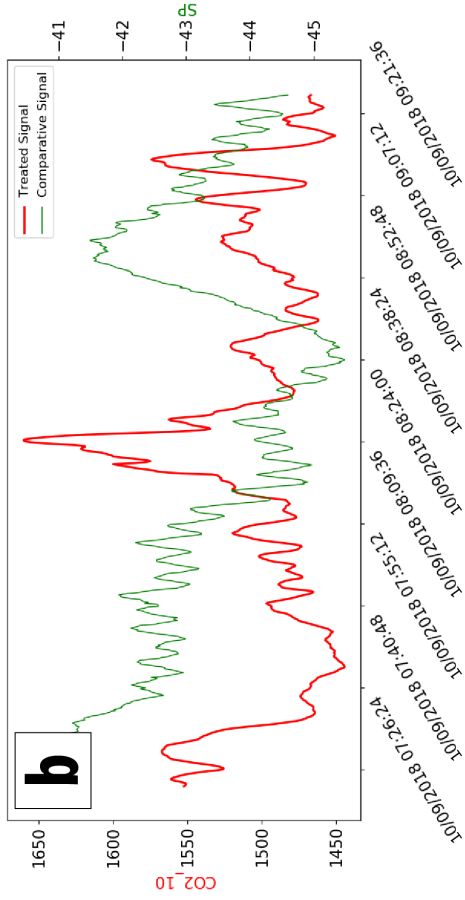
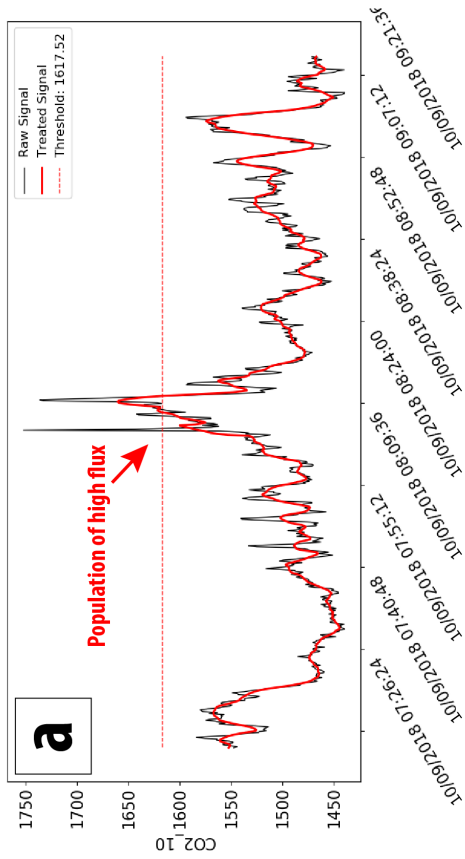
719

720

721

722

723



725 **Fig. 5. Example of data analysis obtained by using the “Time” option at Stromboli (soil**  
726 **CO<sub>2</sub> flux measured during about 2 hours at 0.1 Hz at the same site; cf. Test “Temporal”**  
727 **on Fig. 1b). (a) Comparison between raw and treated data (after applying the moving**  
728 **average). The threshold analysis allows us to detect the highest population of values (often**  
729 **considered as “anomalous” values) during the acquisition. (b) Comparison between**  
730 **treated soil CO<sub>2</sub> flux (e.g. CO2\_10) and self-potential (e.g. SP) time series. Here the**  
731 **detected soil CO<sub>2</sub> flux anomaly is synchronous with low self-potential records. (c) Fast**  
732 **Fourier Transform (FFT) spectrum of the treated soil CO<sub>2</sub> signal. The 3 greatest**  
733 **frequency peaks are labelled with the corresponding period that may be cut using the cut**  
734 **band filter. (d) Spectrogram of the treated soil CO<sub>2</sub> signal (linear detrend; magnitude**  
735 **mode; NFFT=256; noverlap=128). In this example, with about 1200 measurements, there**  
736 **are not enough data available to obtain a smoothed spectrogram considering a NFFT of**  
737 **256.**