

¹ [COR: Number]
² Cover Letter

³ **formikoj: A flexible library for data management and processing in geophysics - Application
4 for seismic refraction data**

⁵ Matthias Steiner, Adrián Flores Orozco

⁶ Dear Editors-in-Chief,

⁷
⁸ we are submitting our manuscript "formikoj: A flexible library for data management and processing in geophysics -
⁹ Application for seismic refraction data", which we consider is a suitable contribution for Computers & Geosciences.
¹⁰ We confirm that the submission follows all the requirements and includes all the items of the submission checklist.

¹¹
¹² The manuscript presents the open-source python library formikoj for managing and processing geophysical data col-
¹³ lected in environmental and engineering investigations. formikoj was specifically implemented for multi-platform
¹⁴ usage to allow the efficient collaboration and exchange of data between different partners in research projects and
¹⁵ academia. In this regard, we believe that this library aids in providing reproducible data and results as well as estab-
¹⁶ lishing and maintaining good research practices. Accordingly, we consider this manuscript relevant to the audience of
¹⁷ Computers & Geosciences, and in general for geoscientists and practitioners working with geophysical methods.

¹⁸
¹⁹ We provide the source codes in a public repository with details listed in the section "Code availability".

²⁰
²¹ We look forward to your decision.

²²
²³ Yours sincerely,

²⁴
²⁵ Matthias Steiner and Adrián Flores Orozco
²⁶ Research Unit Geophysics, Department of Geodesy and Geoinformation, TU Wien; matthias.steiner@geo.tuwien.ac.at
²⁷

²⁸ **Highlights**

²⁹ **formikoj: A flexible library for data management and processing in geophysics - Application**
³⁰ **for seismic refraction data**

³¹ Matthias Steiner, Adrián Flores Orozco

- ³² • flexible open-source and cross-platform library for managing and processing of geophysical data
³³ • possibility to be deployed for different geophysical methods and/or instruments
³⁴ • application for the modeling and processing of seismic refraction datasets
³⁵ • applicable for seismic refraction data collected in 2D and 3D survey geometries
³⁶ • easily scalable for custom requirements

37 **formikoj: A flexible library for data management and processing in**
38 **geophysics - Application for seismic refraction data**

39 Matthias Steiner^{a,*}, Adrián Flores Orozco^a

40 ^aResearch Unit Geophysics, Department of Geodesy and Geoinformation, TU Wien

41

42 **ARTICLE INFO**

43

44 **Keywords:**
45 geophysical data processing
46 seismic refraction
47 first break picking
48 seismic waveform modeling
49 cross-platform application
50 geophysical python library
51 flexible open-source libraries
52 wave based methods

53

54

55

56

57

58

59 **CRediT authorship contribution statement**

60 **Matthias Steiner:** Conceptualization and implementation of the library, creating the figures, preparation of the
61 manuscript. **Adrián Flores Orozco:** Conceptualization of the library, preparation of the manuscript.

62 **1. Introduction**

63 The acquisition of spatially quasi-continuous data in a non-invasive manner renders geophysical methods suitable
64 for engineering and environmental investigations (e.g., Parsekian et al., 2015; Nguyen et al., 2018; Romero-Ruiz et al.,
65 2018). However, the processing of geophysical data often relies on commercial software solutions and the associ-
66 ated licensing costs might render their use prohibitively expensive, which might be the case for academic projects
67 or institutions. The most popular packages are Res2DInv¹ for electrical methods, Halliburton Landmark SeisSpace
68 ProMAX² or ParkSeis³ for seismic methods, or ReflexW⁴ for ground-penetrating radar and seismic methods. A com-
69 mon limitation of the aforementioned software solutions refers to their specific platform requirements mainly related
70 to the type and version of the operating system; moreover, the possibility to adapt the code are limited if possible at
71 all. Considering the substantial changes regarding the market shares of operating systems within the last two decades,
72 platform-specific software packages are becoming particularly obstructive for academic research and teaching. The

*Corresponding author

ORCID(s): 0000-0002-3595-3616 (M. Steiner); 0000-0003-0905-3718 (A. Flores Orozco)

¹<https://www.geometrics.com/software/res2dinv/>, last accessed on July 21, 2022

²<https://www.landmark.solutions/SeisSpace-ProMAX>, last accessed on July 21, 2022

³<https://www.parkseismic.com/parkseis/>, last accessed on July 21, 2022

⁴<https://www.sandmeier-geo.de/reflexw.html>, last accessed on July 21, 2022

73 increasing popularity of the Python programming language led to the development of various cross-platform open-
 74 source software packages for processing, modeling and inverting geophysical data. Available packages can focus on
 75 specific geophysical methods, for instance, ResIPy (Blanchy et al., 2020) for electrical data, GPRPy (Plattner, 2020)
 76 for ground-penetrating radar data, or ObsPy (Beyreuther et al., 2010) and Pyrocko (Heumann et al., 2017) for seis-
 77 mological data. Other packages provide frameworks for the inversion and permit the inclusion of forward models for
 78 different geophysical methods, e.g., SimPEG (Cockett et al., 2015), Fatiando a Terra (Uieda et al., 2013) or pyGIMLi
 79 (Rücker et al., 2017).

80 The seismic refraction tomography (SRT) is a standard technique in environmental and engineering applications.
 81 Often applied together with other geophysical methods, the SRT is routinely used, e.g., in permafrost studies (e.g.,
 82 Draebing, 2016; Steiner et al., 2021), for the investigation of landfills (e.g., Nguyen et al., 2018; Steiner et al., 2022),
 83 or for hydrogeological characterizations (e.g., Bücker et al., 2021). The market for seismic processing software has
 84 long been dominated by software packages designed for the processing of large datasets, e.g., associated to oil or gas ex-
 85 ploration. Accordingly, these seismic processing solutions may not be suited for small-scale projects in environmental
 86 and engineering studies, or for teaching activities. ReflexW overcomes such limitations by providing processing tools
 87 specifically designed for near-surface investigations at substantially lower costs. In terms of licensing costs, Stockwell
 88 (1999) went a step further by making the Seismic Unix framework available entirely free of charge; whereas Guedes
 89 et al. (2022) recently presented RefraPy, a python processing tool for seismic refraction data. Implemented in python,
 90 RefraPy is potentially suitable for cross-platform usage, yet Guedes et al. (2022) developed and tested solely for Win-
 91 dows operating systems. Moreover, RefraPy does not offer the possibility to generate synthetic seismic waveform data,
 92 as required for survey design, as well as teaching and interpretation purposes.

93 The formikoj library presented here is an open-source framework for creating synthetic datasets, as well as for man-
 94 aging and processing numerical and field data independently from the operating system and without licensing costs;
 95 thus, overcoming limitations associated to existing solutions. The design of the library follows the multi-method
 96 concept of pyGIMLi and SimPEG, which allows for the implementation of custom designed tools for different geo-
 97 physical methods. The usage of transparent file formats, e.g., the unified data format (udf⁵), and data management
 98 concepts (SQLite database) within the formikoj framework facilitates a simple data exchange between partners in re-
 99 search projects and academia, which is required to guarantee the repeatability of results and good research practices.
 100 Considering the diverse applications of the SR method we demonstrate the applicability of the proposed library based
 101 on tools implemented within the formikoj framework for the modeling and processing seismic waveform data. In par-
 102 ticular, we present here a series of illustrative use cases based on the formikoj library referring to (i) the modeling of
 103 synthetic seismic refraction (SR) waveform data, (ii) the processing of a 2D SR field dataset collected with a roll-along

⁵http://resistivity.net/bert/data_format.html, last accessed on July 21, 2022

104 survey geometry, and (iii) the processing of a 3D SR field data set.

105 **2. Design and structure of the formikoj library**

106 As illustrated in Figure 1, the formikoj library comprises a modeling and a processing module, which both rely on a
 107 common utilities module. The `DataModeler` and the `MethodManager` class provide the basis to add modeling or pro-
 108 cessing functionalities for specific geophysical methods. In particular, we present here the `SeismicWaveformModeler`
 109 and `SeismicRefractionManager` classes implemented within the formikoj framework, which aim at creating and
 110 processing seismic waveform data, respectively. Similar to RefraPy, these classes are built upon the functionalities
 111 of existing packages such as ObsPy for the processing of seismological data (Beyreuther et al., 2010) and pyGIMLi
 112 for the modeling and inversion of different geophysical data (Rücker et al., 2017). Other important third party de-
 113 pendencies refer to NumPy (Harris et al., 2020) and Pandas (McKinney, 2010) for general data handling, as well as
 114 matplotlib (Hunter, 2007) and PyVista (Sullivan and Kaszynski, 2019) for data visualization. In the current version,
 115 we implemented and tested formikoj primarily on Linux machines, yet the library has been successfully used on all
 116 major operating systems, i.e., Linux, MacOS and Windows.

117 **2.1. Generation of seismic waveform data for synthetic subsurface models: The**

118 `SeismicWaveformModeler`

119 The SR method exploits the ground motion recorded by sensors installed in the surface (e.g., geophones) to char-
 120 acterize the propagation of seismic waves generated at well known locations (i.e., shot stations). The visualization
 121 of the ground motion as function of time yields a so-called seismogram for each geophone position. Accordingly,
 122 the `SeismicWaveformModeler` class provides a flexible way to generate synthetic seismic waveform data either in a
 123 python script, interactively in a jupyter notebook or an ipython shell. To create an instance of the class the path to the
 124 working directory is provided as parameter to the constructor:

```
125 from useis import SeismicWaveformModeler
126
127 swm = SeismicWaveformModeler('..')
128
129 INFO    : Created instance of SeismicWaveformModeler
```

127 The working directory needs to contain a subdirectory *in*, whereas the output directory *out* will be created automati-
 128 cally:

```
129
130 working_directory
131   |
132   +-- in
133   |
134   +-- out
```

129 The required input files need to be provided via the subdirectory *in* of the working directory. The key input file is the
 130 measurement scheme, which contains information regarding the distribution of the shot and geophone stations. If pro-

Table 1

Description of the information to be provided in the columns of a measurement scheme in csv format.

Column	Content	Data type	Description
1	x coordinate	float	Station x coordinate, e.g., given in (m)
2	y coordinate	float	Station y coordinate, e.g., given in (m)
3	z coordinate	float	Station z coordinate, e.g., given in (m)
4	Geophone	bool	1 in case of a receiver station, 0 otherwise
5	Shot	bool	1 in case of a shot station, 0 otherwise

131 vided in the unified data format the measurement scheme is imported directly with pyGIMLi into a DataContainer.
 132 In case the measurement scheme is provided as a csv file the SeismicWaveformModeler reads the information and
 133 writes it to a pyGIMLi DataContainer. In the csv format the measurement scheme contains a single line for each
 134 station in the survey layout, where a station either hosts a geophone or a shot, or both (see Table 1). The values provided
 135 in each line need to be separated by a unique delimiter, and the file must not contain a header.

136 For the modeling of the seismic waveforms, the parameters for the base wavelet, the synthetic subsurface model
 137 and the resulting waveform datasets are provided (see Table 2) in a configuration file following the yaml format:

```
wavelet:
  length: 1.024
  frequency: 100
  sampling_rate: 2000
  pretrigger: 0.02

dataset:
  number: 2
  names: [dataset1, dataset2]
  noise: 0
  noise_level: 1.e-4
  missing_shots: 0
  broken_geophones: 0
  wrong_polarity: 0

model:
  velmap: [[1, 750.], [2, 3000.]]
  layers: [[1, 3.], [2, 10.]]
  quality: 32
  area: 10
  smooth: [1, 10]
  sec_nodes: 3

traveltimes:
  noise_relative: 0.
  noise_absolute: 0.
```

138 In the exemplary configuration file shown above, the first block contains information regarding the wavelet, which con-
 139 trols the modeling of the seismic waveforms as described in Table 2. In the second block, we can set specific names
 140 for the datasets to be created and the number of datasets is automatically determined. Alternatively, it is possible to
 141 set the number of datasets to be created and the dataset names are automatically generated with the prefix *dataset_*.
 142 Various parameters control the random error (*noise*) and systematic errors in the modeled seismic waveform data (see
 143 Table 2). The number and position of the shot and geophone stations affected by the systematic errors are randomly
 144 chosen with a maximum of 5% of the total number of stations in order to avoid a high number of invalid trace data. The
 145 third block contains information regarding the synthetic subsurface model. For each layer the corresponding velocity
 146

147 (velmap) and layer thickness (layers) need to be provided and all layers are considered to be parallel to the surface
 148 topography (geometrical information regarding the stations in the measurement scheme). The remaining parameters,
 149 namely quality, area, smooth and sec_nodes, define the properties of the mesh to be used for the forward mod-
 150 eling (we refer to the respective pyGIMLi resources⁶ for further information). Alternatively, the user can provide a
 151 more complex mesh in the binary mesh format (i.e., a bms file):

```
swm.load('mesh')
```

152 INFO : Mesh loaded

153 The parameters in the final block control the forward modeling of the corresponding seismic traveltimes (see Table 2).
 154 A configuration file located in the input directory can be imported through the load method:

```
swm.load('config')
```

155 INFO : Configuration loaded

156 Once the **SeismicWaveformModeler** instance is parameterized the synthetic seismic waveform data can be created
 157 as follows:

```
swm.create('waveforms')
```

158 INFO : Measurement scheme loaded
 INFO : Velocity model created
 INFO : Wavelet created
 [+++++ 100% +++++] 2048 of 2048 complete
 ...
 [+++++ 100% +++++] 2048 of 2048 complete
 INFO : Dataset 'clean' created

159 The **SeismicWaveformModeler** loads the measurement scheme and creates the velocity model used for the waveform
 160 modeling. Based on the wavelet properties a Ricker wavelet is generated through the pyGIMLi function ricker. Sub-
 161 sequently, mesh, velocity model and Ricker wavelet are used to solve the pressure wave equation for each shot station
 162 defined in the measurement scheme with the pyGIMLi function solvePressureWave. The resultant waveforms at
 163 the corresponding geophone stations are extracted and stored in an ObsPy Stream data structure.

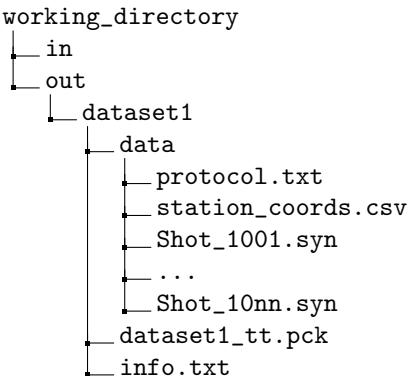
⁶https://www.pygimli.org/pygimliapi/_generated/pygimli.mesh.html#pygimli.mesh.createMesh, last accessed July 21, 2022

Table 2

Description of the parameters, which can be defined in a configuration file used for the modeling of the synthetic seismic data.

Parameter	Unit/ Data type	Description
wavelet		
length	s	Length of the base wavelet, which also defines the length of the synthetic seismic waveform data
frequency	Hz	Frequency of the base wavelet
sampling_rate	Hz	Defines temporal resolution of the seismic waveforms
pretrigger	s	Add buffer to the seismic waveforms before the onset of the actual data
dataset		
number	int	Number of datasets to be created
names	list (string)	Names of the datasets
noise	bool	1 in case noise should be added to the synthetic waveforms, 0 otherwise
noise_level	-	Level of the seismic background noise
missing_shots	bool	1 in case the datasets should be affected by missing shot files, 0 otherwise
broken_geophones	bool	1 in case the datasets should comprise broken geophones (i.e., no data in the corresponding seismograms), 0 otherwise
wrong_polarity	bool	1 in case the datasets should contain traces with inverse polarity, 0 otherwise
model		
velmap	list (float/int)	For each layer the first value defines the marker and the second one the seismic velocity within the layer
layers	list (float/int)	For each layer the first value defines the marker and the second one the thickness
traveltimes		
noise_relative	%/100	Relative noise to be added to the forward modeled seismic traveltimes
noise_absolute	s	Absolute noise to be added to the forward modeled seismic traveltimes

164 A directory for each dataset defined in the configuration file will be created in the output directory (*out*):



165 In the subdirectory *data*, the synthetic seismic waveforms are stored in a separate shot file for each shot position in
 166 the miniseed format (Ahern et al., 2012; Ringler and Evans, 2015) with the file extension *syn* to identify the forward
 167 modeled shot files, e.g., *Shot_1001.syn*. The measurement protocol (*protocol.txt*) and the station coordinates provided
 168 as a csv file (*station_coords.csv*) are also stored in this directory. The header of the measurement protocol contains
 169 the survey parameters required for the processing of the seismic waveforms, namely sampling rate, recording length,
 170 number of geophones and geophone spacing. Moreover, the protocol associates each shot file of the dataset to a spe-
 171 cific location within the survey geometry, i.e., with respect to the geophone positions:

```

#####
Line: SYN_dataset_1
Sampling rate: 2000 Hz
Recording length: 1.024 s
Number of geophones: 48
Geophone spacing: 6 m
#####
File number | Station
1001      | G001
:          | :
1048      | G048

```

172

173 The auxiliary file *info.txt* provided in the dataset directory summarizes the parameters from the configuration file and
 174 information regarding the simulated systematic errors in the synthetic seismic waveform data:

```

Number of geophones: 48
Number of shots: 48
Recording length (s): 1.024
Sampling frequency (Hz): 2000
Wavelet type: Ricker
Frequency of the wavelet (Hz): 100

Missing shot(s): 16
Broken geophone(s): 6, 14
Wrong polarity geophone(s): 35

```

175

176 Synthetic traveltimes (`swm.create('traveltimes')`) are stored in the dataset directory as udf files, e.g., *dataset1_tt.pck*.

177 The file extension *pck* is an abbreviation of the word 'pick' and refers to the first break traveltimes saved in the file.

178 **2.2. Processing of seismic refraction datasets: the SeismicRefractionManager**

179 The SR method is based on the measurement of the traveltimes of seismic waves determined from the the first onset
 180 of the seismic energy recorded by geophones. In the SRT, the inversion of traveltimes gathered from tens to hundreds of
 181 seismograms permits the computation of variations in the seismic velocities in an imaging framework. Measuring the
 182 traveltimes in seismograms (i.e., first break picking) is commonly done manually (or semi-automatically) in an iterative
 183 process. The signal-to-noise ratio in the recorded seismograms substantially influences the quality of the traveltimes
 184 picked in the seismograms, and thus a proper enhancement of the perceptibility of the first onsets is crucial.

185 Accordingly, the `SeismicRefractionManager` class provides functionalities that permit the processing of seis-
 186 mic waveform data for a refraction seismic analysis. These functionalities involve the reading of seismic waveform data,
 187 combining the data with information about the survey geometry, processing of the waveforms as well as the picking of
 188 first break traveltimes. Since the first break picking is a highly interactive process, the `SeismicRefractionManager`
 189 is designed primarily for usage from within an ipython shell.

190 **2.2.1. Compiling the survey information and creating a project**

191 An instance of the `SeismicRefractionManager` can be created by providing the path to the working directory
 192 as parameter to the constructor. Based on the content of the working directory, the `SeismicRefractionManager`
 193 automatically decides whether (i) to start in the data preview mode, (ii) create a new project, or (iii) load an existing
 194 project from disk.

195 The data preview mode is primarily initiated if the provided working directory contains seismic shot files:

```
from useis import SeismicRefractionManager
srm = SeismicRefractionManager('./01_data/raw')

INFO    : Starting in data preview mode
Progress <=====100.0% completed>
INFO    : Read 48 files
```

196 To create or load a project, the working directory needs to contain specific subdirectories:

```
working_directory
├── 01_data
│   └── raw
└── 02_geom
└── 03_proc
```

197 In this directory structure, the seismic shot files are stored in `01_data/raw` and the geometry file (`geometry.csv`) is
 198 provided in `02_geom`.

199 The geometry file is a csv file that provides an abstract representation of the survey layout and must not contain
 200 a header. The fundamental element for the description of the survey layout is the station, which refers either to a
 201 geophone position, a shot position or a position with co-located shot and geophone. For each station the geometric and
 202 semantic information described in Table 3 are provided column-wise, i.e., each line in the geometry file corresponds
 203 to a station.

to a single station with a unique position within the survey layout.

Table 3

Description of the information to be provided in the columns of the geometry file.

Col	Content	Data type	Description
1	x coordinate	float	Station x coordinate, e.g., given in (m)
2	y coordinate	float	Station y coordinate, e.g., given in (m)
3	z coordinate	float	Station z coordinate, e.g., given in (m)
4	Geophone	bool	1 if a geophone was deployed at station, 0 otherwise
5	Shot	int	The numerical part of the file name, e.g., 1001, if a shot was conducted at this station, -1 otherwise
6	First geophone	int	First active geophone (-1 if no shot was conducted at the station)
7	Number of geophones	int	Number of active geophones (-1 if no shot was conducted at the station)

204

In case the shot files as well as the geometry file are provided and a basic sanity check of the geometry file was
205 successful, the `SeismicRefractionManager` creates a new project:

```
from useis import SeismicRefractionManager
srm = SeismicRefractionManager('.')

INFO : Read geometry information from file
INFO : Extracted shot geometry
INFO : Extracted receiver geometry
INFO : Applied geometry
INFO : Standard pickset 'picks' created
INFO : Pickset 'picks' loaded
INFO : 'picks' set as active pickset
Progress <===== 100.0% completed>
INFO : Read 48 files
```

207

In particular, the `SeismicRefractionManager` creates an SQLite database (prj.db in the working directory) to store
208 the geometry information, whereby the stations are consecutively numbered (see Figure 2). To allow for an efficient
209 data selection through the user the `SeismicRefractionManager` links the station numbers to shot index numbers
210 (SIN) and receiver index numbers (RIN) assigned to shot and receiver stations, respectively (see Figure 2). Based on
211 this information, the geometry is applied, i.e., the database tables required for the processing of the seismic waveform
212 data are created. In particular, the first break traveltimes for each SIN-RIN pair are stored in a dedicated database table
213 `fbpicks` together with the name of the corresponding pickset, i.e., the a common label for an entire set of first break
214 traveltimes. By default, each project contains the default pickset 'picks', which is loaded and activated on startup.
215 Once the database is initialized, the waveform data are read from disk and the project is ready for processing.

If a database file is already present in the working directory, the project information, the seismic waveforms as well
216 as the default pickset 'picks' are automatically loaded:

```

from useis import SeismicRefractionManager

srm = SeismicRefractionManager('.')

INFO    : Project information loaded
Progress <===== 100.0% completed
INFO    : Read 48 files
INFO    : Pickset 'picks' loaded
INFO    : 'picks' set as active pickset

```

219

220 For a project with a successfully applied geometry, the `select` method of the `SeismicRefractionManager`
 221 allows to gather the seismic waveforms based on a common shot (`sin`), a common receiver (`rin`) or the common
 222 absolute offset (`aoffset`):

```

srm.select('sin 1')

INFO    : 48 traces selected

srm.select('rin 10')

INFO    : 48 traces selected

srm.select('aoffset 6')

INFO    : 90 traces selected

```

223

224 2.2.2. Visualization and enhancement of the seismic waveform data

225 Calling the `plot` method without passing any parameter allows for the visualization of the currently selected traces:

```
srm.plot()
```

226

227 Once opened, the seismogram plot visualizes the seismic waveforms in a combination of wiggle trace and variable area
 228 mode, i.e., the trace data are shown as curves. The area of the curves is colored red for negative and blue for positive
 229 amplitudes, respectively (not shown for brevity). Pressing the up or down arrow key on the keyboard toggles the
 230 visualization mode between variable area and variable density. In variable density mode the strength of the amplitudes
 231 is additionally reflected by the color saturation, i.e., high amplitudes refer to a stronger shade than low amplitudes (see
 232 Figure 3).

233 The active processing mode and data scaling mode are reported together with the traveltimes at the current cursor
 234 position in the status bar of the seismogram plot window (see Figure 4). The initial processing mode is 'Fb pick',
 235 i.e., first break picking is possible. The user can switch between the different modes by pressing specific keys on the
 236 keyboard. The 'm' key activates the trace mute mode ('Trc mute'), which allows to set the amplitude of a trace to
 237 zero by clicking with the left mouse button; clicking again on the same trace restores the amplitude information. The
 238 trace reverse mode ('Trc rev') is activated by pressing the 'r' key and enables the user to toggle the polarity of a trace
 239 by clicking on it with the left mouse button. The default data scaling mode is 'Zoom', which allows the scaling of
 240 the y-axis by turning the mouse wheel. By pressing the 'a' key the amplitude scaling mode ('Amp scal') is activated.
 241 Turning the mouse wheel increases or decreases the amplitudes of the traces currently shown in the seismogram plot,

242 and thus might help to enhance the perceptibility of the first onsets. By pressing the key of the currently active mode
 243 again, the `SeismicRefractionManager` returns to the default mode; yet, the different modes can be activated in any
 244 arbitrary order (as illustrated in Figure 4).

245 Through the `plot` method the frequency spectrum of the currently selected trace data can be visualized:

```
srm.plot('spectrum')
```

246

247 A frequency spectrum as shown in Figure 5 can be used to discriminate the dominating signal frequencies from those
 248 associated to the background noise, and thus allows for the definition of adequate filter settings. To improve the
 249 signal-to-noise ratio it is possible to apply filters on the selected traces through the `filter` method, which utilizes the
 250 frequency filters implemented in the ObsPy package (lowpass, highpass, bandpass and bandstop; Beyreuther et al.,
 251 2010), e.g.:

```
srm.filter('lp 120')
```

INFO : Applied 120.0 Hz lowpass filter

```
srm.filter('bp 10 120')
```

INFO : Applied bandpass filter (10.0 to 120.0 Hz)

```
srm.filter('hold on')
```

252

INFO : Set filter hold on

253 By default, filters are solely applied to the currently selected traces, yet setting the filter on hold allows the filtering of
 254 all subsequently selected traces with the same filter settings. In case the seismogram plot is opened, the effect of the
 255 applied filter on the seismic waveforms is interactively visualized.

256 2.2.3. Analysis of the seismic waveforms and first break traveltime picking

257 In the seismogram plot, the waveforms can be analyzed and processed to obtain information about the subsurface
 258 conditions. Activating the velocity estimation mode ('Vel est') by pressing the 'v' key on the keyboard, enables the
 259 user to estimate velocities for different wave phases (e.g., originating from a refractor) by pressing the left mouse button
 260 and moving the cursor. Once the left mouse button is released, a line between the start and end point is drawn and
 261 labeled with the corresponding velocity (estimates can be deleted by clicking with the right mouse button).

262 If the processing mode 'Fb pick' is activated, picking of first break traveltimes is done individually by clicking
 263 with the left mouse button on the respective trace. Clicking again on the same trace will set the first break pick to the
 264 new location as there can only be one travelttime for each SIN-RIN pair; whereas clicking with the right mouse button
 265 deletes the pick. Alternatively, first break picks can be set for multiple traces by pressing the left mouse button and
 266 moving the cursor. Once the left mouse button is released, first break picks are defined at the intersections between
 267 the line and the seismograms. In the same way, multiple picks can be deleted if a line is drawn with the right mouse
 268 button pressed. The first break traveltimes determined in the seismogram plot are automatically written to the project

269 database when the window is closed or another set of traces is loaded by pressing the 'left' or 'right' arrow keys on the
 270 keyboard.

271 The traveltimes diagram for the currently active pickset can be created through the plot method:

```
srm.plot('traveltimes')
```

272

273 Figure 6 presents an exemplary traveltimes diagram, which is a common way to examine the quality of the first break
 274 picking. Such illustration of the traveltimes can be used to identify outliers or erroneous measurements, which are
 275 commonly associated to traveltimes with substantial deviations from those observed at adjacent stations such as the
 276 first break pick for the SIN-RIN pair (23, 11). Outliers can be removed by clicking on the corresponding symbol ('x')
 277 in the traveltimes diagram, which is instantly synchronized with the project database. If the seismogram plot and the
 278 traveltimes diagram are used side-by-side, changes made to the first break picks in one window will interactively trigger
 279 an update of the other one and vice versa.

280 The SeismicRefractionManager handles first break picks in picksets, which can be organized and manipulated
 281 through the picksets method of the SeismicRefractionManager:

```
srm.picksets()
```

pickset	loaded	active
picks	Y	Y

```
srm.picksets('import picking_part1.pck pck_p1')
```

```
INFO    : Created new pickset 'pck_p1'  

INFO    : Pickset 'pck_p1' loaded  

INFO    : 'pck_p1' set as active pickset  

INFO    : Imported 'picking_part1.pck' to pickset 'pck_p1'
```

```
srm.picksets('copy pck_p1 pck_all')
```

```
INFO    : Created new pickset 'pck_all'  

INFO    : Pickset 'pck_all' loaded  

INFO    : 'pck_all' set as active pickset  

INFO    : Copied pickset 'pck_p1' to 'pck_all'
```

```
srm.picksets('unload pck_p1 picks')
```

```
INFO    : Pickset 'pck_p1' removed from workflow  

INFO    : Pickset 'picks' removed from workflow
```

```
srm.picksets('delete pck_p1')
```

```
INFO    : Pickset 'pck_p1' deleted
```

```
srm.picksets()
```

pickset	loaded	active
picks	N	N
pck_all	Y	Y

282

283 Calling the pickset method without parameters shows the status of all picksets in the project. From the above use
 284 case, we see that by default, the pickset 'picks' is loaded and activated, i.e., modifications of first breaks are stored in
 285 this pickset. First break picks provided by another source (as udf file) can be imported from *03_proc/picks*. For the

- 286 first break picking, it is sufficient to keep only one pickset in the workflow, i.e., not required picksets can be unloaded.
- 287 A pickset currently not loaded can be deleted permanently, i.e., the corresponding traveltimes are removed from the
- 288 database. The `picksets` method can also be used to load picksets from the database:

```
srm.picksets('load picks')
```

INFO : Pickset 'picks' loaded

```
srm.picksets()
```

pickset	loaded	active
picks	Y	N
pck_all	Y	Y

```
srm.picksets('use picks')
```

INFO : Pickset 'picks' loaded

INFO : 'picks' set as active pickset

```
srm.picksets()
```

pickset	loaded	active
picks	Y	Y
pck_all	Y	N

289

- 290 When a pickset is loaded from the database, it does not become the active pickset automatically; whereas by using the
- 291 parameter `use` the corresponding pickset is loaded and also becomes the active pickset. The first break traveltimes of a
- 292 pickset can be exported to an udf file that is stored in `03_proc/picks` subdirectory with the current timestamp as suffix:

```
srm.picksets('export pck_all')
```

INFO : pickset 'pck_all' saved to pck_all_20220428T145648.pck

293

294 3. Exemplary use cases

295 3.1. Modeling and use of synthetic seismic waveform data

296 To demonstrate the applicability of the `SeismicWaveformModeler` class, we generate two synthetic seismic wave-

297 form datasets assuming two horizontal layers in a half-space without topography and an interface between these layers

298 with a constant depth. Table 4 summarizes the parameterization used for the forward modeling of one dataset without

299 noise (Dataset_1) and another dataset containing random noise and systematic errors (Dataset_2). For the visualization

300 of both synthetic datasets, we provide the corresponding geometry files to the `SeismicRefractionManager` in order

301 to have full control regarding data selection and processing capabilities.

302 The synthetic seismic waveforms for SIN 24 of Dataset_1 presented in Figure 7 reveal clear negative first onsets and

303 allow the identification of crossover points, i.e., inflexion points between the first arrivals from the first layer (RIN 17

304 to 30) and the first onsets associated to the second layer (RIN 1 to 16 and RIN 31 to 48). In contrast the signal-to-noise

305 ratio of Dataset_2 (see Figure 8) is a function of the offset between shot and geophone position, i.e., traces farther away

306 from the shot contain a higher level of random noise. Moreover, Dataset_2 also contains systematic errors, where RIN

307 6 and 14 refer to broken geophones, and RIN 35 is an example for readings with a wrong polarity. We believe that such
 308 synthetic datasets might be useful for investigating the effect of complex survey geometries on seismic data as well as
 the development and evaluation of processing strategies.

Table 4

Measurement scheme and parameters provided in the yaml files used to create synthetic seismic waveform datasets with (Dataset_1) and without added noise (Dataset_2).

Measurement scheme		
Number of stations	48	
Station spacing	2 m	
Number of geophones	48	
Number of shots	48	
Model	Layer 1	Layer 2
Thickness	3 m	10 m
Velocity	750 m/s	3000 m/s
Dataset	Dataset_1	Dataset_2
Noise	False	True
Noise level	0	1e-5
Missing shots	False	True
Broken geophones	False	True
Wrong polarity	False	True
Wavelet		
Length	1.024 s	
Frequency	100 Hz	
Sampling rate	2000 Hz	

309
 310 Accordingly, the concept of the formikoj library allows the implementation of supplementary functionalities. Such
 311 custom extensions should be implemented either as internal methods or as functions in the utilities module, which are
 312 then executed through the `compute` method with a custom keyword. As an illustrative example, we implemented a
 313 simplified version of an automatic first break picking algorithm, which determines the traveltimes based on the energy
 314 ratio method (e.g., Earle and Shearer, 1994). The algorithm is added to the `SeismicRefractionManager` in form
 315 of two internal methods `_manage_autopicking` and `_compute_autopicks`, respectively. The autopicking process
 316 can be started by passing the new keyword `autopick` as the first parameter to the `compute` method:

```
srm.compute('autopick all autopicks')

INFO : Created new pickset 'autopicks'
INFO : Pickset 'autopicks' loaded
INFO : 'autopicks' set as active pickset
Progress <===== 100.0% completed
```

317
 318 The second parameter defines whether traveltimes should be determined for the entire dataset (`all`) or solely for
 319 the currently selected traces (`cur`). The third parameter provides the name of the corresponding pickset in the project
 320 database; thus, making the traveltimes available for visualization and processing with existing functionalities or further
 321 custom implementations.

322 3.2. First break travel time picking for a 2D roll-along field data set: the Danube island example

323 The seismic data used in this example were collected at the Danube island (Vienna) in June 2021, using 48 geo-
 324 phones deployed with 2 m spacing between them and shot locations located between the geophone positions. As
 325 illustrated in Figure 9, the survey geometry refers to a roll-along geometry, i.e., the geophone spread was moved along
 326 a profile with 50% overlap yielding a total of five segments. The objective of the survey was to define the contact
 327 between different sediments within the tertiary and quaternary deposits used to build the man-made Danube island.
 328 Additionally, the survey aimed to identify lateral changes that might indicate the position of a fault, which has been
 329 inferred from sediments recovered from drillings.

Table 5

Extract from the roll-along survey geometry file showing how the information regarding the first geophone assigns the traces in the shot files to the correct stations.

x (m)	y (m)	z (m)	Geo	Shot	1st Geo	# Geo
0.0	0.0	163.5	1	-1	-1	-1
2.0	0.0	163.5	0	1001	1	48
:	:	:	:	:	:	:
94.0	0	163.5	0	1024	1	48
96.0	0	163.5	1	-1	-1	-1
98.0	0.0	163.5	0	1037	25	48
:	:	:	:	:	:	:
194.0	0.0	163.5	0	1049	25	48
196.0	0.0	163.5	1	-1	-1	-1
198.0	0.0	163.5	0	1073	49	48
200.0	0.0	163.5	1	-1	-1	-1
202.0	0.0	163.5	0	1050	25	48
:	:	:	:	:	:	:
286.0	0.0	163.5	0	1084	49	48
288.0	0.0	163.5	1	-1	-1	-1
290.0	0.0	163.5	0	1097	73	48
:	:	:	:	:	:	:
386.0	0.0	163.5	0	1109	73	48
388.0	0.0	163.5	1	-1	-1	-1
390.0	0.0	163.5	0	2025	97	48
:	:	:	:	:	:	:
570.0	0.0	163.5	0	2048	97	48
572.0	0.0	163.5	1	-1	-1	-1

330 In the field, each segment was measured separately, yet for the processing all measurements are combined in a single
 331 profile. We can implement such measurement layout in a geometry file as illustrated in Table 5. The key parameter is
 332 '1st Geo' as it indicates the first active geophone along the profile for each shot file. For the first segment in a roll-along
 333 survey geometry, the first active geophone is always geophone 1. Considering the number of geophones used in each
 334 segment, and an overlap of 50%, the first geophone for segments two to five are 25, 49, 73 and 97, respectively.

335 Based on the shot files and the geometry file stored in the required directory structure the SeismicRefractionManager
 336 creates the project database. A first illustration of the subsurface conditions can be obtained by computing a common

337 offset stack (COS):

```
srm.compute('cos')

Progress <===== 100.0% completed
INFO    : Computed the common offset stack
```

338

339 The COS for the Danube island dataset presented in Figure 10 shows first onsets for absolute offsets up to approxi-
 340 mately 150 m, which suggest a two-layered subsurface model. By using the velocity estimation functionality we can
 341 approximate the seismic velocity in the corresponding layers.

342 For the first break picking a set of traces is selected, filtered if necessary and visualized:

```
srm.select('sin 99')

INFO    : 48 traces selected

srm.filter('lp 120')

INFO    : Applied 120.0 Hz lowpass filter

srm.plot()
```

343

344 In this example, the trace data for SIN 99 are used, which refers to a shot position located in segment four. As can
 345 be seen in Figure 11, the first onsets are easily identifiable despite the substantial seismic background noise at large
 346 offsets (RIN 73 to 90). A pseudosection provides an illustration of the corresponding apparent seismic velocity values
 347 computed from the picked traveltimes and the distance between shots and geophones:

```
srm.plot('pseudosection')
```

348

349 In a pseudosection, as presented in Figure 12 for the Danube island dataset, the apparent velocity values are assigned
 350 to pseudolocations, with the corresponding x- and z-coordinates determined as the midpoint and as 1/3 of the absolute
 351 offset between the shot and geophone, respectively. Accordingly, such plot allows for the identification of outliers in
 352 the data, e.g., stark velocity contrasts for adjacent points, or systematic errors, e.g., velocities erroneously influenced
 353 by a single shot or receiver. The main assumption here is that the pseudosection should reveal smooth transitions
 354 between lateral and vertical neighbors, considering that the data were collected with gradual changes in the position of
 355 the source (i.e., hammer blow) and the receiver (i.e., geophone). The pseudosection will reveal large variations in case
 356 of abrupt changes in the topography or the geometry of the array, yet this can be taken into account by the user during
 357 the identification of outliers and possible systematic errors. For the Danube island dataset, the pseudosection suggests
 358 an increase in the seismic velocity along profile direction in deeper subsurface regions (i.e., larger pseudodepth). Such
 359 pattern could be related to a fault expected in this area of the Danube island, thus indicating that the geophysical survey
 360 was sufficiently designed to detect such feature. Moreover, the pseudosection presented in Figure 12 shows a low
 361 number of data points in the first segment of the Danube island survey. This lack of data points at large pseudodepths
 362 is due to a low number of picked traveltimes at large offsets, and thus might indicate a low signal-noise ratio.

363 To review the data quality along the entire profile it is possible to visualize the picking percentage, i.e., the ratio of
 364 actually picked traveltimes and total number of SIN-RIN pairs:

```
365 srm.plot('pickperc')
```

366 Figure 13 shows that the picking percentage is visualized separately for each SIN. In this way, a low picking percentage
 367 indicates shots affected by a low signal-to-noise ratio. For the Danube island dataset, the picking percentage is low
 368 in the first segment, which confirms the lack of datapoints observed in the pseudosection. Moreover, the picking
 369 percentage plot can be used to track the picking progress, for instance if the traveltimes cannot be determined in one
 370 session or to identify single shots that might have been forgotten during the first break picking. Accordingly, it is
 371 advisable to check the picking percentage prior to exporting the traveltimes for the inversion.

372 3.3. Processing of a 3D seismic refraction dataset: the soda lake example

373 The `SeismicRefractionManager` allows also the visualization and processing of data collected with a 3D survey
 374 layout. To illustrate these capabilities, we present in Figure 14 the geometry of a 3D survey conducted in a soda lake
 375 located close to Vienna. The soda lake corresponds to quaternary sediments where capillary forces have developed a
 376 low permeable layer close to the surface (between 50 and 100 cm) with a high clay and salt content. The seismic survey
 377 aims to support the interpretation of the electrical and electromagnetic models obtained in a monitoring framework.
 378 Accordingly, the survey geometry shown in Figure 14 was specified by previously conducted electrical measurements
 379 with electrodes arranged along two perpendicular lines. The seismic data were collected with 48 geophones deployed
 380 along the North-East to South-West oriented line, and 48 geophones along the North-West to South-East oriented line,
 381 with a spacing of 2 m between the geophones. Shots were generated with an 8 kg sledgehammer at the geophone
 382 positions as well as at positions along the diagonals.

383 The geometry file contains the 3D coordinates for each shot or receiver station, and thus the
 384 `SeismicRefractionManager` automatically configures the project for 3D processing. Figure 15 presents the seismic
 385 waveform data recorded for SIN 1, i.e., the shot position co-located with the first geophone (Station 01 in Figure 14).
 386 The data for RIN 1 to 48 appear familiar as the corresponding SIN-RIN layout refers to the one of conventional
 387 2D profiles; whereas the seismic waveforms for RIN 49 to 96 show an entirely different pattern. To understand such
 388 visualization, we need to take into account that RIN 49 to 96 are deployed perpendicular to the direction of propagation
 389 of the waveform originating from SIN 1. Accordingly, the observed curvature in the first onsets is due to the varying
 390 offset of the different SIN-RIN pairs.

391 Due to the 3D survey geometry, a 2D pseudosection is not suitable to illustrate the apparent velocity values obtained
 392 from the picked first break traveltimes. Hence, the `SeismicRefractionManager` automatically switches to a 3D
 393 representation where the apparent velocity values are illustrated in an interactive 3D pseudosection. This plot can be

394 rotated and the image section can be zoomed and panned allowing the user to easily investigate the data quality for
 395 3D geometries. Figure 16 shows a screenshot of the 3D pseudosection for the salt lake dataset; yet, such screenshot
 396 cannot reveal the full capabilities implemented in the **SeismicRefractionManger** for the interactive analysis and
 397 visualization of 3D pseudosections.

398 Once the first break picking is finished, the corresponding pickset can be exported for the inversion. The inversion
 399 results and their interpretation are not the scope of this manuscript, yet Figures 15 and 16 reveal the capabilities
 400 provided by the proposed framework for the visualization and processing of data collected in 3D survey geometries.

401 4. Conclusions and Outlook

402 We have presented formikoj, a flexible open-source library enabling the development of modeling and processing
 403 tools for geophysical data. Implemented in python and tested on all major operating systems (Linux/Unix, MacOS,
 404 Windows), formikoj is suitable for multi- and cross-platform applications; thus, allowing collaboration between users
 405 free from licensing costs and platform requirements.

406 We demonstrated the capabilities of the formikoj framework to develop versatile and easily scalable classes for
 407 the modeling and processing of waveforms in seismic refraction surveys. The required interaction with the user
 408 is reduced to a minimum as crucial processing steps are automatized within the **SeismicWaveformModeler** and
 409 **SeismicRefractionManager** based on efficient data input strategies, for instance the preparation and import of the
 410 geometry file or the keyboard-based interaction related to the first break picking. In this regard, the user controls the
 411 formikoj library by providing text-based commands preferably through an ipython shell to exploit the full interactive
 412 potential modeling and processing tools. However, applications of the formikoj library can also be automatized by
 413 implementing workflows in python scripts or jupyter notebooks.

414 Based on three exemplary use cases, we illustrated the applicability of both the **SeismicWaveformModeler** and
 415 the **SeismicRefractionManager**. In the first use case, we showed the possibility to forward model seismic wave-
 416 form data based on custom subsurface models and survey geometries with the **SeismicWaveformModeler**. Addi-
 417 tionally, the resulting waveforms can be subjected to systematic and random noise sources. The capabilities of the
 418 **SeismicRefractionManager** were demonstrated through the processing of field datasets collected in complex sur-
 419 vey layout, namely a roll-along and a 3D geometry. Moreover, we showed how the different data visualization options
 420 can assist during the data processing to ensure consistency in the first break traveltimes. In particular, we developed a
 421 visualization of the traveltimes by means of pseudosections illustrating the corresponding apparent seismic velocities.
 422 Such plots allow for a quick identification of systematic errors and outliers in both 2D and 3D datasets.

423 By making the source code of the formikoj library available under the MIT license we intend to spark the develop-
 424 ment of further modeling and processing tools for various geophysical models based on this framework. Our further

425 efforts will focus on implementing tools for other wave-based geophysical methods used in frame of our research
426 activities, such as multi-channel analysis of (seismic) surface waves or magnetotelluric surveys.

427 **5. Acknowledgments**

428 The authors are grateful to Nathalie Roser and Lukas Aigner for benchmarking the formikoj library against estab-
429 lished software packages in frame of their research activities. Furthermore, we would like to thank Clemens Moser,
430 Martin Mayr, Vinzenz Schichl and Harald Pammer for their constructive comments during first tests of the formikoj
431 framework as well as for their help during the seismic surveys.

432 Code and data availability section**433 Name of the code/library: formikoj****434 Contact: matthias.steiner@geo.tuwien.ac.at****435 Hardware requirements: No specific requirements****436 Program language: Python****437 Software required: Anaconda/Miniconda recommended****438 Total program and dataset size: 250 MB****439 The source codes and exemplary data sets are available for downloading at the link:****440 <https://git.geo.tuwien.ac.at/msteiner/formikoj.git>****441 The orthophotos used in Figures 9 and 14 were published by geoland.at under a CC BY 4.0 license.****442 References****443 Ahern, T., Casey, R., Barnes, D., Benson, R., Knight, T., Trabant, C., 2012. SEED Reference Manual, Version 2.4. URL: http://www.fdsn.org/pdf/SEEDManual_V2.4.pdf.****444 Beyreuther, M., Barsch, R., Krischer, L., Megies, T., Behr, Y., Wassermann, J., 2010. Obspy: A python toolbox for seismology. Seismological Research Letters 81, 530–533. doi:10.1785/gssrl.81.3.530.****445 Blanchy, G., Saneian, S., Boyd, J., McLachlan, P., Binley, A., 2020. Resipy, an intuitive open source software for complex geoelectrical inversion/modeling. Computers & Geosciences 137, 104423. URL: <https://www.sciencedirect.com/science/article/pii/S0098300419308192>, doi:<https://doi.org/10.1016/j.cageo.2020.104423>.****446 Bücker, M., Flores Orozco, A., Gallistl, J., Steiner, M., Aigner, L., Hoppenbrock, J., Glebe, R., Morales Barrera, W., Pita de la Paz, C., García García, C.E., et al., 2021. Integrated land and water-borne geophysical surveys shed light on the sudden drying of large karst lakes in southern mexico. Solid Earth 12, 439–461. doi:10.5194/se-12-439-2021.****447 Cockett, R., Kang, S., Heagy, L.J., Pidlisecky, A., Oldenburg, D.W., 2015. Simpeg: An open source framework for simulation and gradient based parameter estimation in geophysical applications. Computers & Geosciences 85, 142–154. URL: <https://www.sciencedirect.com/science/article/pii/S009830041530056X>, doi:<https://doi.org/10.1016/j.cageo.2015.09.015>.****448 Draebing, D., 2016. Application of refraction seismics in alpine permafrost studies: A review. Earth-Science Reviews 155, 136–152. doi:<https://doi.org/10.1016/j.earscirev.2016.02.006>.****449 Earle, P.S., Shearer, P.M., 1994. Characterization of global seismograms using an automatic-picking algorithm. Bulletin of the Seismological Society of America 84, 366–376.****450 Guedes, V.J.C.B., Maciel, S.T.R., Rocha, M.P., 2022. Refrapy: A python program for seismic refraction data analysis. Computers & Geosciences 159, 105020. URL: <https://www.sciencedirect.com/science/article/pii/S0098300421003022>, doi:<https://doi.org/10.1016/j.cageo.2021.105020>.****451 Harris, C.R., Millman, K.J., van der Walt, S.J., Gommers, R., Virtanen, P., Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N.J., Kern, R., Picus, M., Hoyer, S., van Kerkwijk, M.H., Brett, M., Haldane, A., del Río, J.F., Wiebe, M., Peterson, P., Gérard-Marchant, P., Sheppard, K., Reddy, T., Weckesser, W., Abbasi, H., Gohlke, C., Oliphant, T.E., 2020. Array programming with NumPy. Nature 585, 357–362. URL: <https://doi.org/10.1038/s41586-020-2649-2>, doi:10.1038/s41586-020-2649-2.**

- 467 Heimann, S., Kriegerowski, M., Isken, M., Cesca, S., Daout, S., Grigoli, F., Juretzek, C., Megies, T., Nooshiri, N., Steinberg, A., Sudhaus, H.,
468 Vasyura-Bathke, H., Willey, T., Dahm, T., 2017. Pyrocko - an open-source seismology toolbox and library doi:10.5880/GFZ.2.1.2017.001.
- 469 Hunter, J.D., 2007. Matplotlib: A 2d graphics environment. Computing in Science & Engineering 9, 90–95. doi:10.1109/MCSE.2007.55.
- 470 McKinney, W., 2010. Data Structures for Statistical Computing in Python, in: Stéfan van der Walt, Jarrod Millman (Eds.), Proceedings of the 9th
471 Python in Science Conference, pp. 56 – 61. doi:10.25080/Majora-92bf1922-00a.
- 472 Nguyen, F., Ghose, R., Isunza Manrique, I., Robert, T., Dumont, G., 2018. Managing past landfills for future site development: A review of the
473 contribution of geophysical methods, in: Proceedings of the 4th International Symposium on Enhanced Landfill Mining, pp. 27–36.
- 474 Parsekian, A., Singha, K., Minsley, B.J., Holbrook, W.S., Slater, L., 2015. Multiscale geophysical imaging of the critical zone. Reviews of
475 Geophysics 53, 1–26. doi:10.1002/2014RG000465.
- 476 Plattner, A.M., 2020. Gprpy: Open-source ground-penetrating radar processing and visualization software. The Leading Edge 39, 332–337.
477 doi:10.1190/tle39050332.1.
- 478 Ringler, A.T., Evans, J.R., 2015. A quick seed tutorial. Seismological Research Letters 86, 1717–1725. doi:10.1785/0220150043.
- 479 Romero-Ruiz, A., Linde, N., Keller, T., Or, D., 2018. A review of geophysical methods for soil structure characterization. Reviews of Geophysics
480 56, 672–697. doi:10.1029/2018RG000611.
- 481 Rücker, C., Günther, T., Wagner, F.M., 2017. pygimli: An open-source library for modelling and inversion in geophysics. Computers & Geosciences
482 109, 106–123. doi:10.1016/j.cageo.2017.07.011.
- 483 Steiner, M., Katona, T., Fellner, J., Flores Orozco, A., 2022. Quantitative water content estimation in landfills through joint inversion of seismic re-
484 fraction and electrical resistivity data considering surface conduction. Waste Management 149, 21–32. URL: <https://www.sciencedirect.com/science/article/pii/S0956053X22002793>, doi:<https://doi.org/10.1016/j.wasman.2022.05.020>.
- 485 Steiner, M., Wagner, F.M., Maierhofer, T., Schöner, W., Flores Orozco, A., 2021. Improved estimation of ice and water contents in alpine permafrost
486 through constrained petrophysical joint inversion: The hoher sonnblick case study. Geophysics 86, 1–84. doi:10.1190/geo2020-0592.1.
- 487 Stockwell, J.W., 1999. The cwp/su: Seismic unix package. Computers & Geosciences 25, 415–419. URL: <https://www.sciencedirect.com/science/article/pii/S0098300498001459>, doi:10.1016/S0098-3004(98)00145-9.
- 488 Sullivan, C.B., Kaszynski, A., 2019. PyVista: 3d plotting and mesh analysis through a streamlined interface for the visualization toolkit (VTK).
489 Journal of Open Source Software 4, 1450. URL: <https://doi.org/10.21105/joss.01450>, doi:10.21105/joss.01450.
- 490 Uieda, L., Oliveira Jr, V.C., Barbosa, V.C., 2013. Modeling the earth with fatiando a terra, in: Proceedings of the 12th Python in Science Conference,
491 pp. 96–103.

494 List of Figures

495 1	General architecture of the formikoj library comprising a utility, modeling and processing module. The base classes <code>DataModeler</code> and <code>MethodManager</code> can be used to build tools for specific geophysical methods, e.g., seismic refraction.	24
496 2	The <code>SeismicRefractionManager</code> addresses the stations through consecutive station numbers based on the sort order in the geometry file. The shot index numbers (SIN) and receiver index numbers (RIN) are assigned to the shot and receiver stations separately to allow for an intuitive data handling for the user.	25
497 3	The seismogram plot presents the currently selected traces along the x axis, where the sort order is determined through the geometry. The corresponding trace data are illustrated as a function of time along the y axis (solid curves), with positive and negative amplitudes depicted in blue and red, respectively. The selection criterion and the applied filter are shown in the upper left corner of the plot. Green crosses refer to the picked traveltimes stored in the currently active pickset.	26
498 4	The status bar in the interactive seismogram plot displays the currently active processing and data scaling modes as well as the time (in seconds) at the current cursor position. By pressing the keys 'a', 'm', 'r', 'v' on the keyboard the different modes can be activated.	27
499 5	The frequency spectrum illustrates the frequency content of the currently selected traces, which allows for the identification of frequency ranges associated to noise, which can be omitted through the application of corresponding frequency filtering.	28
500 6	The travelttime diagram shows the first break traveltimes stored in the currently active pickset along the y axis (x symbols), where solid lines connect traveltimes assigned to a common shot. The sort order of the stations along the x axis reflects the geometry. Filled circles indicate stations with co-located shot and geophone (receiver), triangles refer to receiver stations (no shot) and stars refer to shot stations (no geophone).	29
501 7	Synthetic seismic waveform data without random or systematic errors created with the <code>SeismicWaveformModeler()</code> class for a shot position in the center of the survey layout.	30
502 8	Synthetic seismic waveform data with added noise created with the <code>SeismicWaveformModeler()</code> class for a shot position in the center of the survey layout. The random noise refers to an offset dependent decrease of the signal-to-noise ratio, while the systematic broken geophones and wrong polarity are systematic errors.	31
503 9	The Danube island field dataset was collected along a single line with a roll-along survey layout; the filled triangles indicate the direction of the measurements. The five segments have an overlap of 50% to ensure an adequate data coverage along the entire profile.	32
504 10	The common offset stack computed for the Danube island dataset clearly showing a two-layered subsurface. The seismic velocities within the layer can be estimated from the gradient of the lines drawn along the corresponding first onsets of the seismic waves.	33
505 11	Exemplary seismic waveforms from the Danube island dataset shown for shot index number (SIN) 99 with a 120 Hz lowpass filter applied to suppress high frequency noise. The receiver index numbers (RIN) start at 73, thus indicating that the data were collected with a roll-along survey geometry.	34
506 12	Pseudosection showing the apparent seismic velocities determined from the first break traveltimes obtained from the Danube island dataset and the corresponding absolute offset between the shot and receiver stations. The apparent velocity for each shot-receiver pair is illustrated at the corresponding midpoint and pseudodepth (1/3 of the absolute offset).	35
507 13	Picking percentage for each shot in the Danube island roll-along dataset. Low values in the first third of the dataset indicate a low signal-to-noise ratio in the seismograms hindering the picking of first break traveltimes for each shot-receiver pair.	36
508 14	The soda lakes field dataset was collected in a 3D survey layout with stations (co-located geophones and shots indicated by filled dots) deployed in form of a cross. Additional shots (yellow stars) were conducted in front of a cross rotated by 45° to increase the coverage of the dataset.	37

- 543 15 Examplary seismic waveforms from the soda lakes dataset shown for shot index number (SIN) 1 with
544 a 100 Hz lowpass filter applied to suppress high frequency noise. The recorded seismic waveforms
545 clearly reflect the geometry of the geophones with RIN 1 to 48 deployed along the direction of wave
546 propagation, while RIN 49 to 96 are deployed perpendicular to the propagating wavefront. 38
- 547 16 3D pseudosection showing the apparent seismic velocities determined from the first break traveltimes
548 obtained from the soda lakes dataset and the corresponding absolute offset between the shot and re-
549 ceiver stations. The apparent velocity for each shot-receiver pair is illustrated at the corresponding 2D
550 midpoint and pseudodepth (1/3 of the absolute offset), thus yielding a 3D representation. 39

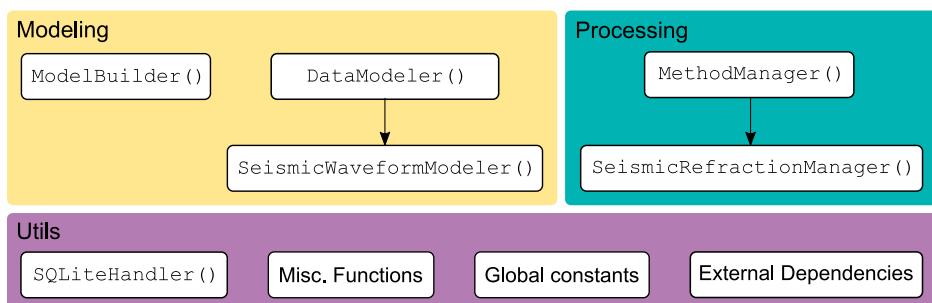


Figure 1: General architecture of the formikoj library comprising a utility, modeling and processing module. The base classes DataModeler and MethodManager can be used to build tools for specific geophysical methods, e.g., seismic refraction.

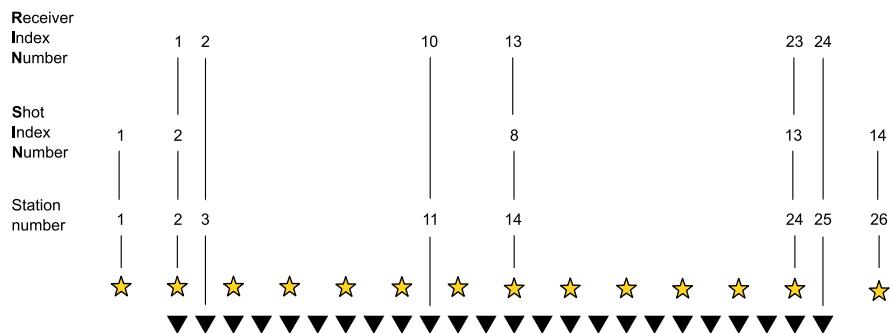


Figure 2: The SeismicRefractionManager addresses the stations through consecutive station numbers based on the sort order in the geometry file. The shot index numbers (SIN) and receiver index numbers (RIN) are assigned to the shot and receiver stations separately to allow for an intuitive data handling for the user.

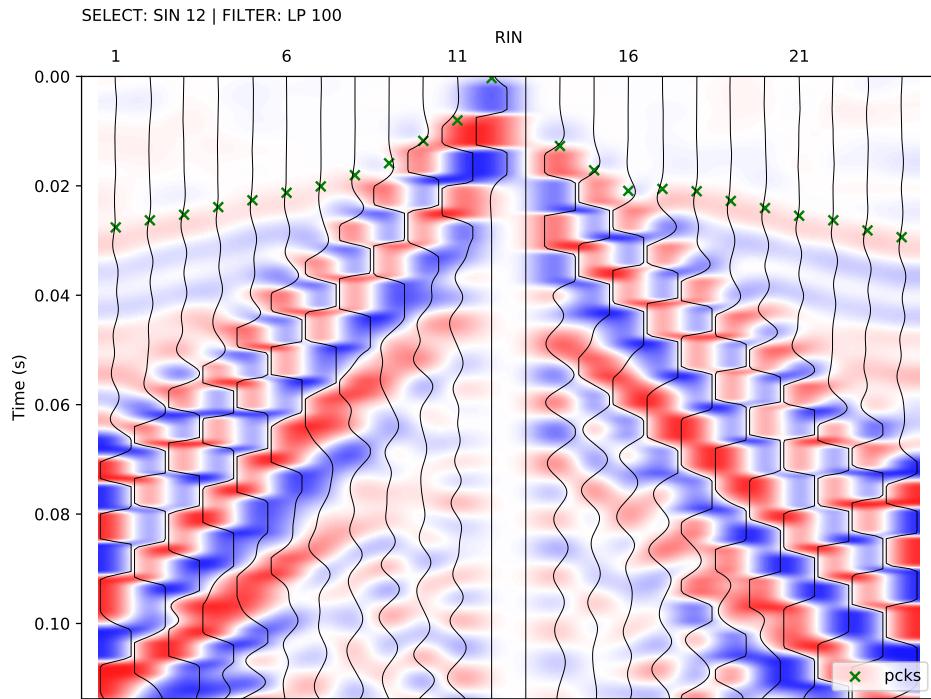


Figure 3: The seismogram plot presents the currently selected traces along the x axis, where the sort order is determined through the geometry. The corresponding trace data are illustrated as a function of time along the y axis (solid curves), with positive and negative amplitudes depicted in blue and red, respectively. The selection criterion and the applied filter are shown in the upper left corner of the plot. Green crosses refer to the picked traveltimes stored in the currently active pickset.

Proc: Fb pick | Scroll: Zoom | Time (s) = 0.000



Figure 4: The status bar in the interactive seismogram plot displays the currently active processing and data scaling modes as well as the time (in seconds) at the current cursor position. By pressing the keys 'a', 'm', 'r', 'v' on the keyboard the different modes can be activated.

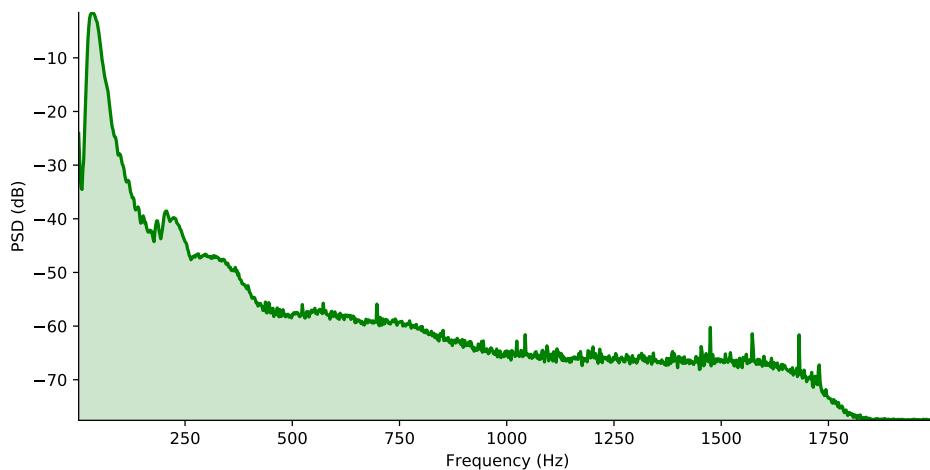


Figure 5: The frequency spectrum illustrates the frequency content of the currently selected traces, which allows for the identification of frequency ranges associated to noise, which can be omitted through the application of corresponding frequency filtering.

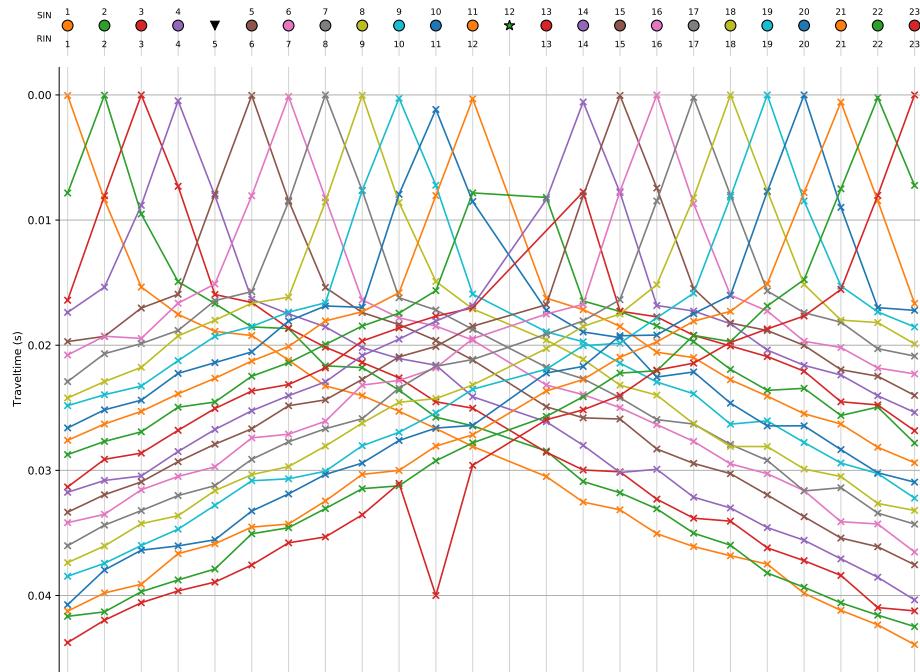


Figure 6: The traveltimes diagram shows the first break traveltimes stored in the currently active pickset along the y axis (x symbols), where solid lines connect traveltimes assigned to a common shot. The sort order of the stations along the x axis reflects the geometry. Filled circles indicate stations with co-located shot and geophone (receiver), triangles refer to receiver stations (no shot) and stars refer to shot stations (no geophone).

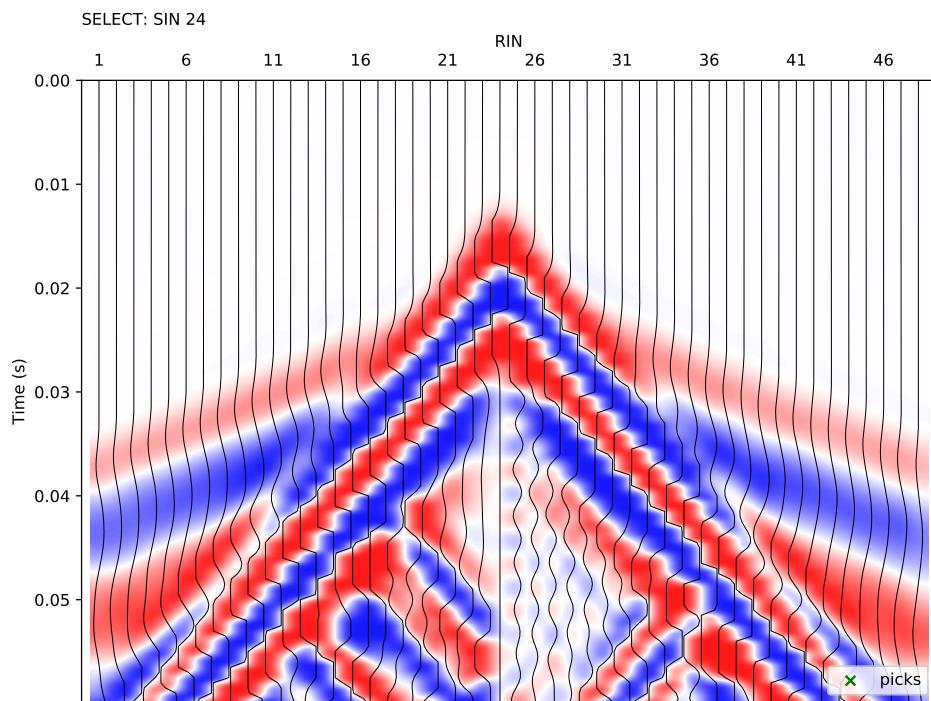


Figure 7: Synthetic seismic waveform data without random or systematic errors created with the `SeismicWaveformModeler()` class for a shot position in the center of the survey layout.

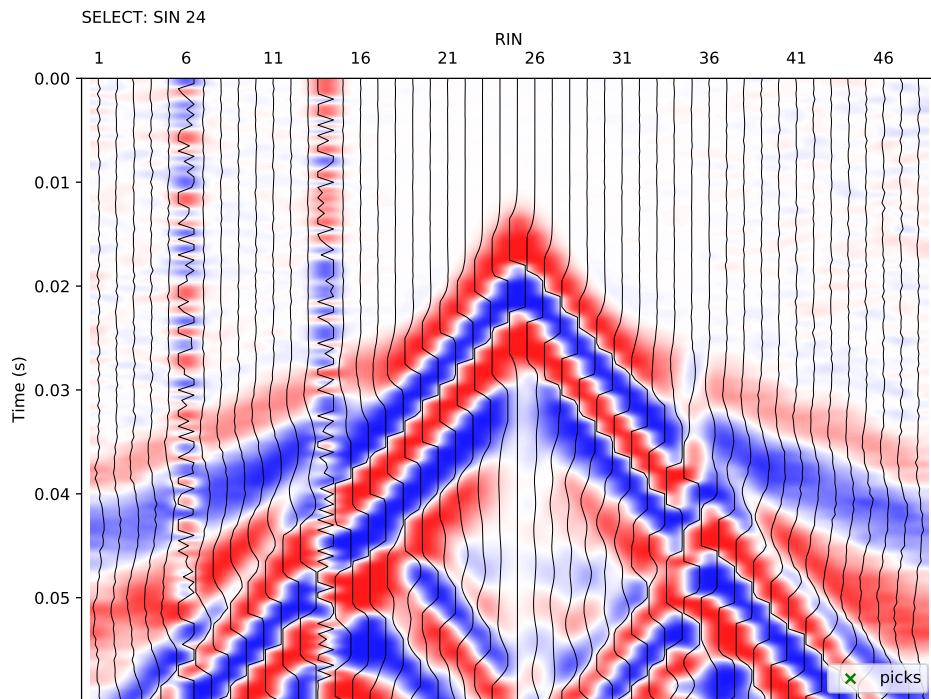


Figure 8: Synthetic seismic waveform data with added noise created with the `SeismicWaveformModeler()` class for a shot position in the center of the survey layout. The random noise refers to an offset dependent decrease of the signal-to-noise ratio, while the systematic broken geophones and wrong polarity are systematic errors.

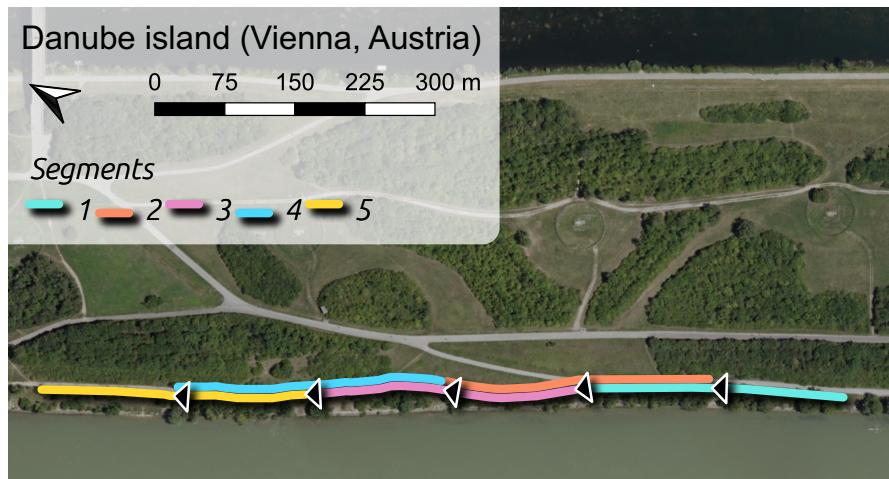


Figure 9: The Danube island field dataset was collected along a single line with a roll-along survey layout; the filled triangles indicate the direction of the measurements. The five segments have an overlap of 50% to ensure an adequate data coverage along the entire profile.

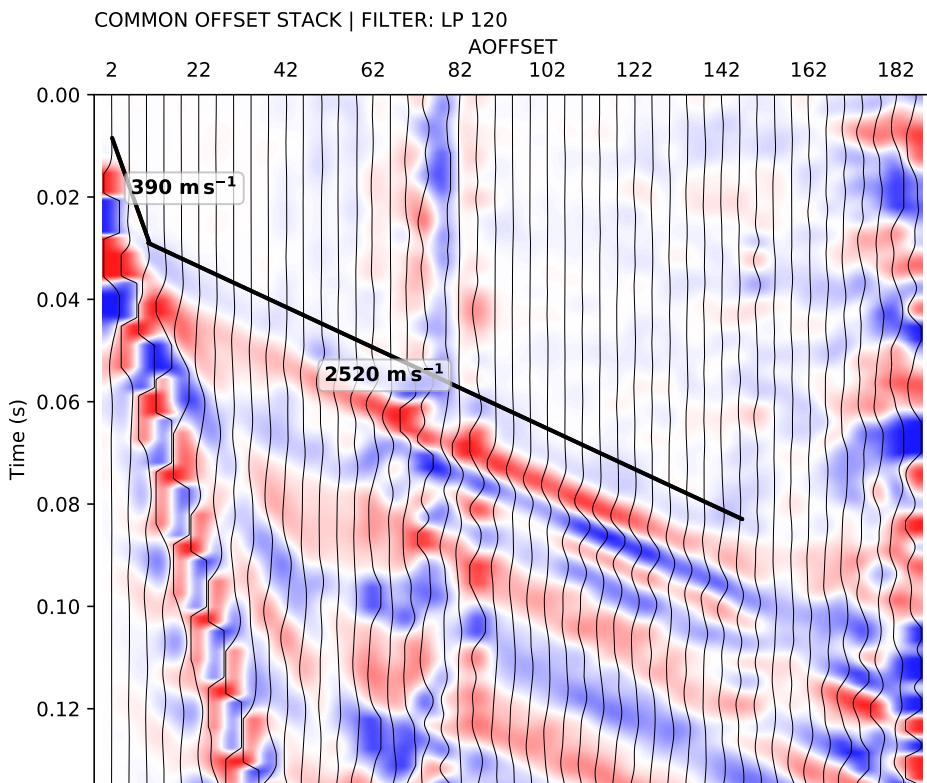


Figure 10: The common offset stack computed for the Danube island dataset clearly showing a two-layered subsurface. The seismic velocities within the layer can be estimated from the gradient of the lines drawn along the corresponding first onsets of the seismic waves.

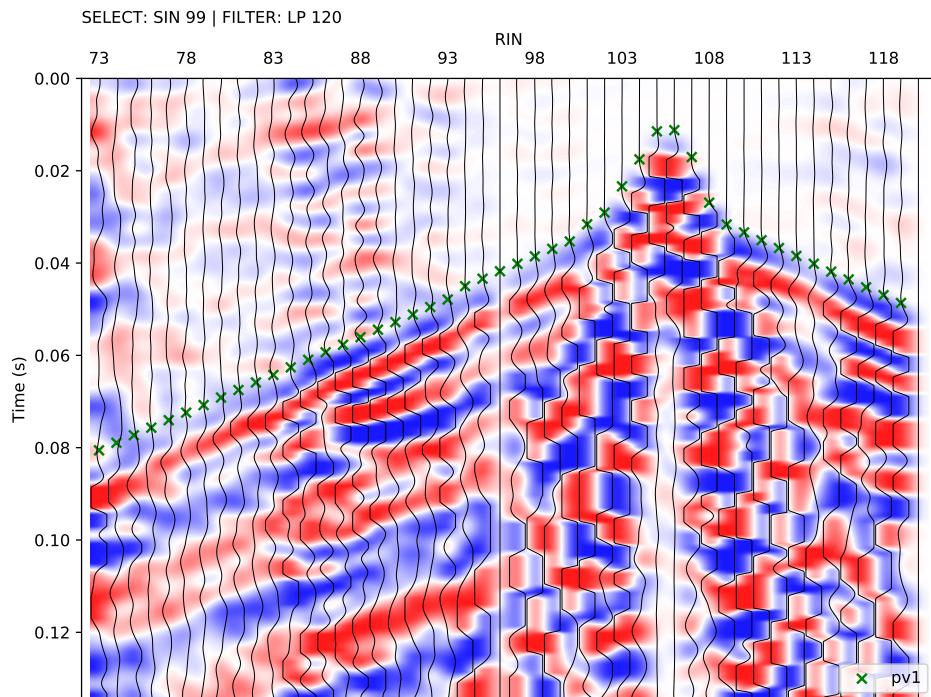


Figure 11: Exemplary seismic waveforms from the Danube island dataset shown for shot index number (SIN) 99 with a 120 Hz lowpass filter applied to suppress high frequency noise. The receiver index numbers (RIN) start at 73, thus indicating that the data were collected with a roll-along survey geometry.

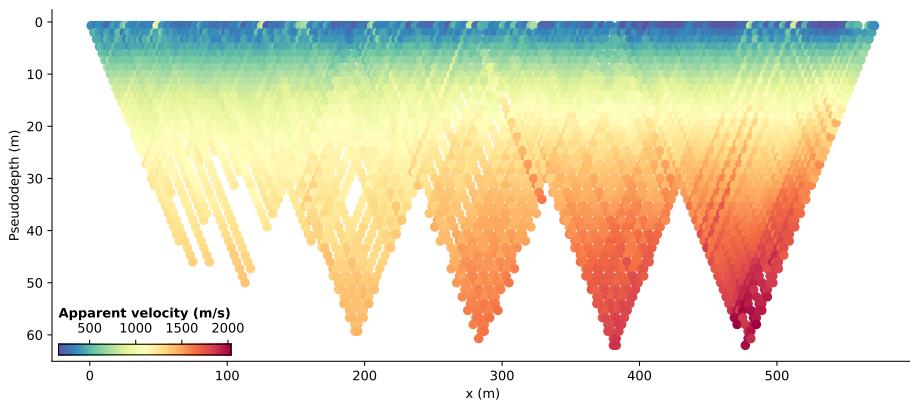


Figure 12: Pseudosection showing the apparent seismic velocities determined from the first break traveltimes obtained from the Danube island dataset and the corresponding absolute offset between the shot and receiver stations. The apparent velocity for each shot-receiver pair is illustrated at the corresponding midpoint and pseudodepth (1/3 of the absolute offset).

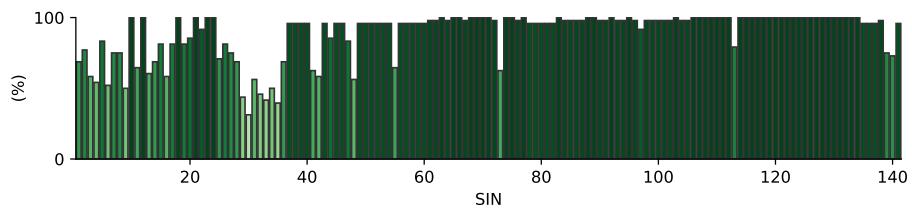


Figure 13: Picking percentage for each shot in the Danube island roll-along dataset. Low values in the first third of the dataset indicate a low signal-to-noise ratio in the seismograms hindering the picking of first break traveltimes for each shot-receiver pair.

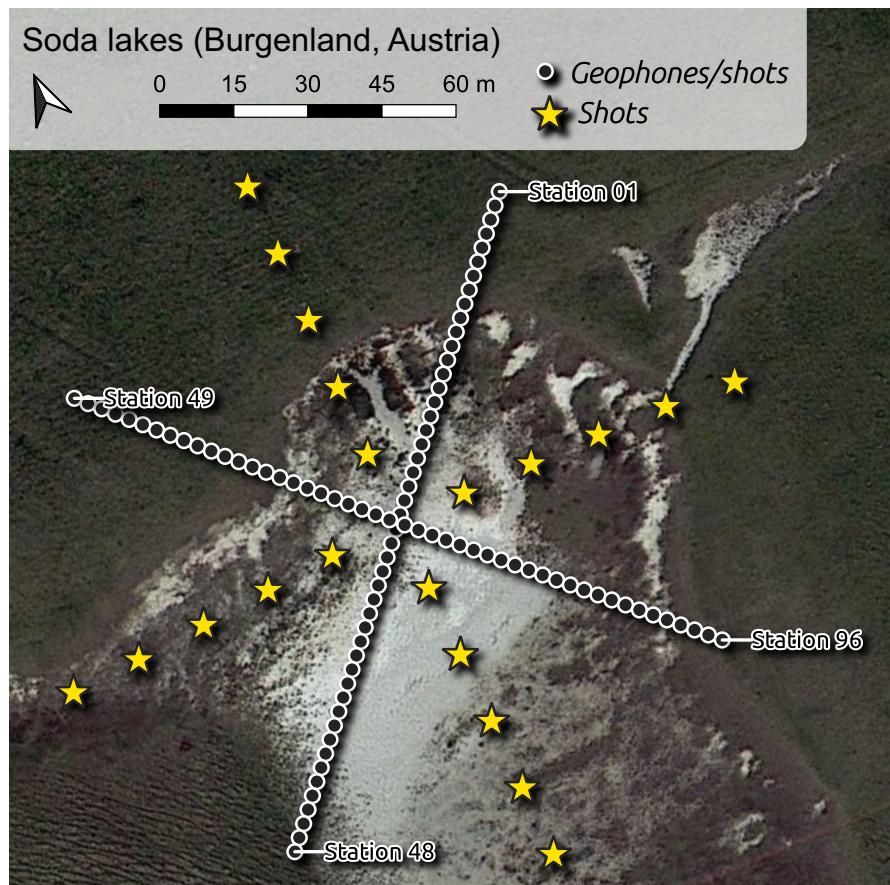


Figure 14: The soda lakes field dataset was collected in a 3D survey layout with stations (co-located geophones and shots indicated by filled dots) deployed in form of a cross. Additional shots (yellow stars) were conducted in front of a cross rotated by 45° to increase the coverage of the dataset.

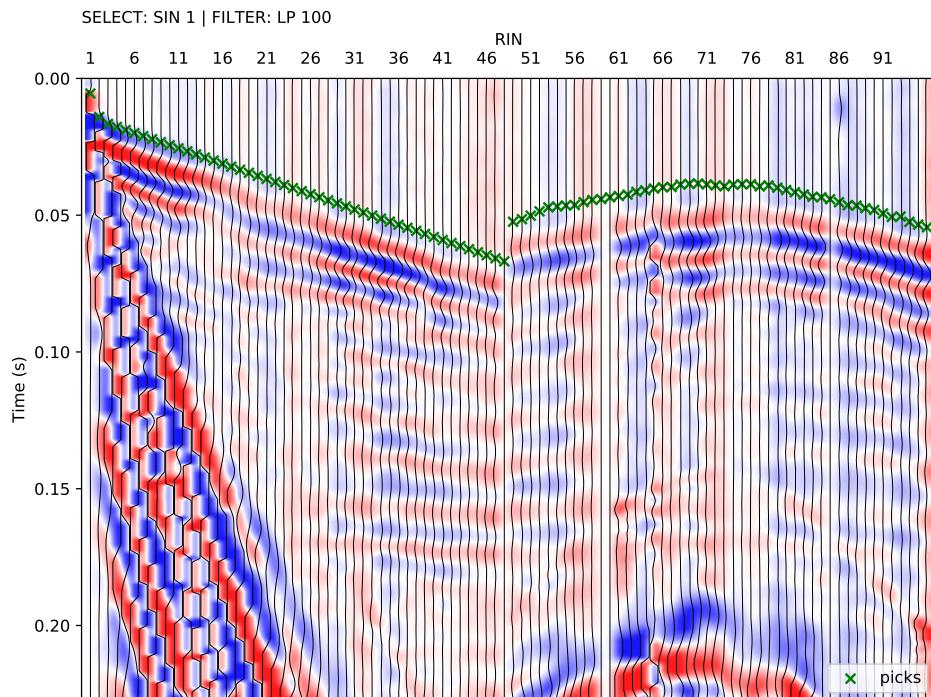


Figure 15: Exemplary seismic waveforms from the soda lakes dataset shown for shot index number (SIN) 1 with a 100 Hz lowpass filter applied to suppress high frequency noise. The recorded seismic waveforms clearly reflect the geometry of the geophones with RIN 1 to 48 deployed along the direction of wave propagation, while RIN 49 to 96 are deployed perpendicular to the propagating wavefront.

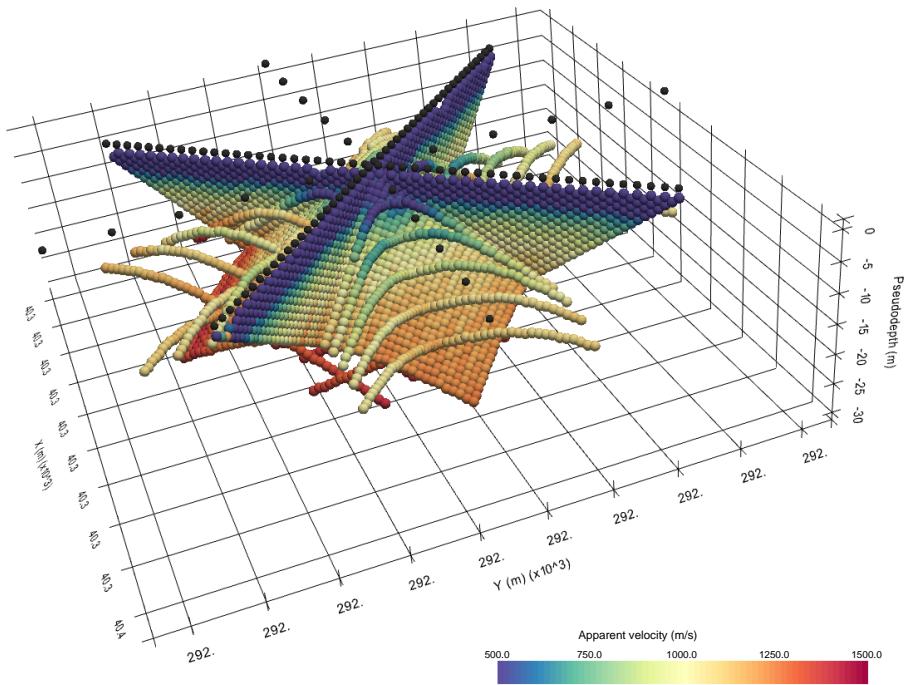


Figure 16: 3D pseudosection showing the apparent seismic velocities determined from the first break traveltimes obtained from the soda lakes dataset and the corresponding absolute offset between the shot and receiver stations. The apparent velocity for each shot-receiver pair is illustrated at the corresponding 2D midpoint and pseudodepth (1/3 of the absolute offset), thus yielding a 3D representation.