

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 March 8, 2023

²<https://www.landmark.solutions/SeisSpace-ProMAX>, last accessed on March 8, 2023

³<https://www.parkseismic.com/parkseis/>, last accessed on March 8, 2023

⁴<https://www.sandmeier-geo.de/reflexw.html>, last accessed on March 8, 2023

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
 102 particular, we present a carefully designed synthetic study highlighting the capabilities of the formikoj library for the
 103 modeling and processing of 2D seismic refraction datasets. Based on a 3D field dataset we illustrate that the formikoj

⁵http://resistivity.net/bert/data_format.html, last accessed on March 8, 2023

104 library also allows for the processing of complex survey geometries, where the obtained first break traveltimes can be
 105 inverted with third party open-source libraries to solve for 3D subsurface models expressed in terms of the seismic
 106 P-wave velocity.

107 2. Design and structure of the formikoj library

108 As illustrated in Figure 1, the formikoj library comprises a modeling and a processing module, which both rely on
 109 a common utilities module. The `DataModeler` and the `MethodManager` class provide the basis to add modeling or
 110 processing functionalities for any kind of geophysical methods.

111 The formikoj library provides a well thought out data management concept based on the SQLite database engine
 112 that does not require⁶. In particular, the information for each (processing) project, such as survey geometry and
 113 first break travel times, are stored in a SQLite database file. Using this application file format facilitates the cross-
 114 platform design of the formikoj library, and provides fast I/O operations through concise SQL queries. The portability
 115 of the application file allows for an easy exchange of projects between partners across institutions with different IT
 116 infrastructures. Accordingly, the formikoj library aims at providing a transparent and customizable framework for the
 117 collaborative design of reproducible workflows.

118 In this manuscript, we present the application of the formikoj library for the modeling and processing of seismic re-
 119 fraction data based on the fundamental use cases presented in Figure 2 and 3, respectively. These flow charts illustrate
 120 the corresponding workflows as well as the required interactions between the user, the formikoj library and third-party
 121 packages. As can be seen from Figures 2 and 3, the formikoj library acts as an interface between the user and more
 122 complex functionalities of third-party libraries, such as pyGIMLi for modeling and inversion of seismic refraction data.
 123 The `SeismicWaveformModeler` and `SeismicRefractionManager` classes implement these fundamental use cases,
 124 yet their actual capabilities are continuously expanded, e.g., to address specific modeling or processing requirements
 125 as well as to enhance the user experience. Similar to RefraPy, these classes are built upon the functionalities of exist-
 126 ing packages such as ObsPy for the processing of seismological data (Beyreuther et al., 2010) and pyGIMLi for the
 127 modeling and inversion of different geophysical data (Rücker et al., 2017). Other important third party dependencies
 128 refer to NumPy (Harris et al., 2020) and Pandas (McKinney, 2010) for general data handling, as well as matplotlib
 129 (Hunter, 2007) and PyVista (Sullivan and Kaszynski, 2019) for data visualization. In the current version, we imple-
 130 mented and tested formikoj primarily on Linux machines, yet the library has been successfully tested and used on all
 131 major operating systems, i.e., Linux, MacOS and Windows.

132 The comprehensive usage of clear error and log messages aims at supporting the user throughout the modeling and
 133 processing workflows. A meticulous exception handling ensures that the data stored in the SQLite project database is

⁶<https://www.sqlite.org/index.html>, last accessed on March 8, 2023

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

134 not corrupted in case of erroneous input. Moreover, documenting the user input and the respective responses of the
 135 formikoj library with the python logging module provides a timestamped command history that further enhances the
 136 transparency and repeatability of the conducted workflows.

137 2.1. Generation of seismic waveform data for synthetic subsurface models

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

```
144 1 # Import the SeismicWaveformModeler from the formikoj library
145 2 from formikoj import SeismicWaveformModeler
146 3
147 4 # Create an instance of the SeismicWaveformModeler
148 5 swm = SeismicWaveformModeler('..')
```

149 INFO : Created instance of SeismicWaveformModeler

150 In the working directory, the required input files are provided via the subdirectory *in*, whereas the modeling out-
 151 put will be stored in the automatically created subdirectory *out*. The key input file is the measurement scheme as it
 152 contains information regarding the distribution of the shot and geophone stations. If provided in the unified data for-
 153 mat, the measurement scheme is imported directly with pyGIMLI into a DataContainer. In case the measurement
 154 scheme is provided as a csv file, the SeismicWaveformModeler reads the information and writes it to a pyGIMLI
 155 DataContainer object. If provided as csv file, the measurement scheme contains a single line for each station in the
 156 survey layout, where a station either hosts a geophone or a shot, or both (see Table 1). The values provided in each
 157 line need to be separated by a unique delimiter, and the file must not contain a header.

158 For the modeling of the seismic waveform data, the parameters characterizing the base wavelet, the synthetic

159 subsurface model and the resulting waveform datasets are provided (see Table 2) in a configuration file following the
 160 yaml format:

```

 161   wavelet:
 162     length: 1.024
 163     frequency: 100
 164     sampling_rate: 2000
 165     pretrigger: 0.02
 166
 167   model:
 168     velmap: [[1, 750], [2, 2500], [3, 4000]]
 169     layers: [[1, 3], [2, 5], [3, 15]]
 170     quality: 32
 171     area: 10
 172     smooth: [1, 10]
 173     sec_nodes: 3
 174
 175   dataset:
 176     number: 1
 177     names: [syn_data]
 178     noise: 1
 179     noise_level: 1e-4
 180     missing_shots: 1
 181     broken_geophones: 1
 182     wrong_polarity: 1
 183
 184   traveltimes:
 185     noise_relative: 0.
 186     noise_absolute: 0.
```

In this exemplary configuration, the first block (`wavelet`) contains the parameterization of the base wavelet, which controls the modeling of the seismic waveforms (see Table 2). The second block (`model`) contains information regarding the synthetic subsurface model. Here, the user can define simple models with all layers considered to be parallel to the surface topography, which is inferred from the station geometry provided in the measurement scheme. The seismic velocity values for the different model regions (`velmap`) and the corresponding thickness of the different layers (`layers`) have to be explicitly defined in the configuration file. The remaining parameters (`quality`, `area`, `smooth` and `sec_nodes`) refer to the properties of the mesh that is eventually used for the forward modeling of the seismic waveform data and corresponding first break travel times (we refer to the respective pyGIMLi resources⁷ for further

⁷https://www.pygimli.org/pygimliapi/_generated/pygimli.mesh.html#pygimli.mesh.createMesh, last accessed March 8, 2023

information). In the third block of the exemplary configuration file (`dataset`), the user can set specific names for the datasets to be created (the number of datasets is automatically determined), or set the number of datasets to be created and the dataset names are automatically generated with the prefix `dataset_`. The remaining parameters provided in the `dataset` block control the random error (`noise` and `noise_level`) as well as systematic errors (`missing_shots`, `broken_geophones` and `wrong_polarity`) in the modeled seismic waveform data (see Table 2 for a detailed description). The number and position of the shot and geophone stations affected by the systematic errors are randomly chosen with a maximum of 5 % of the total number of stations in order to avoid a high number of invalid trace data. The parameters in the final block of the configuration file (`traveltimes`) the user defines values for the absolute error (e_{abs}) and the relative error (e_{rel}) considered for the forward modeling of the corresponding first break traveltimes (see Table 2). In particular, the error is estimated as

$$e = e_{abs} + d e_{rel}, \quad (1)$$

where d denotes the forward modeled data not affected by noise, i.e., here the first break travel times obtained through the Dijkstra algorithm (Dijkstra, 1959). These computed error values are subsequently added to the data to obtain the noisy data as

$$d_{noise} = d (1 + e). \quad (2)$$

187 The configuration file needs to be provided via the input directory from where it can be imported through the `load`
 188 method of the `SeismicWaveformModeler`:

```
189 6 # Load and parse the configuration file
190 7 swm.load(type='config')

191 INFO    : Configuration loaded
```

192 To demonstrate the ability to model seismic waveform data for arbitrary subsurface conditions we do not define a
 193 simple subsurface model in the configuration file but provide a more complex model in the binary mesh format (e.g.,
 194 a `bms` file). We prepare the model and the corresponding forward modeling mesh based on the mesh tools provided
 195 by `pyGIMLi` and save the mesh in the `bms` format (a commented version of the corresponding python script can be
 196 found in the Appendix). Similar to the configuration file, a `bms` file stored in the input directory can be imported into
 197 the workflow through the `load` method as follows:

```
198 8 # Load the mesh into the workflow
199 9 swm.load(type='mesh')
```

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

200 INFO : Mesh loaded

201 The forward modeling process generating the seismic waveform data is then invoked through the `create` method:

```
20210 # Start the modeling of the seismic waveform data
20311 swm.create(type='waveforms')
```

204 INFO : Measurement scheme loaded

205 INFO : Velocity model created

206 INFO : Wavelet created

207 [+++++ 100% +++++] 2048 of 2048 complete

```

208      ...
209  [+++++ 100% +++++] 2048 of 2048 complete
210 INFO    : Dataset 'syn_data' created

```

As can be seen from the log messages, the `SeismicWaveformModeler` first loads the measurement scheme into the workflow. In the second step, the provided mesh and the seismic velocity values defined in the configuration file are combined to create the seismic velocity model considered for the waveform modeling in this study (see Figure 4a). The model consists of a top and a bottom layer with varying thickness characterized by seismic velocity values of 750 ms^{-1} and 4000 ms^{-1} , respectively. In the center, the model contains an irregularly shaped anomaly associated with a seismic velocity of 2500 ms^{-1} , i.e., the model features vertical and lateral variations in the seismic velocity distribution. The third step refers to the generation of a Ricker wavelet through the `pyGIMLi` function `ricker` as

$$u = (1 - 2\pi^2 f^2 t^2) e^{-\pi^2 f^2 (t-t_0)^2} \quad (3)$$

```

211 based on the wavelet properties provided in the configuration file. In Equation 3,  $f$  is the frequency of the wavelet
212 given in Hz,  $t$  refers to the time base definition, i.e., the length and resolution of the wavelet in the time domain, and  $t_0$ 
213 is the time offset of the wavelet.

```

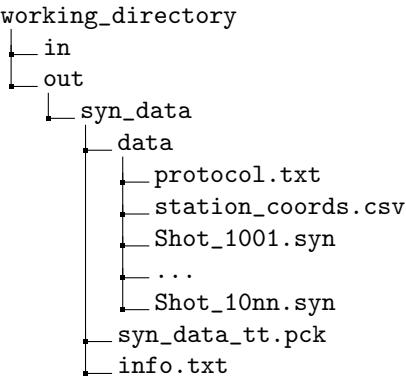
Subsequently, mesh, velocity model and Ricker wavelet are used to solve the pressure wave equation for each shot station defined in the measurement scheme with the `pyGIMLi` function `solvePressureWave`:

$$\frac{\partial^2 u}{\partial t^2} = \nabla \cdot (D \nabla c) \quad (4)$$

```

214 The resultant seismograms are extracted at the corresponding geophone stations, gathered for each shot position in an
215 obspy Stream object and saved in the output directory (out) as shown here:

```



```

216 In particular, the subdirectory data contains a separate file in the miniseed format (Ahern et al., 2012; Ringler and
217 Evans, 2015) for each shot position, where the file extension syn indicates that the shot files contain forward modeled
218 seismic waveform data. The measurement protocol (protocol.txt) and the station coordinates provided as a csv file
219 (station_coords.csv) are also stored in this directory. The header of the measurement protocol contains the survey

```

220 parameters, e.g., sampling rate, recording length, number of geophones and geophone spacing. Moreover, the protocol
 221 associates each shot file of the dataset with a specific location in the survey geometry with respect to the geophone
 222 positions, e.g.:

```
223 #####  

224 Line: SYN_syn_data  

225 Sampling rate: 2000 Hz  

226 Recording length: 0.512 s  

227 Number of geophones: 48  

228 Geophone spacing: 2 m  

229 #####  

230  

231     File number | Station  

232         1001   |   G001  

233         :      |      :  

234         1048   |   G048
```

235 The auxiliary file info.txt exported to the dataset directory summarizes the parameters from the configuration file as
 236 well as information regarding the simulated systematic errors in the synthetic seismic waveform data, e.g.:

```
237 Number of geophones: 48  

238 Number of shots: 48  

239 Recording length (s): 0.512  

240 Sampling frequency (Hz): 2000  

241 Wavelet type: Ricker  

242 Frequency of the wavelet (Hz): 100  

243  

244 Missing shot(s): 11  

245 Broken geophone(s): 13  

246 Wrong polarity geophone(s): 26
```

247 Figure 4b presents the seismic waveform data forward modeled for a shot point located in the center of the profile.
 248 The seismograms are shown as curves, whereas the strength of the amplitudes is added as color-coded information.
 249 In the seismograms we see clear first onsets along the entire profile, yet receivers 3 and 45 were modeled as broken
 250 geophones, i.e., the corresponding seismograms contain solely noise. Crosses overlaid on the valid seismograms at
 251 the respective first onset indicate the first break traveltimes tt between the shot and the receiver. In addition to the
 252 missing first break traveltimes for the broken receivers, we manually added a systematic error, i.e., we set an erroneous
 253 traveltime for receiver 36, to demonstrate how such outliers can be identified in a so-called pseudosection.

A pseudosection as presented in Figure 4c, illustrates apparent velocity (v_{app}) values computed as

$$v_{app} = \frac{aoff\ set}{tt}, \quad (5)$$

where $aoffset$ refers to the absolute offset between shot and receiver. The v_{app} values are plotted at pseudolocations, where the location along profile direction is defined by the midpoint of the corresponding shot-receiver pair and the pseudodepth is computed as $1/3$ of $aoffset$. As demonstrated in Figure 4c, a pseudosection allows for the identification of missing data (e.g., receiver 3 and 45) as well as systematic errors or outliers, i.e., velocities erroneously influenced by a single shot or receiver (see receiver 36). The main assumption here is that the pseudosection should reveal smooth transitions between lateral and vertical neighbors, considering that the data were collected with gradual changes in the position of the source (i.e., hammer blow) and the receiver (i.e., geophone). The pseudosection will reveal large variations in case of abrupt changes in the topography or the geometry of the array, yet this can be taken into account by the user during the identification of outliers and possible systematic errors.

Based on pseudosections erroneous traveltimes can be identified and subsequently removed or corrected, which is a critical processing step prior to the inversion of the data. In particular, the inversion of first break traveltimes aims at resolving a model of the P-wave velocity of the subsurface materials. Since the inversion of geophysical data is not within the scope of the formikoj library we rely for this purpose on the modeling and inversion capabilities of pyGIMLi (Rücker et al., 2017). The default inversion framework of pyGIMLi uses a generalized Gauss-Newton method to solve the inversion problem through the minimization of an objective function given as:

$$\Psi(\mathbf{m}) = \Psi_d(\mathbf{m}) + \lambda \Psi_m(\mathbf{m}) . \quad (6)$$

The first term on the right-hand side of Equation 6 refers to the data misfit, whereas the second term denotes represents the model constraints, i.e., the regularization. The regularization parameter λ controls the influence of the regularization term on the inversion process.

The inversion of the synthetic first break traveltimes considered here resolves the subsurface model presented in Figure 4d. To aid in the evaluation of this inversion result we superimposed the known interfaces between the different subsurface unit of the synthetic model. As can be seen from this plot, the imaging result resolves the fundamental structural features and reflects the P-wave velocity distribution of the synthetic subsurface model. Deviations from the true velocity model are due to the smoothness-constraint inversion scheme applied by pyGIMLi. To obtain sharper contrasts between the different subsurface units in the imaging result structural constraints could be incorporated in the inversion, e.g., as demonstrated by (Steiner et al., 2021). Nonetheless, we consider Figure 4 to demonstrate the applicability of the `SeismicWaveformModeler` class for generating synthetic seismic waveform data to be used in numerical P-wave refraction seismic investigations.

275 2.2. Processing of seismic refraction datasets

276 The SR method is based on the measurement of the traveltimes of seismic waves determined from the the first onset
 277 of the seismic energy recorded by geophones. In the SRT, the inversion of traveltimes gathered from tens to hundreds
 278 of seismograms permits the computation of variations in the seismic velocities in an imaging framework. Measuring
 279 the traveltimes in seismograms (i.e., first break picking) is commonly done manually (or semi-automatically) in an
 280 iterative process. The signal-to-noise ratio in the recorded seismograms substantially influences the quality of the
 281 traveltimes picked in the seismograms, and thus a proper enhancement of the perceptibility of the first onsets is crucial.
 282 Accordingly, the `SeismicRefractionManager` class provides functionalities that permit the processing of seismic
 283 waveform data for a refraction seismic analysis. These functionalities involve the reading of seismic waveform data,
 284 combining the data with information about the survey geometry, processing of the waveforms as well as the picking of
 285 first break traveltimes. Since the first break picking is a highly interactive process, the `SeismicRefractionManager`
 286 is designed primarily for usage from within an ipython shell.

287 For the demonstration of the fundamental `SeismicRefractionManager` capabilities we consider the synthetic
 288 seismic data created with the `SeismicWaveformModeler` above. This allows us to show that both classes can be
 289 combined in subsequent workflows, and to present the forward modeled seismic waveform data, the synthetic first
 290 break traveltimes as well as the seismic P-wave velocity model obtained through the inversion of these traveltimes.

291 2.2.1. Compiling the survey information and creating a project

292 To create a new or load an existing `SeismicRefractionManager` project, the working directory needs to contain
 293 specific subdirectories:

```
working_directory
  +-- 01_data
    +-- raw
  +-- 02_geom
  +-- 03_proc
```

294 In this directory structure, the seismic shot files are stored in `01_data/raw` and the geometry file (`geometry.csv`) is
 295 provided in `02_geom`. The geometry file is a csv file that provides an abstract representation of the survey layout and
 296 must not contain a header. The fundamental element for the description of the survey layout is the station, which refers
 297 either to a geophone position, a shot position or a position with co-located shot and geophone. For each station the
 298 geometric and semantic information described in Table 3 are provided column-wise, i.e., each line in the geometry file
 299 corresponds to a single station with a unique position within the survey layout. For the synthetic dataset considered in
 300 this study, the 2D station coordinates are provided in the geometry file, whereas the z coordinate is assumed to be 0 m
 301 along the entire profile, as illustrated in Table 4; thus, reflecting the lack of surface topography in the synthetic model
 302 (see Figure 4a). In the column **Geo** the geometry file contains the value 1 (True) for each station except the station

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)

Table 4

Excerpt from the geometry file describing the survey geometry of the synthetic dataset generated in this study.

x (m)	y (m)	z (m)	Geo	Shot	1st Geo	# Geo
0.0	0.0	0.0	1	1001	1	48
2.0	0.0	0.0	1	1002	1	48
:	:	:	:	:	:	:
20.0	0.0	0.0	1	-1	1	48
:	:	:	:	:	:	:
24.0	0.0	0.0	0	1013	1	48
:	:	:	:	:	:	:
92.0	0.0	0.0	1	1047	1	48
94.0	0.0	0.0	1	1048	1	48

located at 24 m referring to the broken geophone modeled by the `SeismicWaveformModeler`. When compiling the geometry file we also need to take into account the missing shot at station 11, i.e., the station located at 20 m along profile direction, and set the corresponding value to -1 . The column **1st Geo** indicates the first active geophone along the profile for each shot file, which for the synthetic dataset considered here is the geophone deployed at the first station along the entire profile. In particular, the column **1st Geo** allows the reproduction of roll-along survey geometries, where in the first segment the first active geophone is always the geophone at station 1. Considering 48 geophones deployed in each segment, and an overlap of 50 %, the first geophone for the consecutive segments would be 25, 49, 73, 97, and so on.

An instance of the `SeismicRefractionManager` can be created by providing the path to the working directory as parameter to the constructor. Depending on the content of the working directory, the `SeismicRefractionManager` automatically decides whether (i) to start in the data preview mode (first break picking not possible), (ii) create a new project, or (iii) load an existing project from disk. In case the shot files as well as the geometry file are provided and a basic sanity check of the geometry file was successful, the `SeismicRefractionManager` creates a new project:

```
316 1 # Import the SeismicRefractionManager from the formikoj library
317 2 from formikoj import SeismicRefractionManager
318 3
```

```

319 4 # Create an instance of the SeismicRefractionManager
320 5 srm = SeismicRefractionManager('.')

321 INFO    : Read geometry information from file
322 INFO    : Extracted shot geometry
323 INFO    : Extracted receiver geometry
324 INFO    : Applied geometry
325 INFO    : Standard pickset 'picks' created
326 INFO    : Pickset 'picks' loaded
327 INFO    : 'picks' set as active pickset
328 Progress <===== 100.0% completed
329 INFO    : Read 48 files

```

330 In a first step, the `SeismicRefractionManager` creates an SQLite database `prj.db` in the working directory based on
331 the entity-relationship diagram shown in Figure 5. The geometry information is then read from the geometry file and
332 stored in the database table `geometry` with consecutively numbered stations (see Figure 6). To allow for an efficient data
333 selection for the user the `SeismicRefractionManager` creates database tables `shots` and `receivers`, which link the
334 station numbers to shot index numbers (SIN) and receiver index numbers (RIN), respectively. For each shot-receiver
335 pair the corresponding SIN and RIN are stored in the table `applied_geometry` together with the absolute offset and
336 midpoint between these stations, i.e., the geometry is applied. In the last step, the database table `fbpicks` is created,
337 which stores the first break traveltimes for each SIN-RIN pair together with the name of the corresponding pickset, i.e.,
338 a common label for an entire set of first break traveltimes. By default, each project contains the default pickset 'picks',
339 which is loaded and activated on startup. Once the database is initialized, the waveform data are read from disk and
340 the project is ready for processing.

341 2.2.2. Selecting and visualizing seismic waveform data

342 Once the geometry is applied the `select` method of the `SeismicRefractionManager` allows to gather the seis-
343 mic waveform data based on a common absolute offset (`aoffset`), a common RIN (`rin`), or a common SIN (`sin`)

```

344 6 # Select traces with common absolute offset
345 7 srm.select(by='aoffset', num=6)

```

```

346 INFO    : 88 traces selected

```

```

347 8 # Select traces with receiver
348 9 srm.select(by='rin', num=10)

```

```

349 INFO    : 48 traces selected

```

```
350:0 # Select traces with receiver
351:1 srm.select(by='sin', num=23)
```

352 INFO : 48 traces selected

353 For the currently selected trace data the frequency spectrum can be computed and visualized through the plot
354 method:

```
355:2 # Plot the frequency spectrum
356:3 srm.plot(type='spectrum')
```

357 Examining a frequency spectrum as shown in Figure 7 provides information regarding the dominating signal frequen-
358 cies. Discriminating signal and noise frequency ranges is critical for the selection of adequate filter types and settings
359 in order to enhance the signal-to-noise ratio of the seismograms. Filters can be applied to the currently selected traces
360 through the filter method, which utilizes the frequency filters implemented in the ObsPy package (lowpass, high-
361 pass, bandpass and bandstop; Beyreuther et al., 2010), e.g.:

```
362:4 # Apply bandpass filter
363:5 srm.filter(type='bandpass', freqmin=10, freqmax=120)
```

364 INFO : Applied bandpass filter (10.0 to 120.0 Hz)

365 By default, filters are solely applied to the currently selected traces, yet setting the filter to hold on automatically filters
366 all subsequently selected traces with the same settings:

```
367:6 # Apply lowpass filter
368:7 srm.filter(type='lowpass', freq=200, onhold=True)
```

369 INFO : Applied 200.0 Hz lowpass filter

370 INFO : Set filter hold on

371 The currently selected traces can be visualized by calling the plot method without passing any parameter:

```
372:8 # Plot selected traces
373:9 srm.plot()
```

374 By default, the seismogram plot visualizes the seismic waveforms in a combination of wiggle trace and variable area
375 mode, i.e., the trace data are shown as curves. The area of the curves is colored red for negative and blue for positive
376 amplitudes, respectively (not shown for brevity). Pressing the up or down arrow key on the keyboard toggles between
377 variable area and variable density mode, with the latter reflecting the strength of the amplitudes by the color saturation,
378 i.e., high amplitudes refer to a stronger shade than low amplitudes (see Figure 8).

379 The active processing mode and data scaling mode are reported together with the traveltimes at the current cursor
 380 position in the status bar of the seismogram plot window (see Figure 9). The initial processing mode is 'Fb pick',
 381 i.e., first break picking is possible. The user can switch between the different modes by pressing specific keys on the
 382 keyboard. The 'm' key activates the trace mute mode ('Trc mute'), which allows to set the amplitude of a trace to
 383 zero by clicking with the left mouse button; clicking again on the same trace restores the amplitude information. The
 384 trace reverse mode ('Trc rev') is activated by pressing the 'r' key and enables the user to toggle the polarity of a trace
 385 by clicking on it with the left mouse button. The default data scaling mode is 'Zoom', which allows the scaling of
 386 the y-axis by turning the mouse wheel. By pressing the 'a' key the amplitude scaling mode ('Amp scal') is activated.
 387 Turning the mouse wheel increases or decreases the amplitudes of the traces currently shown in the seismogram plot,
 388 and thus might help to enhance the perceptibility of the first onsets. By pressing the key of the currently active mode
 389 again, the SeismicRefractionManager returns to the default mode; yet, the different modes can be activated in any
 390 arbitrary order (as illustrated in Figure 9).

391 2.2.3. Analysis of the seismic waveforms and first break traveltimes picking

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

397 If the processing mode 'Fb pick' is activated, picking of first break traveltimes is done individually by clicking
 398 with the left mouse button on the respective trace. Clicking again on the same trace will set the first break pick to the
 399 new location as there can only be one traveltimes for each SIN-RIN pair; whereas clicking with the right mouse button
 400 deletes the pick. Alternatively, first break picks can be set for multiple traces by pressing the left mouse button and
 401 moving the cursor. Once the left mouse button is released, first break picks are defined at the intersections between
 402 the line and the seismograms. In the same way, multiple picks can be deleted if a line is drawn with the right mouse
 403 button pressed. The first break traveltimes determined in the seismogram plot are automatically written to the project
 404 database when the window is closed or another set of traces is loaded by pressing the 'left' or 'right' arrow keys on the
 405 keyboard.

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

```
40720 # Plot traveltimes diagram
40821 srm.plot(type='traveltimes')
```

409 Figure 10 presents an exemplary traveltimes diagram, which is a common way to examine the quality of the first break

410 picking. Such illustration of the traveltimes can be used to identify outliers or erroneous measurements, which are
 411 commonly associated to traveltimes with substantial deviations from those observed at adjacent stations such as the
 412 first break pick for the SIN-RIN pair (23, 11). Outliers can be removed by clicking on the respective data point in
 413 the travelttime diagram, which is instantly synchronized with the project database. If the seismogram plot and the
 414 travelttime diagram are used side-by-side, changes made to the first break picks in one window will interactively trigger
 415 an update of the other one and vice versa. Once the user is satisfied with the quality of the first break traveltimes, the
 416 data points of the pickset can be exported in the unified data format, which is compatible with the pyGIMLi framework.
 417 In particular, the udf file is stored in the subdirectory *03_proc/picks* with the current timestamp as suffix:

```
41822 # Export first break traveltimes
41923 srm.picksets(do='export', name='pick')

420 INFO    : Pickset 'pick' saved to pick_20230303T140447.pck'
```

421 2.3. Expanding the capabilities of the formikoj library

422 Making the formikoj library publicly available under an open-source license allows the addition of supplementary
 423 functionalities tailored for the specific requirements of the users. The concept of formikoj suggest that such custom
 424 extensions should be implemented either as internal methods or as functions in the utilities module, which are then
 425 executed through the `compute` method.

We illustrate this possibility for customization by implementing a simplified version of an automatic first break picking algorithm based on the short- and long-time window average ratio (STA/LTA) method (Allen, 1978), which determines the traveltimes following the energy ratio approach (e.g., Earle and Shearer, 1994). In particular, our implementation computes the envelope $E(t)$ of the seismogram $s(t)$ as (e.g., Duan and Zhang, 2020)

$$E(t) = (s(t)^2 + \tilde{s}(t))^{1/2}, \quad (7)$$

where $\tilde{s}(t)$ is the Hilbert transform of $s(t)$. The energy ratio ER is then computed as $ER = STA/LTA$ with

$$STA(i) = \frac{1}{n_{STA}} \sum_{j=i-n_{STA}}^i E(j), \quad (8)$$

and

$$LTA(i) = \frac{1}{n_{LTA}} \sum_{j=i-n_{LTA}}^i E(j), \quad (9)$$

426 where n_{STA} and n_{LTA} refer to the number of data points in $E(t)$ considered for the short- and long-time windows,

427 respectively. For this exemplary implementation we determine the first break traveltimes in the seismograms as the
 428 position of the maximum in the *ER* function, which is automatically saved in the project database.

429 The autopicking algorithm is added to the `SeismicRefractionManager` in form of two internal methods `_manage_autopicks`
 430 and `_compute_autopicks`, respectively. To invoke the autopicking process we modified the `compute` method, which
 431 now accepts the custom-defined keyword `autopicking`. Additionally, the autopicking requires values to be passed
 432 for the parameters `pick` and `pickset`. The former accepts the values `all` or `cur` to determine first break traveltimes
 433 for all traces in the dataset or solely the currently selected traces, respectively. The latter defines the name of the pick-
 434 sets in which the automatically determined traveltimes should be saved to. A typical use case involves conducting the
 435 autopicking and exporting the determined traveltimes:

```
43624 # Automatically pick first break traveltimes
43725 srm.compute(do='autopicking', pick='all', pickset='autopicks')
43826 srm.picksets(do='export', name='autopicks')
```

```
439 INFO    : Created new pickset 'autopicks'
440 INFO    : Pickset 'autopicks' loaded
441 INFO    : 'autopicks' set as active pickset
442 Progress <===== 100.0% completed
443 INFO    : Pickset 'autopicks' saved to autopicks_20230303T140834.pck
```

444 In Figure 11, we compare the automatically determined first break picks with the forward modeled traveltimes
 445 computed above to allow for a basic evaluation of autopicking quality. The inset plot in Figure 11 presents the his-
 446 togram of the autopicked traveltimes, which shows that three traveltimes should be considered outliers, and thus are
 447 removed from the dataset. After the outlier removal the correlation coefficient between forward modeled and au-
 448 topicked traveltimes is 0.99 suggesting that the implemented autopicking algorithm performs well for the synthetic
 449 seismic waveform data. However, the observed deviation from the perfect correlation (i.e., the 1:1 line in Figure 11)
 450 indicates that autopicking and forward modeling algorithm might be sensitive to different seismic wave phases. Further
 451 improvements in terms of the autopicking process could incorporate, for example, machine learning approaches as the
 452 method proposed by Duan and Zhang (2020), which can be easily implemented as an additional functionality in the
 453 `SeismicRefractionManager`.

454 3. Application to field data: Processing a 3D seismic refraction dataset

455 To demonstrate the applicability of the `SeismicRefractionManager` for the processing of seismic refraction data
 456 collected in the field we consider a 3D seismic survey conducted in a soda lake located in the Nationalpark Neusiedler
 457 See-Seewinkel close to Vienna. The soda lake corresponds to quaternary sediments where capillary forces have

Table 5

Excerpt from the 3D survey geometry file.

x (m)	y (m)	z (m)	Geo	Shot	1st Geo	# Geo
40336.056	292470.692	120.0	1	1001	1	96
40334.751	292469.177	120.0	1	1002	1	96
:	:	:	:	:	:	:
40284.472	292455.431	120.0	1	1058	1	96
40285.896	292454.027	120.0	1	1059	1	96
:	:	:	:	:	:	:
40339.421	292402.681	120.0	1	1096	1	96
40305.228	292445.050	120.0	0	1097	1	96
:	:	:	:	:	:	:
40265.415	292431.957	120.0	0	1115	1	96
40255.453	292431.395	120.0	0	1116	1	96

458 developed a low permeable layer close to the surface (between 50 and 100 cm) with a high clay and salt content. The
 459 seismic survey aims to support the interpretation of electrical and electromagnetic models obtained in a monitoring
 460 framework. Accordingly, the geometry of the seismic survey shown in Figure 12 was specified by previously conducted
 461 electrical measurements with electrodes arranged along two perpendicular lines. The seismic data were collected
 462 with 48 geophones deployed along the North-East to South-West oriented line, and 48 geophones along the North-
 463 West to South-East oriented line, with a spacing of 2 m between the geophones. Shots were generated with an 8 kg
 464 sledgehammer at the geophone positions as well as at positions along the diagonals to obtain a sufficient coverage, as
 465 indicated in Figure 12.

466 As in case of the synthetic data set presented above, the geometry file contains the coordinates of the survey
 467 stations, yet for the soda lake dataset 3D coordinates we provided as illustrated in Table 5. Since geophones were not
 468 deployed at each survey station the column **Geo** contains the value 1 (True) for the first 96 stations and 0 (False) for
 469 the remaining 20 stations. Based on the shot files and the geometry file stored in the required directory structure the
 470 SeismicRefractionManager creates the project database, automatically infers the 3D survey layout from the 3D
 471 survey geometry, and thus configures the project for 3D processing. Then we can select seismograms the same way as
 472 for the synthetic data presented above. For this example we select seismic waveform data recorded for SIN 1, i.e., the
 473 shot position co-located with the first geophone (Station 01 in Figure 12):

```
474 10 # Select traces with receiver
475 11 srm.select(by='sin', num=1)
```

476 INFO : 96 traces selected

477 In Figure 13, we can see that the data for RIN 1 to 48 appear familiar as the corresponding SIN-RIN layout refers
 478 to the one of conventional 2D profiles; whereas the seismic waveforms for RIN 49 to 96 show an entirely different
 479 pattern. To understand such visualization, we need to take into account that RIN 49 to 96 are deployed perpendicular

480 to the direction of propagation of the wavefront originating from SIN 1. Accordingly, the observed curvature in the
 481 first onsets is due to the varying offset of the different SIN-RIN pairs. Also for the 3D seismic refraction data the first
 482 break traveltimes can be determined either individually for each trace or by

483 Due to the 3D survey geometry, a 2D pseudosection is not suitable for assessing the quality of the first break trav-
 484 eltimes. However, the `SeismicRefractionManager` project is configured for 3D processing, and thus the apparent
 485 velocity values are illustrated in an interactive 3D pseudosection. This plot can be rotated and the image section can
 486 be zoomed and panned allowing the user to easily investigate the data quality for 3D geometries. Figure 14 shows
 487 a screenshot of the 3D pseudosection for the salt lake dataset; yet, such screenshot cannot reveal the full capabilities
 488 implemented in the `SeismicRefractionManger` for the interactive analysis and visualization of 3D pseudosections.

489 To review the data quality for the entire dataset it is possible to visualize the picking percentage, i.e., the ratio of
 490 actually picked traveltimes and total number of SIN-RIN pairs:

```
491 10 # Plot the picking percentage
492 11 srm.plot(type='pickperc')
```

493 Figure 15 presents the picking percentage visualized separately for each SIN. Accordingly, a low picking percentage
 494 indicates shots affected by a low signal-to-noise ratio, whereas clear first onsets yield a correspondingly high picking
 495 percentage. For the soda lake dataset we observe a consistently high picking percentage for all shot positions; thus,
 496 indicating a good data quality. Moreover, the picking percentage plot can be used to track the picking progress, for
 497 instance if the traveltimes cannot be determined in one session or to identify single shots that might have been forgotten
 498 during the first break picking. Accordingly, it is advisable to check the picking percentage prior to exporting the
 499 traveltimes for the inversion.

500 Once the first break picking is finished, the corresponding pickset can be exported for the inversion. We inverted the
 501 first break traveltimes with pyGIMLi and present the resolved 3D subsurface model in Figure 16. However, to facilitate
 502 the representation and interpretation of the seismic velocity distribution we decided to show selected slices through the
 503 3D subsurface model. From this representation we can see, that the inversion solves for low seismic velocities (600 to
 504 800 ms^{-1}) in the near-surface in the center of the investigated area, which corresponds to the still intact part of the soda
 505 lake, i.e., the part that is covered by water on a seasonal basis. Seismic velocity values above 800 ms^{-1} are resolved
 506 at depth as well as outside of the soda lake. A more detailed interpretation is beyond the scope of this manuscript, yet
 507 the presented figures reveal the capabilities provided by the proposed framework for the visualization and processing
 508 of data collected in 3D survey geometries.

509 4. Conclusions and Outlook

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

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

522 By making the source code of the formikoj library available under the MIT license we intend to spark the develop-
 523 ment of further modeling and processing tools for various geophysical models based on this framework. Our further
 524 efforts will focus on implementing tools for other wave-based geophysical methods used in frame of our research
 525 activities, such as multi-channel analysis of (seismic) surface waves or magnetotelluric surveys.

526 5. Acknowledgments

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

531 Code and data availability section**532** Name of the code/library: formikoj**533** Contact: matthias.steiner@geo.tuwien.ac.at**534** Hardware requirements: No specific requirements**535** Program language: Python**536** Software required: Anaconda/Miniconda recommended**537** Total program and dataset size: 250 MB**538** The source codes and exemplary data sets are available for downloading at the link:**539** <https://git.geo.tuwien.ac.at/msteine1/formikoj.git>**540** The orthophotos used in Figures ?? and 12 were published by geoland.at under a CC BY 4.0 license.**541 A. Source code for generating the subsurface model considered in this study**

```

542 1 # Import required packages
543 2 import numpy as np
544 3 import pygimli as pg
545 4 import pygimli.meshTools as mt
546 5
547 6 # Create the model geometry
548 7 # - top layer
549 8 top = mt.createPolygon([[0, 0], [94, 0],
550 9                 [94, -3.5], [72, -3.5],
551 10                [20, -2], [0, -2]],
552 11                isClosed=True, marker=1, area=0.1)
553 12
554 13 # - bottom layer
555 14 bottom = mt.createPolygon([[0, -2], [20, -2],
556 15                 [22, -6], [70, -6],
557 16                 [72, -3.5], [94, -3.5],
558 17                 [94, -10], [0, -10]],
559 18                isClosed=True, marker=3, area=0.1)
560 19
561 20 # - anomaly/infill
562 21 infill = mt.createPolygon([[20, -2], [72, -3.5],
563 22                 [70, -6], [22, -6]],
564 23                isClosed=True, marker=2, area=0.1)
565 24
566 25 geom = top + infill + bottom

```

```

56726
56827 # Define shot and receiver stations and create corresponding nodes
56928 nstats = 48
57029 spc = 2
57130
57231 stations = np.vstack((np.linspace(0, (nstats-1)*spc, nstats),
57332                         np.zeros(nstats))).T
57433
57534 for p in stations:
57635     geom.createNode(p)
57736
57837 # Create mesh for the finite element modeling
57938 mesh = mt.createMesh(geom, quality=34)
58039
58140 # Save the mesh in the binary mesh format for later use
58241 # with the SeismicWaveformModeler
58342 mesh.save('mesh.bms')

```

584 References

- 585 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.
- 586 Allen, R.V., 1978. Automatic earthquake recognition and timing from single traces. Bulletin of the seismological society of America 68, 1521–1532. doi:10.1785/bssa0680051521.
- 587 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.
- 588 Blanchy, G., Saneiyan, 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>.
- 589 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.
- 590 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>.
- 591 Dijkstra, E.W., 1959. A note on two problems in connexion with graphs. Numerische Mathematik , 269–271doi:10.1007/bf01386390.
- 592 Draebig, 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>.
- 593 Duan, X., Zhang, J., 2020. Multitrace first-break picking using an integrated seismic and machine learning methodpicking based on machine

- 604 learning. *Geophysics* 85, WA269–WA277. doi:10.1190/GEO2019-0422.1.
- 605 Earle, P.S., Shearer, P.M., 1994. Characterization of global seismograms using an automatic-picking algorithm. *Bulletin of the Seismological*
606 Society of America 84, 366–376.
- 607 Guedes, V.J.C.B., Maciel, S.T.R., Rocha, M.P., 2022. Refrappy: A python program for seismic refraction data analysis. *Computers & Geo-*
608 *sciences* 159, 105020. URL: <https://www.sciencedirect.com/science/article/pii/S0098300421003022>, doi:<https://doi.org/10.1016/j.cageo.2021.105020>.
- 610 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,
611 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,
612 K., Reddy, T., Weckesser, W., Abbasi, H., Gohlke, C., Oliphant, T.E., 2020. Array programming with NumPy. *Nature* 585, 357–362. URL:
613 <https://doi.org/10.1038/s41586-020-2649-2>, doi:10.1038/s41586-020-2649-2.
- 614 Heimann, S., Kriegerowski, M., Isken, M., Cesca, S., Daout, S., Grigoli, F., Juretzek, C., Megies, T., Nooshiri, N., Steinberg, A., Sudhaus, H.,
615 Vasyura-Bathke, H., Willey, T., Dahm, T., 2017. Pyrocko - an open-source seismology toolbox and library doi:10.5880/GFZ.2.1.2017.001.
- 616 Hunter, J.D., 2007. Matplotlib: A 2d graphics environment. *Computing in Science & Engineering* 9, 90–95. doi:10.1109/MCSE.2007.55.
- 617 McKinney, W., 2010. Data Structures for Statistical Computing in Python, in: Stéfan van der Walt, Jarrod Millman (Eds.), *Proceedings of the 9th*
618 *Python in Science Conference*, pp. 56 – 61. doi:10.25080/Majora-92bf1922-00a.
- 619 Nguyen, F., Ghose, R., Isunza Manrique, I., Robert, T., Dumont, G., 2018. Managing past landfills for future site development: A review of the
620 contribution of geophysical methods, in: *Proceedings of the 4th International Symposium on Enhanced Landfill Mining*, pp. 27–36.
- 621 Parsekian, A., Singha, K., Minsley, B.J., Holbrook, W.S., Slater, L., 2015. Multiscale geophysical imaging of the critical zone. *Reviews of*
622 *Geophysics* 53, 1–26. doi:10.1002/2014RG000465.
- 623 Plattner, A.M., 2020. Gprpy: Open-source ground-penetrating radar processing and visualization software. *The Leading Edge* 39, 332–337.
624 doi:10.1190/tle39050332.1.
- 625 Ringler, A.T., Evans, J.R., 2015. A quick seed tutorial. *Seismological Research Letters* 86, 1717–1725. doi:10.1785/0220150043.
- 626 Romero-Ruiz, A., Linde, N., Keller, T., Or, D., 2018. A review of geophysical methods for soil structure characterization. *Reviews of Geophysics*
627 56, 672–697. doi:10.1029/2018RG000611.
- 628 Rücker, C., Günther, T., Wagner, F.M., 2017. pygimli: An open-source library for modelling and inversion in geophysics. *Computers & Geosciences*
629 109, 106–123. doi:10.1016/j.cageo.2017.07.011.
- 630 Steiner, M., Katona, T., Fellner, J., Flores Orozco, A., 2022. Quantitative water content estimation in landfills through joint inversion of seismic re-
631 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>.
- 632 Steiner, M., Wagner, F.M., Maierhofer, T., Schöner, W., Flores Orozco, A., 2021. Improved estimation of ice and water contents in alpine permafrost
633 through constrained petrophysical joint inversion: The hoher sonnblick case study. *Geophysics* 86, 1–84. doi:10.1190/geo2020-0592.1.
- 634 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.
- 635 Sullivan, C.B., Kaszynski, A., 2019. PyVista: 3d plotting and mesh analysis through a streamlined interface for the visualization toolkit (VTK).
636 *Journal of Open Source Software* 4, 1450. URL: <https://doi.org/10.21105/joss.01450>, doi:10.21105/joss.01450.
- 637 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*,
638 pp. 96–103.

641 List of Figures

642 1	Fundamental use case illustrating the modeling of seismic refraction data within the formikoj ecosystem.	25
643 2	Fundamental use case illustrating the processing of seismic refraction data within the formikoj ecosystem.	26
644 3	Typical formikoj workflow for the picking of first break traveltimes from seismic waveform data.	27
645 4	Synthetic seismic data generated with the SeismicWaveformModeler:	
646 (a)	Two-layered subsurface model without topography characterized by a distinct anomaly in the center. Triangles along the surface indicate the positions of geophones.	
647 (b)	Synthetic seismic waveform data computed from the subsurface model for a shot point (star-shaped symbol) located in the center of the profile.	
648 (c)	Pseudosection illustrating the apparent velocity computed from the seismic travel times based on the survey geometry.	
649 (d)	Subsurface model of the seismic P-wave velocity distribution obtained through the inversion of the synthetic P-wave travel times.	28
650 5	Entity-relationship diagram illustrating the SQLite database used in a SeismicRefractionManager project to store and manage processing related information.	29
651 6	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.	30
652 7	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.	31
653 8	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.	32
654 9	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.	33
655 10	The traveltime 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).	34
656 11	Comparison of forward modeled synthetic and automatically picked first break traveltimes for the synthetic seismic waveform data created in this study.	35
657 12	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.	36
658 13	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.	37
659 14	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.	38
660 15	Picking percentage for each shot in the soda lakes dataset indicating a high signal-to-noise ratio allowing for the picking of first break traveltimes for nearly all shot-receiver pairs.	39
661 16	Subsurface model of the soda lake in terms of the seismic P-wave velocity sliced along the lines formed by the geophones. The solid curve indicates the shore line of the soda lake.	40

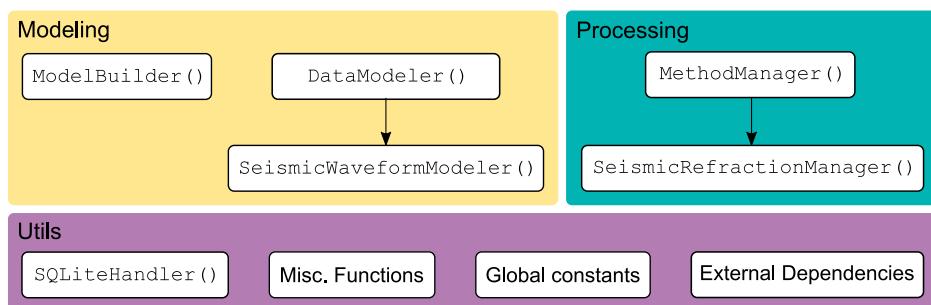


Figure 1: Fundamental use case illustrating the modeling of seismic refraction data within the formikoj ecosystem.

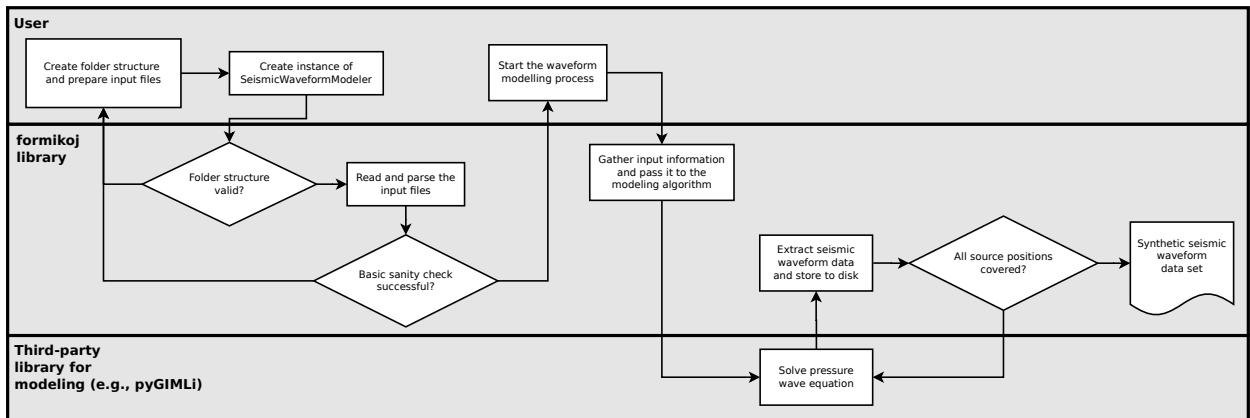


Figure 2: Fundamental use case illustrating the processing of seismic refraction data within the formikoj ecosystem.

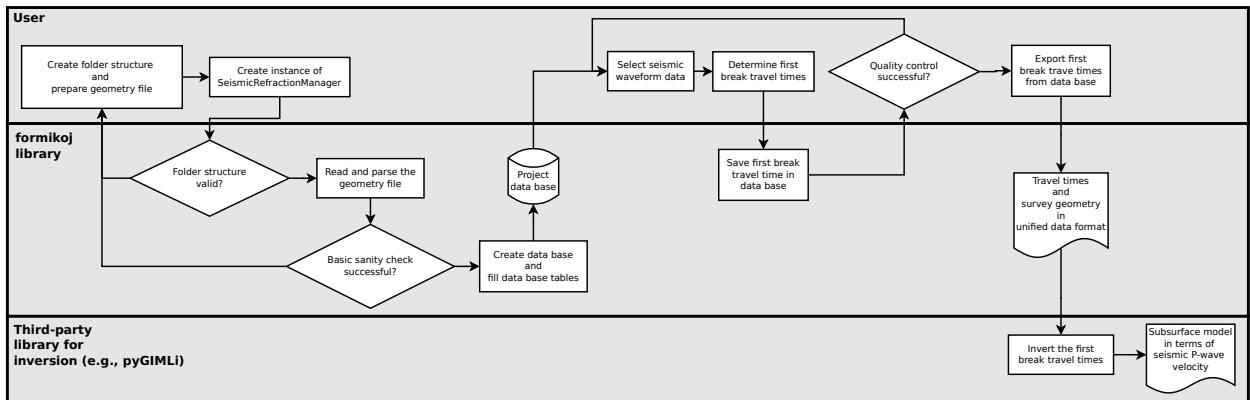


Figure 3: Typical formikoj workflow for the picking of first break traveltimes from seismic waveform data.

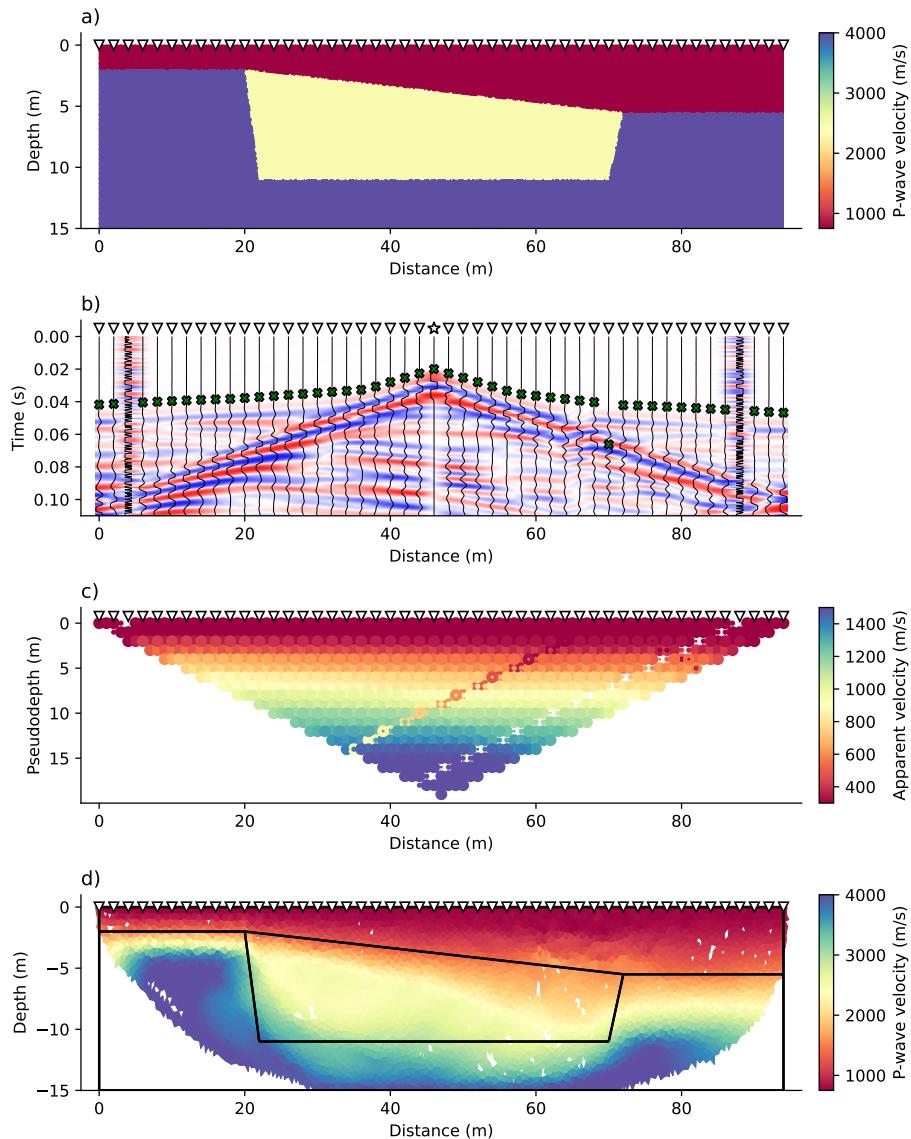


Figure 4: Synthetic seismic data generated with the `SeismicWaveformModeler`:

- (a) Two-layered subsurface model without topography characterized by a distinct anomaly in the center. Triangles along the surface indicate the positions of geophones.
- (b) Synthetic seismic waveform data computed from the subsurface model for a shot point (star-shaped symbol) located in the center of the profile.
- (c) Pseudosection illustrating the apparent velocity computed from the seismic travel times based on the survey geometry.
- (d) Subsurface model of the seismic P-wave velocity distribution obtained through the inversion of the synthetic P-wave travel times.

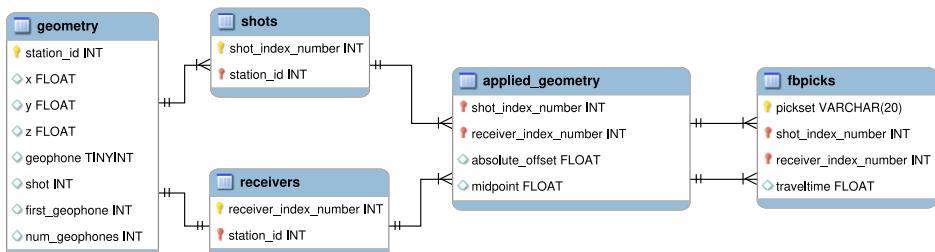


Figure 5: Entity-relationship diagram illustrating the SQLite database used in a SeismicRefractionManager project to store and manage processing related information.

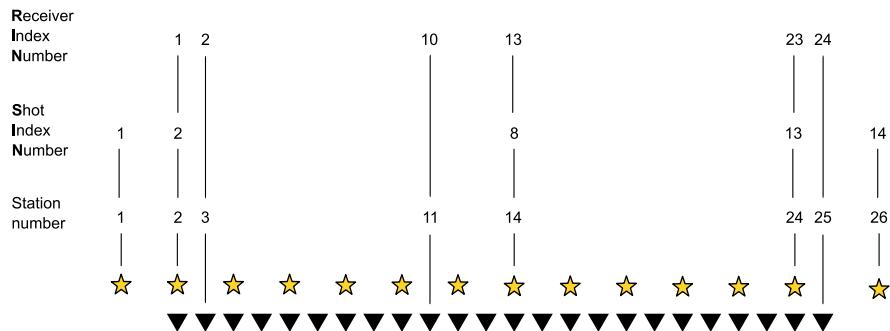


Figure 6: The SeismicRefractionManager addresses the stations through consecutive station numbers based on the sort order in the geometry file. The shot index numbers (SIN) and receiver index numbers (RIN) are assigned to the shot and receiver stations separately to allow for an intuitive data handling for the user.

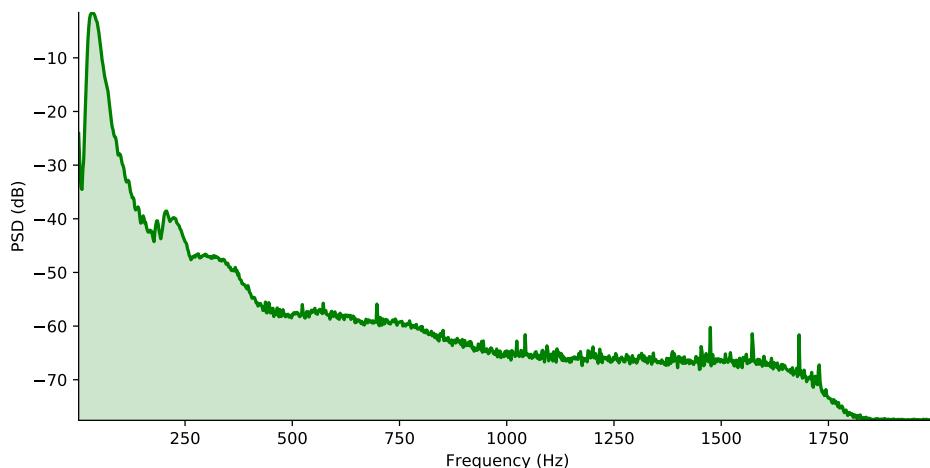


Figure 7: The frequency spectrum illustrates the frequency content of the currently selected traces, which allows for the identification of frequency ranges associated to noise, which can be omitted through the application of corresponding frequency filtering.

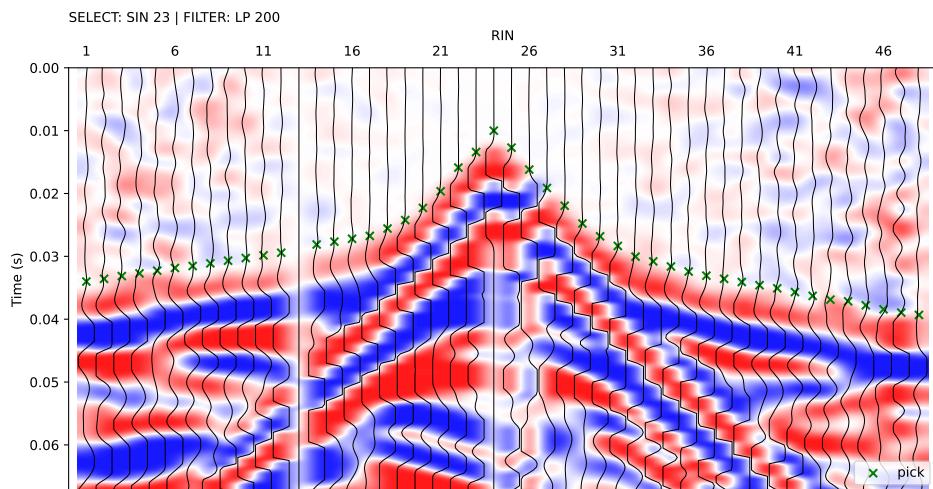


Figure 8: The seismogram plot presents the currently selected traces along the x axis, where the sort order is determined through the geometry. The corresponding trace data are illustrated as a function of time along the y axis (solid curves), with positive and negative amplitudes depicted in blue and 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 9: 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.

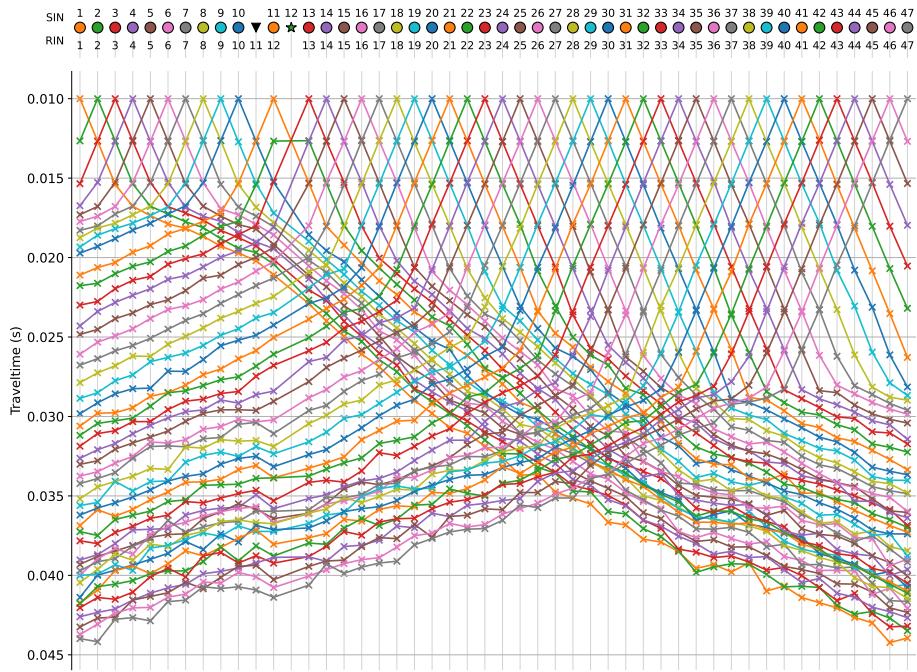


Figure 10: The travelttime diagram shows the first break traveltimes stored in the currently active pickset along the y axis (x symbols), where solid lines connect traveltimes assigned to a common shot. The sort order of the stations along the x axis reflects the geometry. Filled circles indicate stations with co-located shot and geophone (receiver), triangles refer to receiver stations (no shot) and stars refer to shot stations (no geophone).

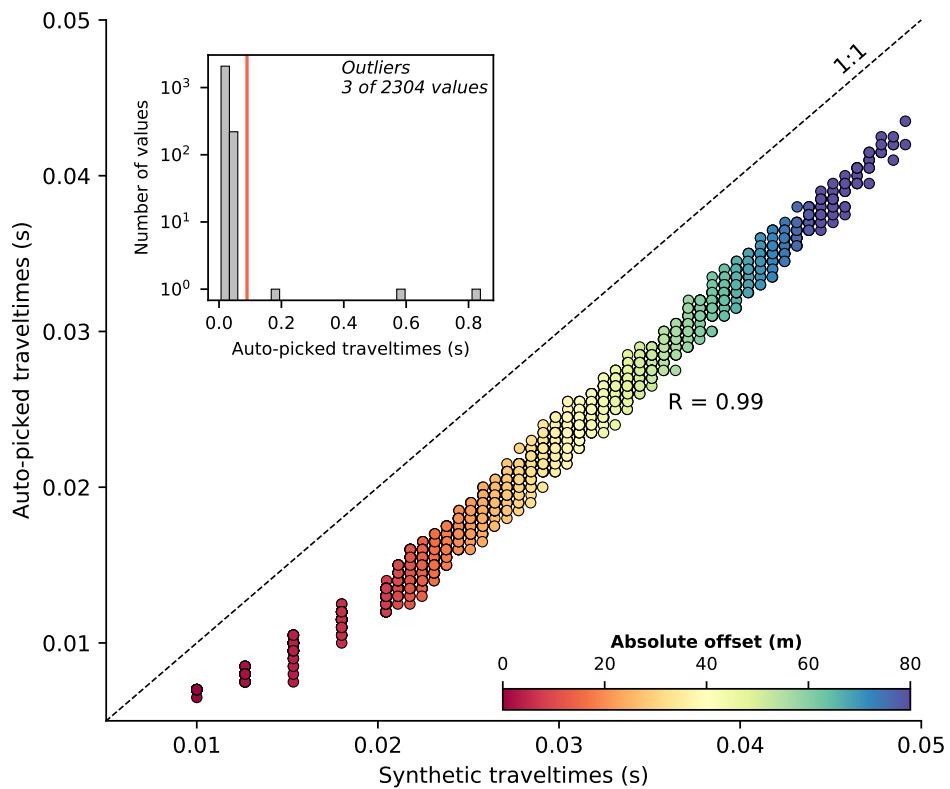


Figure 11: Comparison of forward modeled synthetic and automatically picked first break traveltimes for the synthetic seismic waveform data created in this study.

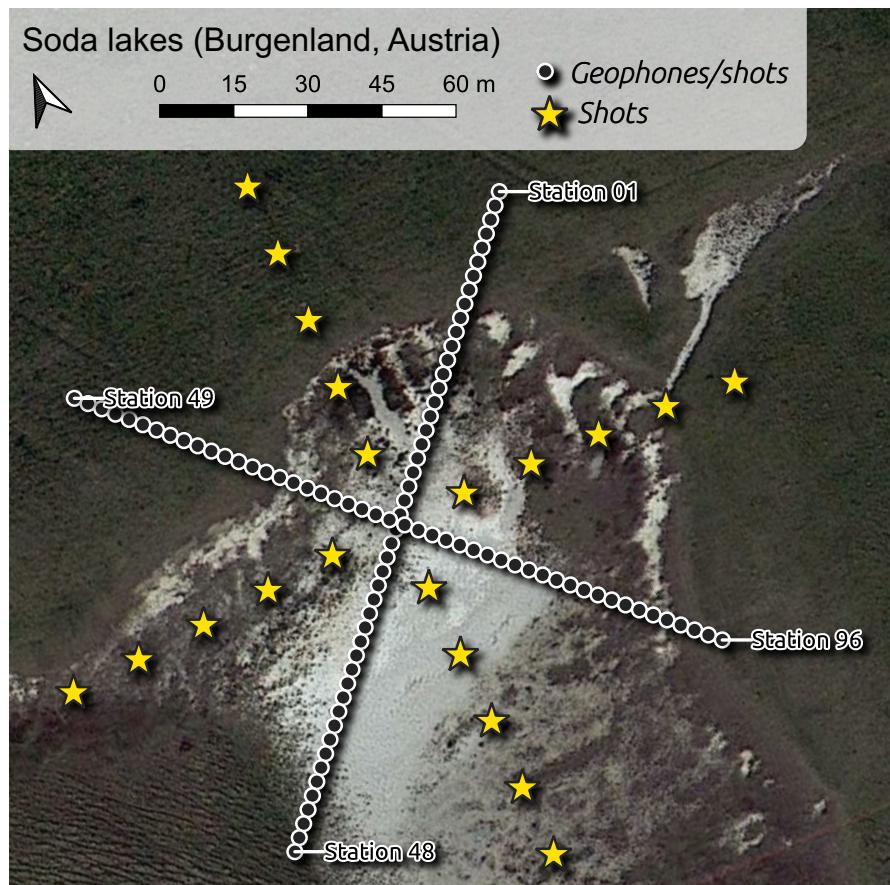


Figure 12: The soda lakes field dataset was collected in a 3D survey layout with stations (co-located geophones and shots indicated by filled dots) deployed in form of a cross. Additional shots (yellow stars) were conducted in front of a cross rotated by 45° to increase the coverage of the dataset.

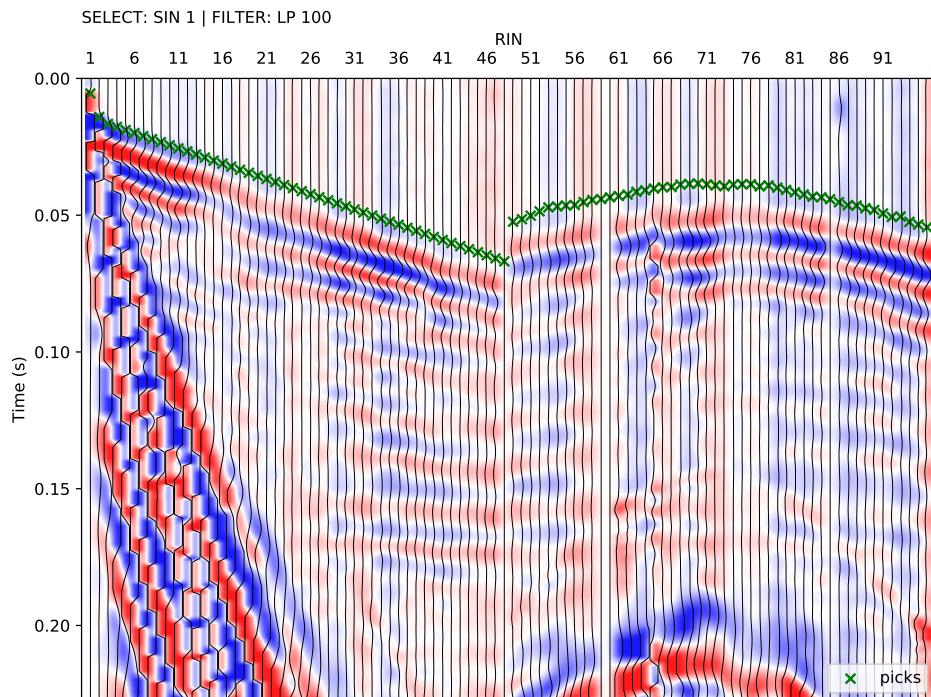


Figure 13: Exemplary seismic waveforms from the soda lakes dataset shown for shot index number (SIN) 1 with a 100 Hz lowpass filter applied to suppress high frequency noise. The recorded seismic waveforms clearly reflect the geometry of the geophones with RIN 1 to 48 deployed along the direction of wave propagation, while RIN 49 to 96 are deployed perpendicular to the propagating wavefront.

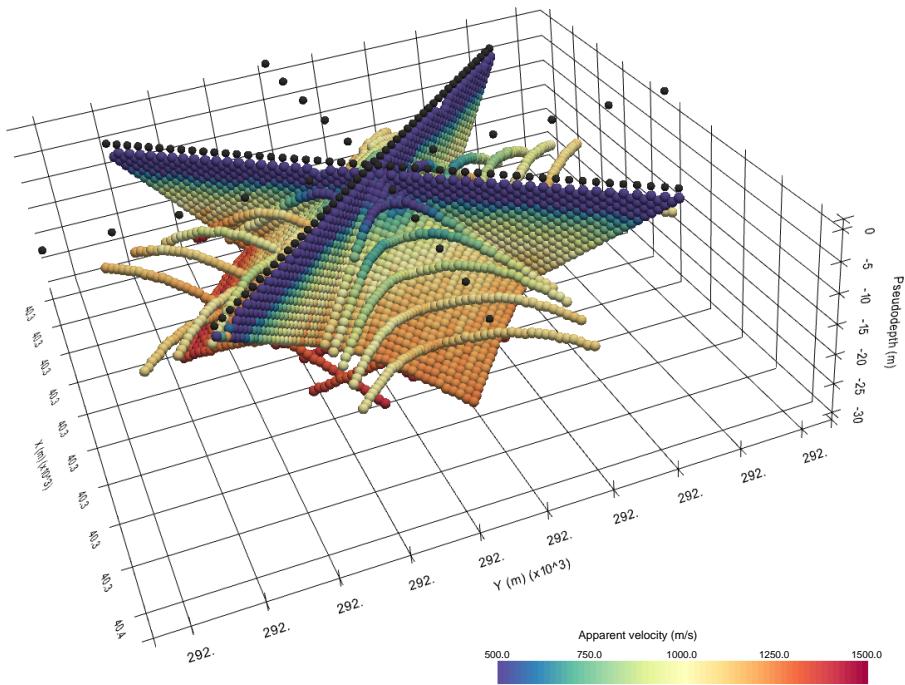


Figure 14: 3D pseudosection showing the apparent seismic velocities determined from the first break traveltimes obtained from the soda lakes dataset and the corresponding absolute offset between the shot and receiver stations. The apparent velocity for each shot-receiver pair is illustrated at the corresponding 2D midpoint and pseudodepth (1/3 of the absolute offset), thus yielding a 3D representation.

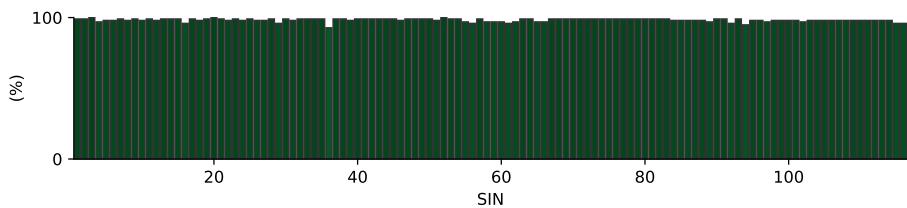


Figure 15: Picking percentage for each shot in the soda lakes dataset indicating a high signal-to-noise ratio allowing for the picking of first break traveltimes for nearly all shot-receiver pairs.

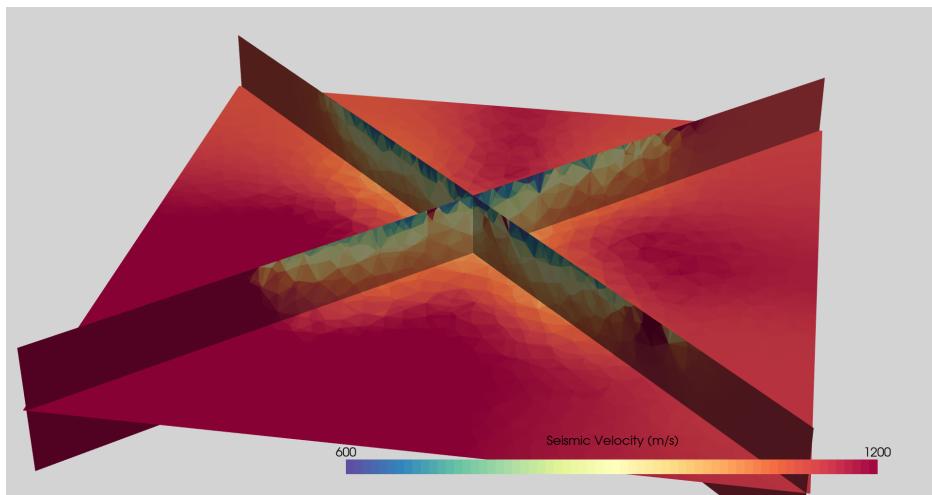


Figure 16: Subsurface model of the soda lake in terms of the seismic P-wave velocity sliced along the lines formed by the geophones. The solid curve indicates the shore line of the soda lake.