

<sup>1</sup> [COR: Number]  
<sup>2</sup> Cover Letter

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

<sup>5</sup> Matthias Steiner, Adrián Flores Orozco

<sup>6</sup> Dear Editors-in-Chief, Editor-in-Chief, dear reviewers

<sup>7</sup>  
<sup>8</sup> ~~we are submitting our manuscript~~ thank you very much for the positive evaluation of our manuscript as well as for  
<sup>9</sup> providing numerous constructive comments and suggestions for improvement. In the revised version of our manuscript  
<sup>10</sup> "formikoj: A flexible library for data management and processing in geophysics - Application for seismic refraction  
<sup>11</sup> data", ~~which we consider is a suitable contribution for Computers & Geosciences. We confirm that the submission~~  
<sup>12</sup> ~~follows all the requirements and includes all the items of the submission checklist~~<sup>we have attempted to implement</sup>  
<sup>13</sup> ~~every suggestion and address all comments.~~

<sup>14</sup>  
<sup>15</sup> The manuscript presents the open-source python library formikoj for managing and processing geophysical data collected  
<sup>16</sup> in environmental and engineering investigations. formikoj was specifically implemented for multi-platform usage to  
<sup>17</sup> allow the efficient collaboration and exchange of data between different partners in research projects and academia.  
<sup>18</sup> In this regard, we believe that this library aids in providing reproducible data and results as well as establishing and  
<sup>19</sup> maintaining good research practices. Accordingly, we consider this manuscript relevant to the audience of Computers  
<sup>20</sup> & Geosciences, and in general for geoscientists and practitioners working with geophysical methods. In particular, we  
<sup>21</sup> present theoretical details regarding the considered and implemented methodologies, provide a self-contained synthetic  
<sup>22</sup> study demonstrating the validity of the forward modeled seismic waveform data, and the fundamental processing  
<sup>23</sup> capabilities of the proposed library for seismic refraction data. Moreover, we discuss a single field data application  
<sup>24</sup> in order to provide a more concise demonstration of the library functionalities with regard to the processing of real  
<sup>25</sup> seismic survey data.

<sup>26</sup>  
<sup>27</sup> We provide the ~~Moreover, we modified the library according to the suggestions provided by the reviewers and provide~~  
<sup>28</sup> the revised source codes in a public repository with details listed in the section "Code availability".

<sup>29</sup>  
<sup>30</sup> We hope that our revisions and replies alleviate the concerns of the reviewers and that you find our study suitable for  
<sup>31</sup> publication in Computers & Geosciences. We look forward to your decision.

<sup>32</sup>  
<sup>33</sup> Yours sincerely,

<sup>34</sup>  
<sup>35</sup> Matthias Steiner and Adrián Flores Orozco  
<sup>36</sup> Research Unit of Geophysics, Department of Geodesy and Geoinformation, TU Wien; matthias.steiner@geo.tuwien.ac.at

<sup>37</sup>

<sup>38</sup> **Highlights**

<sup>39</sup> **formikoj: A flexible library for data management and processing in geophysics - Application**  
<sup>40</sup> **for seismic refraction data**

<sup>41</sup> Matthias Steiner, Adrián Flores Orozco

- <sup>42</sup>     • flexible open-source and cross-platform library for managing and processing of geophysical data  
<sup>43</sup>     • possibility to be deployed for different geophysical methods and/or instruments  
<sup>44</sup>     • application for the modeling and processing of seismic refraction datasets  
<sup>45</sup>     • applicable for seismic refraction data collected in 2D and 3D survey geometries  
<sup>46</sup>     • easily scalable for custom requirements

47 formikoj: A flexible library for data management and processing in  
48 geophysics - Application for seismic refraction data

49 Matthias Steiner<sup>a,\*</sup>, Adrián Flores Orozco<sup>a</sup>

50 <sup>a</sup>Research Unit of Geophysics, Department of Geodesy and Geoinformation, TU Wien

51

---

52 ARTICLE INFO

53

54 **Keywords:**  
geophysical data processing  
seismic refraction  
first break picking  
seismic waveform modeling  
cross-platform application  
geophysical python library  
flexible open-source libraries  
wave based methods

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

---

73 CRediT authorship contribution statement

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

76 1. Introduction

77 The acquisition of spatially quasi-continuous data in a non-invasive manner renders geophysical methods suitable  
78 for engineering and environmental investigations (e.g., Parsekian et al., 2015; Nguyen et al., 2018; Romero-Ruiz et al.,  
79 2018). However, the processing of geophysical data often relies on commercial software solutions and the associated  
80 licensing costs might render their use prohibitively expensive, which might be the case for academic projects  
81 or institutions. The most popular packages are Res2DInv<sup>1</sup> for electrical methods, Halliburton Landmark SeisSpace  
82 ProMAX<sup>1</sup> or ParkSeis<sup>1</sup> for seismic methods, or ReflexW<sup>1</sup> for ground-penetrating radar and seismic methods, and even  
83 for teaching in developing countries. A common limitation of the aforementioned existing software solutions refers  
84 to their specific platform requirements mainly related to the type and version of the operating system; moreover, the

---

\*Corresponding author

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

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

<sup>1</sup><https://www.landmark.solutions/SeisSpace-ProMAX>, last accessed on

<sup>1</sup><https://www.parkseismic.com/parkseis/>, last accessed on

<sup>1</sup><https://www.sandmeier-geo.de/reflexw.html>, last accessed on

85 possibility to adapt the code are limited if possible at all. Considering the substantial changes regarding the market  
 86 shares of operating systems within the last two decades, platform-specific software packages are becoming particularly  
 87 obstructive for academic research and teaching. The increasing popularity of the Python programming language led  
 88 to the development of various cross-platform open-source software packages for processing, modeling and inverting  
 89 geophysical data. Available packages can focus on specific geophysical methods, for instance, ResIPy (Blanchy et al.,  
 90 2020) for electrical data, GPRPy (Plattner, 2020) for ground-penetrating radar data, or ObsPy (Beyreuther et al., 2010)  
 91 and Pyrocko (Heimann et al., 2017) for seismological data. Other packages provide frameworks for the inversion and  
 92 permit the inclusion of forward models for different geophysical methods, e.g., SimPEG (Cockett et al., 2015), Fatiando  
 93 a Terra (Uieda et al., 2013) or pyGIMLi (Rücker et al., 2017).

94 The seismic refraction tomography (SRT) is a standard technique in environmental and engineering applications.  
 95 Often applied together with other geophysical methods, the SRT is routinely used, e.g., in permafrost studies (e.g.,  
 96 Draebing, 2016; Steiner et al., 2021), for the investigation of landfills (e.g., Nguyen et al., 2018; Steiner et al., 2022),  
 97 or for hydrogeological characterizations (e.g., Bücker et al., 2021). The market for seismic processing software has  
 98 long been dominated by software packages designed for the processing of large datasets, e.g., associated to oil or gas  
 99 exploration. Accordingly, these seismic processing solutions may not be suited for small-scale projects in environmen-  
 100 tal and engineering studies, or for teaching activities. ReflexW overcomes such limitations by providing processing  
 101 tools specifically designed for near-surface investigations at substantially lower costs. In terms of licensing costs,  
 102 Stockwell (1999) went a step further by making the Seismic Unix framework available entirely free of charge; whereas  
 103 Guedes et al. (2022) recently presented RefraPy, a python processing tool for seismic refraction data. Implemented in  
 104 python, RefraPy is potentially suitable for cross-platform usage, yet Guedes et al. (2022) developed and tested solely for  
 105 Windows operating systems. Moreover, RefraPy does not offer the possibility to generate synthetic seismic waveform  
 106 data, as required for survey design, as well as for teaching and interpretation purposes, testing of research hypotheses,  
evaluating the accuracy of different measurement configurations or inversion strategies.

108 The formikoj library presented here is an open-source framework for creating synthetic datasets, as well as for man-  
 109 aging and processing numerical and field data independently from the operating system and without licensing costs;  
 110 thus, overcoming limitations associated to existing solutions. The design of the library follows the multi-method  
 111 concept of pyGIMLi and SimPEG, which allows for the implementation of custom designed tools for different geo-  
 112 physical methods. The usage of transparent file formats, e.g., the unified data format (udf<sup>1</sup>), and data management  
 113 concepts (SQLite database) within the formikoj framework facilitates a simple data exchange between partners in re-  
 114 search projects and academia, which is required to guarantee the repeatability of results and good research practices.  
 115 Considering the diverse applications of the **SR-SRT** method we demonstrate the applicability of the proposed library

<sup>1</sup>[http://resistivity.net/bert/data\\_format.html](http://resistivity.net/bert/data_format.html), last accessed on March 11, 2023

116 based on tools implemented within the formikoj framework for the modeling and processing seismic waveform data. In  
 117 particular, we present here a series of illustrative use cases based on a carefully designed synthetic study highlighting  
 118 the capabilities of the formikoj library referring to (i) the modeling of synthetic seismic refraction (SR) waveform data,  
 119 (ii) the processing of a 2D SR field dataset collected with a roll-along survey geometry, and (iii) the processing of a  
 120 for the forward modeling and processing of 2D seismic refraction datasets. Based on a 3D field dataset we illustrate  
 121 that the formikoj library also allows for the processing of complex survey geometries, where the obtained first break  
 122 traveltimes can be inverted with third party open-source libraries to solve for 3D SR field data set subsurface models  
 123 expressed in terms of the seismic P-wave velocity.

## 124 2. Design and structure of the formikoj library

125 As illustrated in Figure 1, the formikoj library comprises a modeling and a processing module, which both rely  
 126 on a common utilities module. The DataModeler and the MethodManager class provide the basis to add mod-  
 127 eling or processing functionalities for specific any kind of geophysical methods. The formikoj library provides a  
 128 well thought out data management concept based on the SQLite database engine that does not require a separate  
 129 server process<sup>2</sup>. In particular, we present here the SeismicWaveformModeler and SeismicRefractionManager  
 130 classes implemented within the formikoj framework, which aim at creating and processing seismic waveform data the  
 131 information for each project, such as survey geometry and first break traveltimes, are stored in a SQLite database file.  
 132 Using such an application file format facilitates the cross-platform design of the formikoj library, and provides fast  
 133 I/O operations through concise SQL queries. The portability of the application file allows for an easy exchange of  
 134 projects between partners across institutions relying on different IT infrastructures. Accordingly, the formikoj library  
 135 aims at providing a transparent and customizable framework for the collaborative design of reproducible workflows.  
 136 The comprehensive usage of clear error and log messages aims at supporting the user throughout the modeling and  
 137 processing workflows. A meticulous exception handling ensures that the data stored in the SQLite project database is  
 138 not corrupted in case of erroneous input. Moreover, documenting the user input and the respective responses of the  
 139 formikoj library with the python logging module provides a timestamped command history that further enhances the  
 140 transparency and repeatability of the conducted workflows.

141 We present the application of the formikoj library for the modeling and processing of seismic refraction data based  
 142 on the fundamental use cases presented in Figure 2 and Figure 3, respectively. Similar to RefraPy, These flow charts  
 143 illustrate the corresponding workflows as well as the required interactions between the user, the formikoj library and  
 144 third-party packages. As can be seen from Figure 2 and Figure 3, the formikoj library acts as an interface between  
 145 the user and more complex functionalities of third-party libraries, such as pyGIMLi for modeling and inversion of

<sup>2</sup><https://www.sqlite.org/index.html>, last accessed on March 11, 2023

146 seismic refraction data. The `SeismicWaveformModeler` and `SeismicRefractionManager` classes implement these  
 147 fundamental use cases, yet their actual capabilities are continuously expanded, e.g., to address specific modeling or  
 148 processing requirements as well as to enhance the user experience. To avoid redundancies in the implementation  
 149 we built these classes are built upon the functionalities of existing packages such as ObsPy for the processing of  
 150 seismological data (Beyreuther et al., 2010) and pyGIMLi for the modeling and inversion of different geophysical  
 151 data (Rücker et al., 2017). Other important third party dependencies refer to NumPy (Harris et al., 2020) and Pandas  
 152 (McKinney, 2010) for general data handling, as well as matplotlib (Hunter, 2007) and PyVista (Sullivan and Kaszynski,  
 153 2019) for data visualization. In the current version, we implemented and tested formikoj primarily on Linux machines  
 154 , yet the library has been successfully used on all major operating systems, i.e., Linux, Mac OS and Windows.

155 The formikoj library is primarily developed on machines running Linux Mint 20 or Kubuntu 22.4, respectively.  
 156 Testing of the library refers to carrying out the fundamental use cases for modeling and processing of seismic refraction  
 157 data presented in Figure 2 and 3. To ensure the cross-platform applicability of the formikoj library we test these  
 158 fundamental use cases also on machines running on Windows 10, Windows 11 as well as macOS versions 12 and 13.

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

## 162 2.1. Generation of seismic waveform data for synthetic subsurface models

## 163 2.2. Generation of seismic waveform data for synthetic subsurface models: The

### 164 SeismicWaveformModeler

165 The SR method exploits the ground motion recorded by sensors installed in the surface (e.g., geophones) to char-  
 166 acterize the propagation of seismic waves generated at well known locations (i.e., shot stations). The visualization  
 167 of the ground motion as function of time yields a so-called seismogram for each geophone position. Accordingly,  
 168 the `SeismicWaveformModeler` class provides a flexible way to generate synthetic seismic waveform data for P-wave  
 169 refraction modeling either in a python script, interactively in a jupyter notebook or an ipython shell. To create an  
 170 instance of the class the user provides the absolute or relative path to the working directory is provided as parameter  
 171 to the constructor: The working directory needs to contain a subdirectory `in`, whereas the output directory `out` will be  
 172 created automatically:-

```
173 1 # Import the SeismicWaveformModeler from the formikoj library
174 2 from formikoj import SeismicWaveformModeler
175 3 %DIF >
176 4 # Create an instance of the SeismicWaveformModeler
177 5 swm = SeismicWaveformModeler('..')
```

**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

178 1 INFO : Created instance of SeismicWaveformModeler

179     Description of the information to be provided in the columns of a measurement scheme in csv format. Column  
 180     **Content Data type Description** 1 x coordinate float Station x coordinate, e.g., given in (m) 2 y coordinate float Station  
 181     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  
 182     of a receiver station, 0 otherwise 5 Shot bool 1 in case of a shot station, 0 otherwise The In the working directory,  
 183     the required input files need to be are provided via the subdirectory *in* of the working directory, whereas the modeling  
 184     output will be stored in the automatically created subdirectory *out*. The key input file is the measurement scheme  
 185     , which as it contains information regarding the distribution of the shot and geophone stations. If provided in the  
 186     unified data format, the measurement scheme is imported directly with pyGIMLi into a DataContainer. In case the  
 187     measurement scheme is provided as a csv file, the SeismicWaveformModeler reads the information and writes it to  
 188     a pyGIMLi DataContainer . In the csv format object. If provided as csv file, the measurement scheme contains a  
 189     single line for each station in the survey layout, where a station either hosts a geophone or a shot, or both (see Table 1).  
 190     The values provided in each line need to be separated by a unique delimiter, and the file must not contain a header.

191     For the modeling of the seismic waveform data, the parameters for characterizing the base wavelet,  
 192     the synthetic subsurface model and the resulting waveform datasets are provided (see Table 2) in a configuration file  
 193     following the yaml format: In the exemplary configurationfile shown above

```
194 1 wavelet:  

195 2   length: 1.024  

196 3   frequency: 100  

197 4   sampling_rate: 2000  

198 5   pretrigger: 0.02  

199 6 %DIF >  

200 7 model:  

201 8   velmap: [[1, 750], [2, 2500], [3, 4000]]  

202 9   layers: [[1, 3], [2, 5], [3, 15]]  

203 10  quality: 32  

204 11  area: 10
```

```

205 12    smooth: [1, 10]
206 13    sec_nodes: 3
207 14 %DIF >
208 15 dataset:
209 16     number: 1
210 17     names: [syn_data]
211 18     noise: 1
212 19     noise_level: 1e-4
213 20     missing_shots: 1
214 21     broken_geophones: 1
215 22     wrong_polarity: 1
216 23 %DIF >
217 24 traveltimes:
218 25     noise_relative: 0.
219 26     noise_absolute: 0.

```

220 In this exemplary configuration, the first block contains information regarding the (wavelet) contains the parameterization  
 221 of the base wavelet, which controls the modeling of the seismic waveforms as described in (see Table 2). In the second  
 222 block). The second block (model) contains information regarding the synthetic subsurface model. Here, the user can  
 223 define simple models with all layers considered to be parallel to the surface topography, which is inferred from the  
 224 station geometry provided in the measurement scheme. The seismic velocity values for the different model regions  
 225 (velmap) and the corresponding thickness of the different geological units (layers) have to be explicitly defined in  
 226 the configuration file. The remaining parameters (quality, we\_area, smooth and sec\_nodes) refer to the properties  
 227 of the mesh that is eventually used for the forward modeling of the seismic waveform data and corresponding first  
 228 break traveltimes (we refer to the respective pyGIMLi resources<sup>3</sup> for further information). In the third block of the  
 229 exemplary configuration file (dataset), the user can set specific names for the datasets to be created and (the number  
 230 of datasets is automatically determined. Alternatively, it is possible to), or set the number of datasets to be created and  
 231 the dataset names are automatically generated with the prefix dataset\_. Various parameters-The remaining parameters  
 232 provided in the dataset block control the random error (noise )and systematic errors and noise\_level) as well as  
 233 systematic errors (missing\_shots, broken\_geophones and wrong\_polarity) in the modeled seismic waveform  
 234 data (see Table 2 for a detailed description). The number and position of the shot and geophone stations affected  
 235 by the systematic errors are randomly chosen with a maximum of 5%–5 % of the total number of stations in order to  
 236 avoid a high number of invalid trace data. The third block contains information regarding the synthetic subsurface  
 237 model. For each layer the corresponding velocity (final block of the configuration file (velmaptraveltimes) and  
 238 layer thickness (layers) need to be provided and all layers are considered to be parallel to the surface topography

<sup>3</sup>[https://www.pygimli.org/pygimliapi/\\_generated/pygimli.meshTools.html#pygimli.meshTools.createMesh](https://www.pygimli.org/pygimliapi/_generated/pygimli.meshTools.html#pygimli.meshTools.createMesh), last accessed March 11, 2023

239 (geometrical information regarding the stations in the measurement scheme). The remaining parameters, namely  
 240 quality, area, smooth contains data-error parameters for both the absolute error ( $e_{abs}$ ) and sec\_nodes, define the  
 241 properties of the mesh to be used for the forward modeling (we refer to the respective pyGIMLi resources<sup>4</sup> for further  
 242 information). Alternatively, the user can provide a more complex mesh in the binary mesh format (the relative error  
 243 ( $e_{rel}$ ) defined by the user. The data error model is then estimated as

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

244 where  $\underline{d}$  denotes the forward modeled data not affected by noise, i.e., a bms file): The parameters in the final block  
 245 control the forward modeling of the corresponding seismic traveltimes (see Table 2). A configuration file located in  
 246 the input directory can be imported through the load method: here the first break traveltimes obtained through the  
 247 Dijkstra algorithm (Dijkstra, 1959). These computed error values are subsequently added to the data to obtain the  
 248 noisy data as

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

249 Once the  
 250 The configuration file needs to be provided via the input directory from where it can be imported through the load  
 251 method of the SeismicWaveformModeler instance is parameterized the synthetic:

```
252 6 # Load and parse the configuration file
253 7 swm.load(type='config')

254 1 INFO : Configuration loaded
```

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

```
261 8 # Load the mesh into the workflow
262 9 swm.load(type='mesh')
```

<sup>4</sup>[https://www.pygimli.org/pygimliapi/\\_generated/pygimli.mesh.html#pygimli.mesh.createMesh](https://www.pygimli.org/pygimliapi/_generated/pygimli.mesh.html#pygimli.mesh.createMesh), last accessed-

**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
<b>wavelet</b>		
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
<b>dataset</b>		
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
<b>model</b>		
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
<b>traveltimes</b>		
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

263 1 INFO : Mesh loaded

264 The forward modeling process generating the seismic waveform data is then invoked through the `create` method:265 10 # Start the modeling of the seismic waveform data  
266 11 swm.create(type='waveforms')

267 1 INFO : Measurement scheme loaded

268 2 INFO : Velocity model created

269 3 INFO : Wavelet created

270 4 [+++++ 100% +++++] 2048 of 2048 complete

```

271 5 ...
272 6 [+++++ 100% ++++++ 2048 of 2048 complete
273 7 INFO : Dataset 'syn\_data' created

```

274 As can be seen from the log messages, the SeismicWaveformModeler first loads the measurement scheme and creates  
 275 the velocity model used into the workflow. In the second step, the provided mesh and the seismic velocity values defined  
 276 in the configuration file are combined to create the seismic velocity model considered for the waveform modeling.  
 277 Based on the wavelet properties a Ricker wavelet is generated in this study (see Figure 4a). The model consists of a  
 278 top and a bottom layer with varying thickness characterized by seismic velocity values of  $750 \text{ ms}^{-1}$  and  $4000 \text{ ms}^{-1}$ ,  
 279 respectively. In the center, the model contains an irregularly shaped anomaly associated with a seismic velocity of  
 280  $2500 \text{ ms}^{-1}$ , i.e., the model features vertical and lateral variations in the seismic velocity distribution. The third step  
 281 refers to the generation of a Ricker wavelet through the pyGIMLi function ricker. Subsequently, meshas

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

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

285 Based on the mesh, the velocity model and Ricker wavelet are used the Ricker wavelet we use pyGIMLi to  
 286 solve the pressure wave equation for each shot station defined in the measurement scheme with the pyGIMLi function  
 287 solvePressureWave. The resultant waveforms seismograms are extracted at the corresponding geophone stations are  
 288 extracted and stored, gathered for each shot position in an ObsPy Stream data structure. A directory for each dataset  
 289 defined in the configuration file will be created in the object and saved in the output directory ( $out$ ) :

290 In as shown here:  
 working\_directory  
 |  
 | in  
 | out  
 | | syn\_data  
 | | | data  
 | | | | protocol.txt  
 | | | | station\_coords.csv  
 | | | | Shot\_1001.syn  
 | | | | ...  
 | | | | Shot\_10nn.syn  
 | | | | syn\_data\_tt.pck  
 | | | | info.txt

291 In particular, the subdirectory data, the synthetic seismic waveforms are stored in a separate shot file for each shot  
 292 position contains a separate file in the miniseed format (Ahern et al., 2012; Ringler and Evans, 2015) with for each shot

293 position, where the file extension *syn* to identify the forward modeled shot files, e.g., *Shot\_1001.syn*, indicates that  
 294 the shot files contain forward modeled seismic waveform data. The measurement protocol (protocol.txt) and the station  
 295 coordinates provided as a csv file (station\_coords.csv) are also stored in this directory. The header of the measurement  
 296 protocol contains the survey parameters required for the processing of the seismic waveforms, namely, e.g., sampling  
 297 rate, recording length, number of geophones and geophone spacing. Moreover, the protocol associates each shot file  
 298 of the dataset to with a specific location within in the survey geometry, i.e., with respect to the geophone positions.  
 299 e.g.:

```
300 1 #####  

301 2 Line: SYN_syn_data  

302 3 Sampling_rate: 2000 Hz  

303 4 Recording_length: 0.512 s  

304 5 Number_of_geophones: 48  

305 6 Geophone_spacing: 2 m  

306 7 #####  

307 8 %DIF >  

308 9 File_number | Station  

309 10 1001 | G001  

310 11 : | :  

311 12 1048 | G048
```

312 The auxiliary file info.txt provided in exported to the dataset directory summarizes the parameters from the configura-  
 313 tion file and as well as information regarding the simulated systematic errors in the synthetic seismic waveform data:  
 314 Synthetic traveltimes (`swn.create('traveltimes')`) are stored in the dataset directory as udf files, e.g.:

```
315 1 Number_of_geophones: 48  

316 2 Number_of_shots: 48  

317 3 Recording_length_(s): 0.512  

318 4 Sampling_frequency_(Hz): 2000  

319 5 Wavelet_type: Ricker  

320 6 Frequency_of_the_wavelet_(Hz): 100  

321 7 %DIF >  

322 8 Missing_shot(s): 11  

323 9 Broken_geophone(s): 13  

324 10 Wrong_polarity_geophone(s): 26
```

325 Figure 4b presents the seismic waveform data forward modeled for a shot point located in the center of the profile.  
 326 The seismograms are shown as curves, whereas the strength of the amplitudes is added as color-coded information.  
 327 In the seismograms we see clear first onsets along the entire profile, yet receivers 3 and 45 were modeled as broken  
 328 geophones, i.e., the corresponding seismograms contain solely noise. Crosses overlaid on the valid seismograms at  
 329 the respective first onset indicate the first break traveltimes *tt* between the shot and the receiver. In addition to the  
 330 missing first break traveltimes for the broken receivers, we manually added a systematic error, i.e., we set an erroneous

331 traveltimes for receiver 36, to demonstrate how such outliers can be identified in a so-called pseudosection.  
 332 Pseudosections are a useful tool for the analysis of the data quality by presenting the data based on the positions of  
 333 shots and receivers. Such a visualization allows for the detection of outliers and their spatial distribution as necessary  
 334 to understand possible sources of error. In case of the SRT, a pseudosection as presented in Figure 4c, illustrates  
 335 apparent velocity ( $v_{app}$ ) values computed as

$$v_{app} = \frac{aoffset}{tt}, \quad (4)$$

336 where *aoffset* refers to the absolute offset between shot and receiver. The  $v_{app}$  values are plotted at pseudolocations,  
 337 where the location along profile direction is defined by the midpoint of the corresponding shot-receiver pair and  
 338 the pseudodepth is computed as 1/3 of *aoffset*. As demonstrated in Figure 4c, a pseudosection allows for the  
 339 identification of missing data (e.g., dataset1\_tt.pck). The file extension *pck* is an abbreviation of the word 'pick' and  
 340 receiver 3 and 45) as well as systematic errors or outliers, i.e., velocities erroneously influenced by a single shot  
 341 or receiver (see receiver 36). The main assumption here is that the pseudosection should reveal smooth transitions  
 342 between lateral and vertical neighbors, considering that the data were collected with gradual changes in the position  
 343 of the source (i.e., hammer blow) and the receiver (i.e., geophone). The pseudosection will reveal large variations in  
 344 case of abrupt changes in the topography or the geometry of the array, yet this can be taken into account by the user  
 345 during the identification of outliers and possible systematic errors.

346 Based on pseudosections erroneous traveltimes can be identified and subsequently removed or corrected, which is  
 347 a critical processing step prior to the inversion of the data. The inversion of first break traveltimes aims at resolving a  
 348 model of the P-wave velocity of the subsurface materials, where systematic errors in the data would lead to the creation  
 349 of artifacts in the imaging results or hinder the convergence of the inversion. Considering the multitude of available  
 350 inversion frameworks we do not include a new inversion strategy in the formikoj library. Instead the processed data can  
 351 be exported in any format required for the use of commercial inversion algorithms, or can be linked directly to other  
 352 open-source libraries. In our case, we refer to the modeling and inversion capabilities of pyGIMLi (Rücker et al., 2017)  
 353 . pyGIMLi uses a generalized Gauss-Newton method to solve the inversion problem through the minimization of an  
 354 objective function given as:

$$\| \mathbf{W}_d (\mathcal{F}(\mathbf{m}) - \mathbf{d}) \|_2^2 + \lambda \| \mathbf{W}_m (\mathbf{m} - \mathbf{m}_0) \|_2^2 \rightarrow \min \quad (5)$$

355 The first term on the right-hand side of Equation 5 refers to the data misfit.  $\mathbf{W}_d$  is the data weighting matrix holding  
 356 the reciprocals of the data errors,  $\mathbf{d}$  denotes the data vector holding the input data (first break traveltimes saved in the  
 357 file),  $\mathbf{m}$  is the model vector representing the target parameter (seismic P-wave velocity values), and  $\mathcal{F}(\mathbf{m})$  is the model  
 358 response. The second term represents the regularization, where  $\mathbf{W}_m$  and  $\mathbf{m}_0$  are the model constraint matrix  
 359 and the reference model, respectively. The regularization parameter  $\lambda$  balances the influence of the data misfit and the  
 360 regularization term on the inversion process.

361 The inversion of the synthetic first break traveltimes considered here resolves the subsurface model presented in  
 362 Figure 4d, which is given in terms of the spatial variations of the seismic P-wave velocity, i.e., reflecting the associated  
 363 lateral and vertical variations. Depending on the survey geometry the inversion can resolve subsurface models in both  
 364 2D or 3D. To aid in the evaluation of this inversion result we superimposed the known interfaces between the different  
 365 subsurface units of the synthetic model. As can be seen from this plot, the imaging result resolves the fundamental  
 366 structural features and reflects the P-wave velocity distribution of the synthetic subsurface model. Deviations from the  
 367 true velocity model are due to the smoothness-constraint inversion scheme applied by pyGIMLi. To obtain sharper  
 368 contrasts between the different subsurface units in the imaging result structural constraints could be incorporated in the  
 369 inversion, e.g., as demonstrated by (Steiner et al., 2021); yet, the exploration of different inversion strategies is beyond  
 370 the scope of this study. We consider Figure 4 to demonstrate the applicability of the `SeismicWaveformModeler` class  
 371 for generating synthetic seismic waveform data to be used in numerical P-wave refraction seismic investigations.

## 372 2.2. Processing of seismic refraction datasets: the `SeismicRefractionManager`

373 The SR seismic refraction (SR) method is based on the measurement of the traveltimes of seismic waves determined  
 374 from the the first onset of the seismic energy recorded by geophones. In the SRT, the inversion of traveltimes gathered  
 375 from tens to hundreds of seismograms permits the computation of variations in the seismic velocities in an imaging  
 376 framework. Measuring the traveltimes in seismograms (i.e., first break picking) is commonly done manually (or semi-  
 377 automatically) in an iterative process. The signal-to-noise ratio in the recorded seismograms substantially influences  
 378 the quality of the traveltimes picked in the seismograms, and thus the seismic velocity model obtained through the  
 379 inversion. Accordingly, a proper enhancement of the perceptibility of the first onsets is crucial.

380 Accordingly, the `The SeismicRefractionManager` class provides functionalities that permit the processing of  
 381 seismic waveform data for a refraction seismic analysis. These functionalities involve the reading of seismic waveform  
 382 data, combining the data with information about the survey geometry, processing of the waveforms as well as the pick-  
 383 ing of first break traveltimes. Since the first break picking is a highly interactive process, the `SeismicRefractionManager`  
 384 is designed primarily for usage from within an ipython shell. The first break traveltimes determined with the  
 385 `SeismicRefractionManager` represent the input data, e.g., for the `TravelTimeManager` of pyGIMLi (Rücker et al., 2017)

386 . This is particularly relevant as the `TravelTimeManager` provides an inversion framework for first break traveltimes,  
 387 yet not a framework for the processing of seismic waveform data. For the demonstration of the fundamental `SeismicRefraction`  
 388 capabilities we consider the synthetic seismic data created with the  
 389 `SeismicWaveformModeler` above, which illustrates that both classes can be combined in subsequent workflows.

390 **2.2.1. Compiling the survey information and creating a project**

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

395 The data preview mode is primarily initiated if the provided working directory contains seismic shot files. To create  
 396 or load a project `SeismicRefractionManager` project, the working directory needs to contain specific subdirectories:

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

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

399 The geometry file is a csv file that provides an abstract representation of the survey layout and must not contain  
 400 a header. The fundamental element for the description of the survey layout is the station, which refers either to a  
 401 geophone position, a shot position or a position with co-located shot and geophone. For each station the geometric and  
 402 semantic information described in Table 3 are provided column-wise, i.e., each line in the geometry file corresponds to  
 a single station with a unique position within the survey layout. For the synthetic dataset considered in this study, the

**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)

403  
 404 2D station coordinates are provided in the geometry file, whereas the z coordinate is assumed to be 0 m along the entire  
 405 profile, as illustrated in Table 4; thus, reflecting the lack of surface topography in the synthetic model (see Figure 4a).

**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

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

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

```
420 1 # Import the SeismicRefractionManager from the formikoj library
421 2 from formikoj import SeismicRefractionManager
422 3 %DIF >
423 4 # Create an instance of the SeismicRefractionManager
424 5 srm = SeismicRefractionManager('')
```

```
425 1 INFO    : Read geometry information from file
426 2 INFO    : Extracted shot geometry
427 3 INFO    : Extracted receiver geometry
428 4 INFO    : Applied geometry
429 5 INFO    : Standard pickset 'picks' created
430 6 INFO    : Pickset 'picks' loaded
```

```

431 7 INFO    : 'picks' set as active pickset
432 8 Progress <===== 100.0% completed
433 9 INFO    : Read 48 files

```

434 In a first step, the SeismicRefractionManager creates an SQLite database ([prj.db](#)-[prj.db](#)) in the working directory  
 435 ~~to store the geometry information, whereby the stations are consecutively numbered based on the entity-relationship~~  
 436 ~~diagram shown in Figure 5. The geometry information is then read from the geometry file and stored in the database~~  
 437 ~~table *geometry* with consecutively numbered stations (see Figure 6). To allow for an efficient data selection through~~  
 438 ~~for the user the SeismicRefractionManager links creates database tables *shots* and *receivers*, which link the sta-~~  
 439 ~~tion numbers to shot index numbers (SIN) and receiver index numbers (RIN) assigned to shot and receiver stations,~~  
 440 ~~respectively (see Figure 6). Based on this information, the geometry is applied, respectively. For each shot-receiver~~  
 441 ~~pair the corresponding SIN and RIN are stored in the table *applied\_geometry* together with the absolute offset and~~  
 442 ~~midpoint between these stations, i.e., the database tables required for the processing of the seismic waveform data are~~  
 443 ~~created. In particular, geometry is applied. In the last step, the database table *fbpicks* is created, which stores the first~~  
 444 ~~break traveltimes for each SIN-RIN pair are stored in a dedicated database table *fbpicks* together with the name of the~~  
 445 ~~corresponding pickset, i.e., the a common label for an entire set of first break traveltimes. By default, each project~~  
 446 ~~contains the default pickset 'picks', which is loaded and activated on startup. Once the database is initialized, the~~  
 447 ~~waveform data are read from disk and the project is ready for processing.~~

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

450 For a project with a successfully applied geometry,

#### 451 2.2.2. Selecting and visualizing seismic waveform data

452 Once the geometry is applied the select method of the SeismicRefractionManager allows to gather the seis-  
 453 mic ~~waveforms~~ waveform data based on a common shot (absolute offset ([sin](#)-[aoffset](#)), a common receiver RIN  
 454 ([rin](#)) or the common absolute offset (, or a common SIN ([aoffset](#)-[sin](#)) -)

```

455 6 # Select traces with common absolute offset
456 7 srm.select(by='aoffset', num=6)

```

```

457 1 INFO    : 88 traces selected

```

```

458 8 # Select traces with receiver
459 9 srm.select(by='rin', num=10)

```

```

460 1 INFO    : 48 traces selected

```

```
461:0 # Select traces with receiver
462:1 srm.select(by='sin', num=23)
```

```
463:1 INFO : 48 traces selected
```

### 2.2.3. *Visualization and enhancement of the seismic waveform data*

Calling the Figure 7 shows the frequency spectrum of the currently selected traces, which can be computed and visualized through the `plot` method without passing any parameter allows for the visualization of the:

```
467:2 # Plot the frequency spectrum
468:3 srm.plot(type='spectrum')
```

The SeismicRefractionManager resolves the frequency spectrum by computing the stacking the power spectral density (*psd*) of the seismograms  $s(t)$  as

$$psd = 10 \log_{10} \left( \frac{|\text{fft}(s(t))|^2}{N} \right), \quad (6)$$

where `fft` refers to the fast fourier transformation (FFT) and  $N$  is the number of the considered seismograms. The frequency spectrum provides information regarding the amplitudes associated to different frequencies in the seismograms, which can be used for the identification of the dominating frequencies as required for the selection and application of adequate filters such as lowpass, highpass, bandpass and bandstop implemented in ObsPy (Beyreuther et al., 2010). To enhance the signal-to-noise ratio of the seismograms the user can apply the respective filters through the `filter` method, which utilizes the frequency filters implemented in the ObsPy package (lowpass, highpass, bandpass and bandstop; Beyr

e.g.:

```
478:4 # Apply bandpass filter
479:5 srm.filter(type='bandpass', freqmin=10, freqmax=120)
```

```
480:1 INFO : Applied bandpass filter (10.0 to 120.0 Hz)
```

In this example, the filter is solely applied to the currently selected traces: Once opened, yet setting the parameter `onhold` as `True` automatically filters all subsequently selected traces with the same filter settings, e.g.:

```
483:6 # Apply lowpass filter
484:7 srm.filter(type='lowpass', freq=200, onhold=True)
```

```
485:1 INFO : Applied 200.0 Hz lowpass filter
```

```
486:2 INFO : Set filter hold on
```

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

```
488 18 # Plot selected traces
489 19 srm.plot()
```

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

496 The active processing mode and data scaling mode are reported together with the traveltimes at the current cursor  
 497 position in the status bar of the seismogram plot window (see Figure 9). The initial processing mode is 'Fb pick',  
 498 i.e., first break picking is possible. The user can switch between the different modes by pressing specific keys on the  
 499 keyboard. Additional modes, accessible through the keyboard, allow for an enhanced visualization of the seismograms,  
 500 e.g., associated to broken geophones or seismograms with wrong polarity. The 'm' key activates the trace mute mode  
 501 ('Trc mute'), which allows to set the amplitude of a trace to zero by clicking with the left mouse button; clicking again  
 502 on the same trace restores the amplitude information. The trace reverse mode ('Trc rev') is activated by pressing the 'r'  
 503 key and enables the user to toggle the polarity of a trace by clicking on it with the left mouse button. The default data  
 504 scaling mode is 'Zoom', which allows the scaling of the y-axis by turning the mouse wheel. By pressing the 'a' key the  
 505 amplitude scaling mode ('Amp scal') is activated. Turning the mouse wheel increases or decreases the amplitudes of  
 506 the traces currently shown in the seismogram plot, and thus might help to enhance the perceptibility of the first onsets.  
 507 By pressing the key of the currently active mode again, the SeismicRefractionManager returns to the default mode;  
 508 yet, the different modes can be activated in any arbitrary order (as illustrated in Figure 9).

509 Through the `plot` method the frequency spectrum of the currently selected trace data can be visualized. A  
 510 frequency spectrum as shown in Figure 7 can be used to discriminate the dominating signal frequencies from those  
 511 associated to the background noise, and thus allows for the definition of adequate filter settings. The frequency  
 512 spectrum illustrates the frequency content of the currently selected traces, which allows for the identification of frequency  
 513 ranges associated to noise, which can be omitted through the application of corresponding frequency filtering. To  
 514 improve the signal-to-noise ratio it is possible to apply filters on the selected traces through the `filter` method, which  
 515 utilizes the frequency filters implemented in the ObsPy package (lowpass, highpass, bandpass and bandstop; Beyreuther et al., 20  
 516 ; e.g.: By default, filters are solely applied to the currently selected traces, yet setting the filter on hold allows the filtering  
 517 of all subsequently selected traces with the same filter settings. In case the seismogram plot is opened, the effect of the  
 518 applied filter on the seismic waveforms is interactively visualized.

519 **2.2.3. Analysis of the seismic waveforms and first break traveltimes picking**

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

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

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

```
53520 # Plot traveltimes diagram
53621 srm.plot(type='traveltimes')
```

537 Figure 10 presents an exemplary traveltimes diagram, which is a common way to examine the quality of the first break  
 538 picking. Such illustration of the traveltimes can be used to identify outliers or erroneous measurements, which are  
 539 commonly associated to traveltimes with substantial deviations from those observed at adjacent stations such as the  
 540 first break pick for the SIN-RIN pair (23, 11). Outliers can be removed by clicking on the ~~corresponding symbol~~  
 541 (~~'x'~~) ~~respective data point~~ in the traveltimes diagram, which is instantly synchronized with the project database. If the  
 542 seismogram plot and the traveltimes diagram are used side-by-side, changes made to the first break picks in one window  
 543 will interactively trigger an update of the other one and vice versa.

544 ~~The SeismicRefractionManager handles first-break picks in picksets, which can be organized and manipulated~~  
 545 ~~through the `picksets` method of the `SeismicRefractionManager`: Calling the `pickset` method without parameters~~  
 546 ~~shows the status of all picksets in the project. From the above use case, we see that by default, the pickset 'picks' is~~  
 547 ~~loaded and activated, i.e., modifications of first breaks are stored in this pickset. First break picks provided by another~~  
 548 ~~source (as udf file) can be imported from `03_proc/picks`. For the first break picking, it is sufficient to keep only one~~  
 549 ~~pickset in the workflow, i. e., not required picksets can be unloaded. A pickset currently not loaded can be deleted~~  
 550 ~~permanently, i.e., the corresponding traveltimes are removed from the database. The `picksets` method can also be~~

551 used to load picksets from the database: When a pickset is loaded from the database, it does not become the active  
 552 pickset automatically; whereas by using the parameter `use` the corresponding pickset is loaded and also becomes the  
 553 active pickset. The first break traveltimes of a pickset can be exported to an udf file that Once the user is satisfied  
 554 with the quality of the first break traveltimes, the data points of the pickset can be exported in the unified data format,  
 555 which is compatible with the pyGIMLi framework. In particular, the udf file is stored in `03_proc/picks` subdirectory  
 556 the subdirectory `03_proc/picks` with the current timestamp as suffix:

```
557 22 # Export first break traveltimes
558 23 srm.picksets(do='export', name='pick')
```

```
559 1 INFO : Pickset 'pick' saved to pick_20230303T140447.pck'
```

### 560 3. Exemplary use cases

#### 561 2.1. Modeling and use Expanding the capabilities of synthetic seismic waveform data the formikoj

##### 562 library

563 To demonstrate the applicability of the `SeismicWaveformModeler` class, we generate two synthetic seismic  
 564 waveform datasets assuming two horizontal layers in a half-space without topography and an interface between these  
 565 layers with a constant depth. Table ?? summarizes the parameterization used for the forward modeling of one dataset  
 566 without noise (Dataset\_1) and another dataset containing random noise and systematic errors (Dataset\_2). For the  
 567 visualization of both synthetic datasets, we provide the corresponding geometry files to the `SeismicRefractionManager`  
 568 in order to have full control regarding data selection and processing capabilities.

569 The synthetic seismic waveforms for SIN-24 of Dataset\_1 presented in Figure ?? reveal clear negative first onsets  
 570 and allow the identification of crossover points, i.e., inflection points between the first arrivals from the first layer  
 571 (RIN 17 to 30) and the first onsets associated to the second layer (RIN 1 to 16 and RIN 31 to 48). In contrast the  
 572 signal-to-noise ratio of Dataset\_2 (see Figure ??) is a function of the offset between shot and geophone position, i.e.,  
 573 traces farther away from the shot contain a higher level of random noise. Moreover, Dataset\_2 also contains systematic  
 574 errors, where RIN 6 and 14 refer to broken geophones, and RIN 35 is an example for readings with a wrong polarity.  
 575 We believe that such synthetic datasets might be useful for investigating the effect of complex survey geometries on  
 576 seismic data as well as the development and evaluation of processing strategies. Measurement scheme and parameters  
 577 provided in the yaml files used to create synthetic seismic waveform datasets with (Dataset\_1) and without added noise  
 578 (Dataset\_2). **Measurement scheme** Number of stations 48 Station spacing 2 m Number of geophones 48 Number of  
 579 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**  
 580 **False True Noise level** 0 1e-5 **Missing shots** False True **Broken geophones** False True **Wrong polarity** False True

581 ~~Wavelet Length 1.024 s Frequency 100 Hz Sampling rate 2000 Hz Synthetic seismic waveform data without random  
582 or systematic errors created with the SeismicWaveformModeler() class for a shot position in the center of the survey  
583 layout.~~

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

587 ~~Accordingly, the concept of the formikoj library allows the implementation of supplementary functionalities. Such  
588 Making the formikoj library publicly available under an open-source license allows the addition of supplementary  
589 functionalities tailored for the specific requirements of the users. The concept of formikoj suggest that such custom  
590 extensions should be implemented either as internal methods or as functions in the utilities module, which are then  
591 executed through the compute method with a custom keyword. As an illustrative example, we implemented~~

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

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

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

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

598 ~~and~~

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

599 ~~where  $n_{STA}$  and  $n_{LTA}$  refer to the number of data points in  $E(t)$  considered for the short- and long-time windows,  
600 respectively. For this exemplary implementation we determine the first break traveltimes in the seismograms as the~~

601 position of the maximum in the *ER* function, which is automatically saved in the project database.

602 The autopicking algorithm is added to the SeismicRefractionManager in form of two internal methods

603 \_manage\_autopicking and \_compute\_autopicks, respectively. The autopicking process can be started by passing

604 the new keyword To invoke the autopicking process we modified the compute method, which now accepts the custom-defined

605 keyword autopickautopicking as the first parameter to the . Additionally, the autopicking requires values to be

606 passed for the parameters computepick method: The second parameter defines whether traveltimes should be determined

607 for the entire dataset (and pickset. The former accepts the values all ) or solely or cur to determine first break

608 traveltimes for all traces in the dataset or solely the currently selected traces(eur). The third parameter provides

609 respectively. The latter defines the name of the corresponding pickset in the project database; thus, making the

610 traveltimes available for visualization and processing with existing functionalities or further custom implementations.

611 picksets in which the automatically determined traveltimes should be saved to. A typical use case involves conducting

612 the autopicking and exporting the determined traveltimes:

```
613 24 # Automatically pick first break traveltimes
614 25 srm.compute(do='autopicking', pick='all', pickset='autopicks')
615 26 srm.picksets(do='export', name='autopicks')
```

```
616 1 INFO : Created new pickset 'autopicks'
617 2 INFO : Pickset 'autopicks' loaded
618 3 INFO : 'autopicks' set at active pickset
619 4 Progress <=====> 100.0% completed
620 5 INFO : Pickset 'autopicks' saved to autopicks_20230303T140834.pck
```

## 621 2.2. First break travel time picking for a 2D roll-along field data set: the Danube island example

622 The seismic data used in this example were collected at the Danube island (Vienna) in June 2021, using 48

623 geophones deployed with 2 m spacing between them and shot locations located between the geophone positions.

624 As illustrated In Figure 11, we compare the automatically determined first break picks with the forward modeled

625 traveltimes computed above to allow for a basic evaluation of autopicking performance. The inset plot in Figure 11

626 presents the histogram of the autopicked traveltimes, which shows that three traveltimes should be considered outliers,

627 and thus are removed from the dataset. After the outlier removal the correlation coefficient between forward modeled

628 and autopicked traveltimes is 0.99 suggesting that the implemented autopicking algorithm performs well for the synthetic

629 seismic waveform data. However, the observed deviation from the perfect correlation (i.e., the 1:1 line in Figure ??, 11)

630 indicates that autopicking and forward modeling algorithm might be sensitive to different seismic wave phases. Further

631 improvements in terms of the autopicking process could incorporate, for example, machine learning approaches as the

632 method proposed by Duan and Zhang (2020), which can be easily implemented as an additional functionality in the

633 SeismicRefractionManager. We point out here, that following the same procedure, the user can implement further  
 634 autopicking algorithms as well as other data processing strategies and have a simple framework for the evaluation of  
 635 their performance through the analysis of synthetic data (e.g., clean and contaminated with Gaussian noise).

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

637 To demonstrate the applicability of the SeismicRefractionManager for the processing of real field data, we  
 638 present here the analysis and inversion results for a seismic refraction survey conducted in a soda lake located  
 639 in the Nationalpark Neusiedler See–Seewinkel close to Vienna. The seismic survey aims at solving for the geometry  
 640 of the confined aquifer below this soda lake with the required information referring to the survey geometry refers to  
 641 a roll-along geometry, i. e., thickness of the aquifer and the geophone spread was moved along a profile with 50%  
 642 overlap yielding a total of five segments. The objective of the survey was to define the contact between different  
 643 sediments within the tertiary and quaternary deposits used to build the man-made Danube island. Additionally, the  
 644 survey aimed to identify lateral changes that might indicate the position of a fault, which has been inferred from  
 645 sediments recovered from drillings. depth of the groundwater table. As presented in Figure 12, the seismic data were  
 646 collected with 48 geophones deployed along the North-East to South-West oriented line, and 48 geophones along the  
 647 North-West to South-East oriented line, with a spacing of 2 m between the geophones. Shots were generated with an  
 648 8 kg sledgehammer at the geophone positions as well as at positions along the diagonals to obtain a sufficient coverage.

649

650 As in case of the synthetic dataset presented above, the geometry file contains the coordinates of the survey stations,  
 651 yet for the soda lake dataset 3D coordinates we provided as illustrated in Table 5. Since geophones were not deployed  
 652 at each survey station the column **Geo** contains the value 1 (True) for the first 96 stations and 0 (False) for the remaining  
 653 20 stations.

654 In the field, each segment was measured separately, yet for the processing all measurements are combined in a single  
 655 profile. We can implement such measurement layout in a geometry file as illustrated in Table ???. The key parameter is  
 656 '1st Geo' as it indicates the first active geophone along the profile for each shot file. For the first segment in a roll-along  
 657 survey geometry, the first active geophone is always geophone 1. Considering the number of geophones used in each  
 658 segment, and an overlap of 50%, the first geophone for segments two to five are 25, 49, 73 and 97, respectively.

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

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

**Table 5**

~~Extract Excerpt from the roll-along 3D 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
0.0 40336.056	0.0 292470.692	163.5 120.0	1	-1 1 1 2 0 0 0 163.5 0 1001
+40334.751	+292469.177	+120.0	+1	+1002
98.0 0.0 163.5 0 1037 25 48	:	:	:	:
194.0 40284.472	0.0 292455.431	163.5 120.0	0 1	1049 1058
196.0 40285.896	0.0 292454.027	163.5 120.0	1	-1 1059
202.0 0.0 163.5 0 1050 25 48	:	:	:	:
286.0 40339.421	0.0 292402.681	163.5 120.0	0 1	1084 1096
290.0 40305.228	0.0 292445.050	163.5 120.0	0	1097
:	:	:	:	:
386.0 40265.415	0.0 292431.957	163.5 120.0	0	1109 1115
390.0 40255.453	0.0 292431.395	163.5 120.0	0	2025 1116 97 48 ::::: 570.0

664 offset stack (COS): The COS for the Danube island dataset presented in Figure ?? shows first onsets for absolute  
 665 offsets up to approximately 150 m, which suggest a two-layered subsurface model. By using the velocity estimation  
 666 functionality we can approximate the seismic velocity in the corresponding layers.

667 For the first break picking a set of traces is selected, filtered if necessary and visualized. In this example, the  
 668 trace data for SIN 99 are used, which refers to a shot position located in segment four. As can be seen in Figure ??,  
 669 the first onsets are easily identifiable despite the substantial seismic background noise at large offsets (RIN 73 to  
 670 90). Exemplary seismic waveforms from the Danube island dataset shown for shot index number (SIN) 99 with a  
 671 120 Hz lowpass filter applied to suppress high frequency noise. The receiver index numbers (RIN) start at 73, thus  
 672 indicating that the data were collected with a roll-along survey geometry. A pseudosection provides an illustration of  
 673 the corresponding apparent seismic velocity values computed from the picked traveltimes and the distance between  
 674 shots and geophones:

675 In a pseudosection, as presented in Figure ?? for the Danube island dataset, the apparent velocity values are  
 676 assigned to pseudolocations, with the corresponding x- and z-coordinates determined as the midpoint and as 1/3  
 677 of the absolute offset between the shot and geophone, respectively. Accordingly, such plot allows for the identification  
 678 of outliers in the data, e.g., stark velocity contrasts for adjacent points, or systematic errors, e.g., velocities erroneously  
 679 influenced by a single shot or receiver. The main assumption here is that the pseudosection should reveal smooth  
 680 transitions between lateral and vertical neighbors, considering that the data were collected with gradual changes in  
 681 the position of the source (i.e., hammer blow) and the receiver (i.e., geophone). The pseudosection will reveal large  
 682 variations in case of abrupt changes in the topography or the geometry of the array, yet this can be taken into account  
 683 by the user during the identification of outliers and possible systematic errors. For the Danube island dataset, the  
 684 pseudosection suggests an increase in the seismic velocity along profile direction in deeper subsurface regions (i.e.,  
 685 larger pseudodepth). Such pattern could be related to a fault expected in this area of the Danube island, thus indicating

686 that the geophysical survey was sufficiently designed to detect such feature. Moreover, the pseudosection presented  
 687 in Figure ?? shows a low number of data points in the first segment of the Danube island survey. This lack of data  
 688 points at large pseudodepths is due to a low number of picked traveltimes at large offsets, and thus might indicate a low  
 689 signal-to-noise ratio. Pseudosection showing the apparent seismic velocities determined from the first break traveltimes  
 690 obtained from the Danube island dataset and the corresponding absolute offset between the shot and receiver stations.  
 691 The apparent velocity for each shot-receiver pair is illustrated at the corresponding midpoint and pseudodepth (1/3 of  
 692 the absolute offset).

693 To review the data quality along the entire profile it is possible to visualize the picking percentage, i.e., the ratio  
 694 of actually picked traveltimes and total number of SIN-RIN pairs: Figure ?? shows that the picking percentage is  
 695 visualized separately for each SIN. In this way, a low picking percentage indicates shots affected by a low signal-to-noise  
 696 ratio. For the Danube island dataset, the picking percentage is low in the first segment, which confirms the lack of  
 697 datapoints observed in the pseudosection. Moreover, the picking percentage plot can be used to track the picking  
 698 progress, for instance if the traveltimes cannot be determined in one session or to identify single shots that might have  
 699 been forgotten during the first break picking. Accordingly, it is advisable to check the picking percentage prior to  
 700 exporting the traveltimes for the inversion.

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

### 704 3.1. Processing of a 3D seismic refraction dataset: the soda lake example

705 The SeismicRefractionManager allows also the visualization and processing of data collected with a, automatically,  
 706 infers the 3D survey layout . To illustrate these capabilities, we present in Figure 12 the geometry of a from the 3D  
 707 survey conducted in a soda lake located close to Vienna. The soda lake corresponds to quaternary sediments where  
 708 capillary forces have developed a low permeable layer close to the surface (between 50 geometry, and 100 cm) with a  
 709 high clay and salt content. The seismic survey aims to support the interpretation of the electrical and electromagnetic  
 710 models obtained in a monitoring framework. Accordingly, the survey geometry shown in Figure 12 was specified by  
 711 previously conducted electrical measurements with electrodes arranged along two perpendicular lines. The seismic  
 712 data were collected with 48 geophones deployed along the North-East to South-West oriented line, and 48 geophones  
 713 along the North-West to South-East oriented line, with a spacing of 2 m between the geophones. Shots were generated  
 714 with an 8 kg sledgehammer at the geophone positions as well as at positions along the diagonals.

715 The soda lakes field dataset was collected in a 3D survey layout with stations (co-located geophones and shots  
 716 indicated by filled dots) deployed in form of a cross. Additional shots (yellow stars) were conducted in front of a cross

717 rotated by 45° to increase the coverage of the dataset.

718 The geometry file contains the 3D coordinates for each shot or receiver station, and thus the SeismicRefractionManager  
 719 automatically thus configures the project for 3D processing. Figure 13 presents the Then we can select seismograms  
 720 the same way as for the synthetic data presented above. For this example we select seismic waveform data recorded  
 721 for SIN 1, i.e., the shot position co-located with the first geophone (Station 01 in Figure 12). The:

```
722:0 # Select traces with receiver
723:1 srm.select(by='sin', num=1)
```

724:1 INFO : 96 traces selected

725 In Figure 13, we can see that the data for RIN 1 to 48 appear familiar as the corresponding SIN-RIN layout refers  
 726 to the one of conventional 2D profiles; whereas the seismic waveforms for RIN 49 to 96 show an entirely different  
 727 pattern. To understand such visualization, we need to take into account that RIN 49 to 96 are deployed perpendicular  
 728 to the direction of propagation of the wavefront originating from SIN 1. Accordingly, the observed curvature in the  
 729 first onsets is due to the varying offset of the different SIN-RIN pairs.

730 Due to the 3D survey geometry, a 2D pseudosection is not suitable to illustrate the apparent velocity values obtained  
 731 from the picked for assessing the quality of the first break traveltimes. HenceHowever, the SeismicRefractionManager  
 732 automatically switches to a project is configured for 3D representation where processing, and thus the apparent velocity  
 733 values are illustrated in an interactive 3D pseudosection. This plot can be rotated and the image section can be zoomed  
 734 and panned allowing the user to easily investigate the data quality for 3D geometries. Figure 14 shows a screenshot of  
 735 the 3D pseudosection for the salt lake dataset; yet, such screenshot cannot reveal the full capabilities implemented in  
 736 the SeismicRefractionManger for the interactive analysis and visualization of 3D pseudosections.

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

```
739:0 # Plot the picking percentage
740:1 srm.plot(type='pickperc')
```

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

748 Once the first break picking is finished, the corresponding pickset can be exported for the inversion. ~~The inversion~~  
 749 ~~results and their interpretation are not~~ We inverted the first break traveltimes with pyGIMLi and present the resolved  
 750 3D subsurface model in Figure 16a. From this representation we can see, that the inversion solves for low seismic  
 751 velocities (600 to 800 ms<sup>-1</sup>) in the near-surface in the center of the investigated area, which corresponds to the still  
 752 intact part of the soda lake, i.e., the part that is covered by water on a seasonal basis. Seismic velocity values above  
 753 800 ms<sup>-1</sup> are resolved at depth as well as outside of the soda lake. To facilitate the interpretation of the resolved  
 754 seismic velocity distribution in terms of the aquifer geometry we show the two vertical slices in Figure 16b and c,  
 755 respectively. In particular, we relate seismic velocity values > 800 ms<sup>-1</sup> to the transition from the unsaturated to the  
 756 saturated zone due to the higher seismic velocity of water ( $\approx$  1500 ms<sup>-1</sup>) compared to air ( $\approx$  340 ms<sup>-1</sup>) filling the pore  
 757 space. Accordingly, our 3D subsurface model indicates the depth of the groundwater table at a depth of approximately  
 758 8 m. ~~A more detailed interpretation is beyond~~ the scope of this manuscript, yet ~~Figures 13 and 14~~ the presented figures  
 759 reveal the capabilities provided by the proposed framework for the visualization and processing of data collected in  
 760 3D survey geometries.

## 761 4. Conclusions and Outlook

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

766 We demonstrated the capabilities of the formikoj framework to develop versatile and easily scalable classes for the  
 767 modeling and processing of waveforms in seismic refraction surveys. ~~The required interaction with the user is reduced~~  
 768 ~~to a minimum as~~ By standardizing the data input in combination with a thorough sanity check aims at reducing the  
 769 risk of corrupting the information stored in the project database. Moreover, crucial processing steps are automatized  
 770 within the SeismicWaveformModeler and the  
 771 SeismicRefractionManager ~~based on~~ classes facilitated by efficient data input strategies, for instance the prepara-  
 772 tion and import of the geometry file or the keyboard-based interaction related to the first break picking. In this regard,  
 773 the user controls the formikoj library by providing text-based commands preferably through an ipython shell to exploit  
 774 the full interactive potential modeling and processing tools. However, applications of the formikoj library can also be  
 775 automatized by implementing workflows in python scripts or jupyter notebooks.

776 Based on three exemplary use cases, we illustrated the applicability of both the SeismicWaveformModeler  
 777 and the SeismicRefractionManager. In the first use case, we showed the possibility to forward model seismic  
 778 waveform data based on custom subsurface models and survey geometries with the SeismicWaveformModeler.

779 Additionally, the resulting waveforms can be subjected to systematic and random noise sources. The capabilities of the  
 780 SeismicRefractionManager were demonstrated through the processing of field datasets collected in complex survey  
 781 layout, namely a roll-along and a 3D geometry. Moreover, we showed how the different data visualization options can  
 782 assist during the data processing to ensure consistency in the first break traveltimes. In particular, we developed a  
 783 visualization of the traveltimes by means of pseuodsections illustrating the corresponding apparent seismic velocities.  
 784 Such plots allow for a quick identification of systematic errors and outliers in both 2D and 3D datasets.

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

## 789 5. Acknowledgments

790 The authors are grateful to Nathalie Roser and Lukas Aigner for benchmarking using and testing the formikoj library  
 791 against established software packages in frame of their research activities, and for providing valuable suggestions for  
 792 improvement. Furthermore, we would like to thank Clemens Moser, Martin Mayr, Vinzenz Schichl and Harald Pammer  
 793 for their constructive comments during first tests of the formikoj framework as well as for their help during the seismic  
 794 surveys.

795    **Code and data availability section**

796    Name of the code/library: formikoj

797    Contact: matthias.steiner@geo.tuwien.ac.at

798    Hardware requirements: No specific requirements

799    Program language: Python

800    Software required: Anaconda/Miniconda recommended

801    Total program and dataset size: 250 MB

802    The source codes and exemplary ~~data-sets~~datasets are available for downloading at the link:803    <https://git.geo.tuwien.ac.at/msteiner/formikoj.git>804    The ~~orthophotos used in Figures ?? and 12 were~~orthophoto used in Figure 12 was published by geoland.at under

805    a CC BY 4.0 license.

806    **A. Source code for generating the subsurface model considered in this study**

```

807 1 # Import required packages
808 2 import numpy as np
809 3 import pygimli as pg
810 4 import pygimli.meshTools as mt
811 5 %DIF >
812 6 # Create the model geometry
813 7 # - top layer
814 8 top = mt.createPolygon([[0, 0], [94, 0],
815 9 [94, -3.5], [72, -3.5],
816 10 [20, -2], [0, -2]],
817 11 isClosed=True, marker=1, area=0.1)
818 12 %DIF >
819 13 # - bottom layer
820 14 bottom = mt.createPolygon([[0, -2], [20, -2],
821 15 [22, -6], [70, -6],
822 16 [72, -3.5], [94, -3.5],
823 17 [94, -10], [0, -10]],
824 18 isClosed=True, marker=3, area=0.1)
825 19 %DIF >
826 20 # - anomaly/infill
827 21 infill = mt.createPolygon([[20, -2], [72, -3.5],
828 22 [70, -6], [22, -6]],
829 23 isClosed=True, marker=2, area=0.1)
830 24 %DIF >
```

```

83125 geom = top + infill + bottom
83226 %DIF >
83327 # Define shot and receiver stations and create corresponding nodes
83428 nstats = 48
83529 spc = 2
83630 %DIF >
83731 stations = np.vstack((np.linspace(0, (nstats-1)*spc, nstats),
83832 np.zeros(nstats))).T
83933 %DIF >
84034 for p in stations:
84135     geom.createNode(p)
84236 %DIF >
84337 # Create mesh for the finite element modeling
84438 mesh = mt.createMesh(geom, quality=34)
84539 %DIF >
84640 # Save the mesh in the binary mesh format for later use
84741 # with the SeismicWaveformModeler
84842 mesh.save('mesh.bms')

```

## 849 References

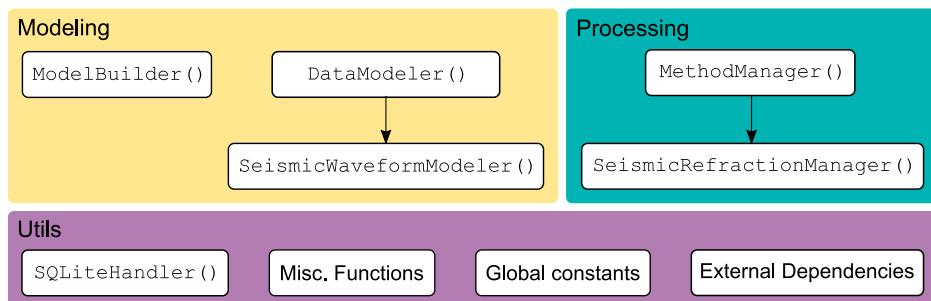
- 850 Ahern, T., Casey, R., Barnes, D., Benson, R., Knight, T., Trabant, C., 2012. SEED Reference Manual, Version 2.4. URL: [http://www.fdsn.org/pdf/SEEDManual\\_V2.4.pdf](http://www.fdsn.org/pdf/SEEDManual_V2.4.pdf).
- 851 Allen, R.V., 1978. Automatic earthquake recognition and timing from single traces. Bulletin of the seismological society of America 68, 1521–1532. doi:10.1785/bssa0680051521.
- 852 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.
- 853 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>.
- 854 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.
- 855 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>.
- 856 Dijkstra, E.W., 1959. A note on two problems in connexion with graphs. Numerische Mathematik , 269–271doi:10.1007/bf01386390.
- 857 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>.

- 868 Duan, X., Zhang, J., 2020. Multitrace first-break picking using an integrated seismic and machine learning methodpicking based on machine  
869 learning. *Geophysics* 85, WA269–WA277. doi:10.1190/GEO2019-0422.1.
- 870 Earle, P.S., Shearer, P.M., 1994. Characterization of global seismograms using an automatic-picking algorithm. *Bulletin of the Seismological  
871 Society of America* 84, 366–376.
- 872 Guedes, V.J.C.B., Maciel, S.T.R., Rocha, M.P., 2022. Refrapy: A python program for seismic refraction data analysis. *Computers & Geo-  
873 sciences* 159, 105020. URL: <https://www.sciencedirect.com/science/article/pii/S0098300421003022>, doi:<https://doi.org/10.1016/j.cageo.2021.105020>.
- 875 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,  
876 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,  
877 K., Reddy, T., Weckesser, W., Abbasi, H., Gohlke, C., Oliphant, T.E., 2020. Array programming with NumPy. *Nature* 585, 357–362. URL:  
878 <https://doi.org/10.1038/s41586-020-2649-2>, doi:10.1038/s41586-020-2649-2.
- 879 Heimann, S., Kriegerowski, M., Isken, M., Cesca, S., Daout, S., Grigoli, F., Juretzek, C., Megies, T., Nooshiri, N., Steinberg, A., Sudhaus, H.,  
880 Vasyura-Bathke, H., Willey, T., Dahm, T., 2017. Pyrocko - an open-source seismology toolbox and library doi:10.5880/GFZ.2.1.2017.001.
- 881 Hunter, J.D., 2007. Matplotlib: A 2d graphics environment. *Computing in Science & Engineering* 9, 90–95. doi:10.1109/MCSE.2007.55.
- 882 McKinney, W., 2010. Data Structures for Statistical Computing in Python, in: Stéfan van der Walt, Jarrod Millman (Eds.), *Proceedings of the 9th  
883 Python in Science Conference*, pp. 56 – 61. doi:10.25080/Majora-92bf1922-00a.
- 884 Nguyen, F., Ghose, R., Isunza Manrique, I., Robert, T., Dumont, G., 2018. Managing past landfills for future site development: A review of the  
885 contribution of geophysical methods, in: *Proceedings of the 4th International Symposium on Enhanced Landfill Mining*, pp. 27–36.
- 886 Parsekian, A., Singha, K., Minsley, B.J., Holbrook, W.S., Slater, L., 2015. Multiscale geophysical imaging of the critical zone. *Reviews of  
887 Geophysics* 53, 1–26. doi:10.1002/2014RG000465.
- 888 Plattner, A.M., 2020. Gprpy: Open-source ground-penetrating radar processing and visualization software. *The Leading Edge* 39, 332–337.  
889 doi:10.1190/tle39050332.1.
- 890 Ringler, A.T., Evans, J.R., 2015. A quick seed tutorial. *Seismological Research Letters* 86, 1717–1725. doi:10.1785/0220150043.
- 891 Romero-Ruiz, A., Linde, N., Keller, T., Or, D., 2018. A review of geophysical methods for soil structure characterization. *Reviews of Geophysics*  
892 56, 672–697. doi:10.1029/2018RG000611.
- 893 Rücker, C., Günther, T., Wagner, F.M., 2017. pygimli: An open-source library for modelling and inversion in geophysics. *Computers & Geosciences*  
894 109, 106–123. doi:10.1016/j.cageo.2017.07.011.
- 895 Steiner, M., Katona, T., Fellner, J., Flores Orozco, A., 2022. Quantitative water content estimation in landfills through joint inversion of seismic re-  
896 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>.
- 898 Steiner, M., Wagner, F.M., Maierhofer, T., Schöner, W., Flores Orozco, A., 2021. Improved estimation of ice and water contents in alpine permafrost  
899 through constrained petrophysical joint inversion: The hoher sonnblick case study. *Geophysics* 86, 1–84. doi:10.1190/geo2020-0592.1.
- 900 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.
- 902 Sullivan, C.B., Kaszynski, A., 2019. PyVista: 3d plotting and mesh analysis through a streamlined interface for the visualization toolkit (VTK).  
903 Journal of Open Source Software 4, 1450. URL: <https://doi.org/10.21105/joss.01450>, doi:10.21105/joss.01450.
- 904 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*,  
905 pp. 96–103.

## 906 List of Figures

907 1	Fundamental use case illustrating the modeling of seismic refraction data within the formikoj ecosystem.	33
908 2	Fundamental use case illustrating the processing of seismic refraction data within the formikoj ecosystem.	34
909 3	Typical formikoj workflow for the picking of first break traveltimes from seismic waveform data . . . . .	35
910 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 traveltimes based on the survey geometry. (d) Subsurface model of the seismic P-wave velocity distribution obtained through the inversion of the synthetic P-wave traveltimes. . . . .	36
919 5	Entity-relationship diagram illustrating the SQLite database used in a SeismicRefractionManager project to store an manage processing related information. . . . .	37
920 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. . . . .	38
925 7	The frequency spectrum illustrates the frequency content of the currently selected traces, which allows for the identification of frequency ranges associated to noise allowing for the definition of corresponding frequency filter settings. . . . .	39
928 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. . . . .	40
933 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. . . . .	41
936 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). . . . .	42
941 11	Comparison of forward modeled synthetic and automatically picked first break traveltimes for the synthetic seismic waveform data created in this study. . . . .	43
943 12	The Danube island soda lakes field dataset was collected along in a single line with a roll-along 3D survey layout ; the with stations (co-located geophones and shots indicated by filled triangles indicate the direction dots) deployed in form of the measurements a cross. The five segments have an overlap Additional shots (yellow stars) were conducted in front of 50% a cross rotated by 45° to ensure an adequate data increase the coverage along of the entire profile dataset. . . . .	44
948 13	Examplary 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. . . . .	45
952 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. . . . .	46
956 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. . . . .	47

958	16	Subsurface model of the soda lake in terms of the seismic P-wave velocity. (a) 3D representation of orthogonal slices through the resolved subsurface model. (b) 2D representation of the NE–SW slice. (c) 2D representation of the NW–SE slice. . . . .	48
959			
960			



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

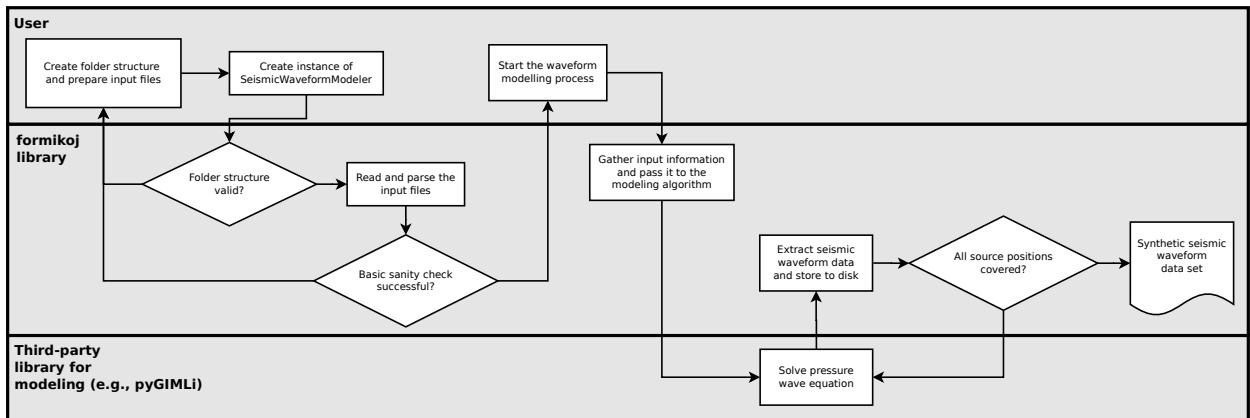


Figure 2: Fundamental use case illustrating the processing of seismic refraction data within the formikoj ecosystem.

## formikoj - Flexible geophysical data processing

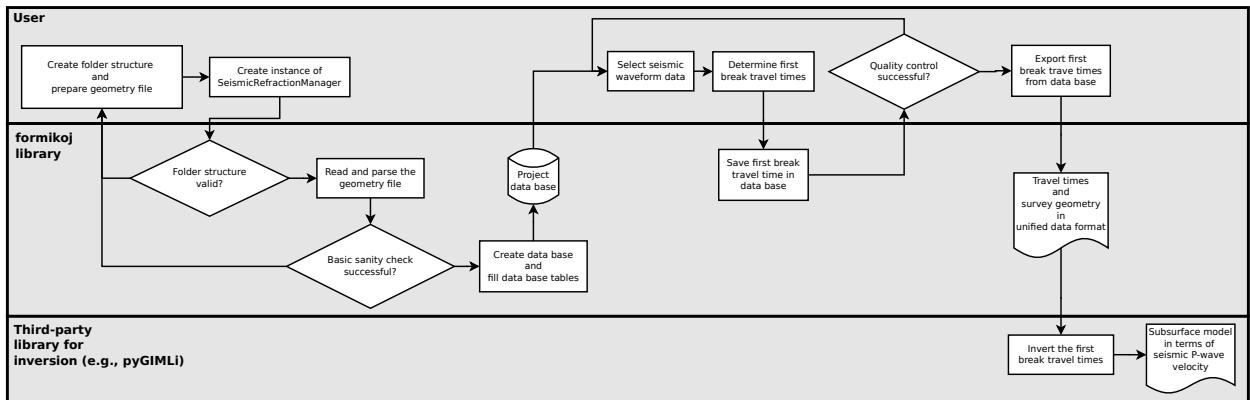
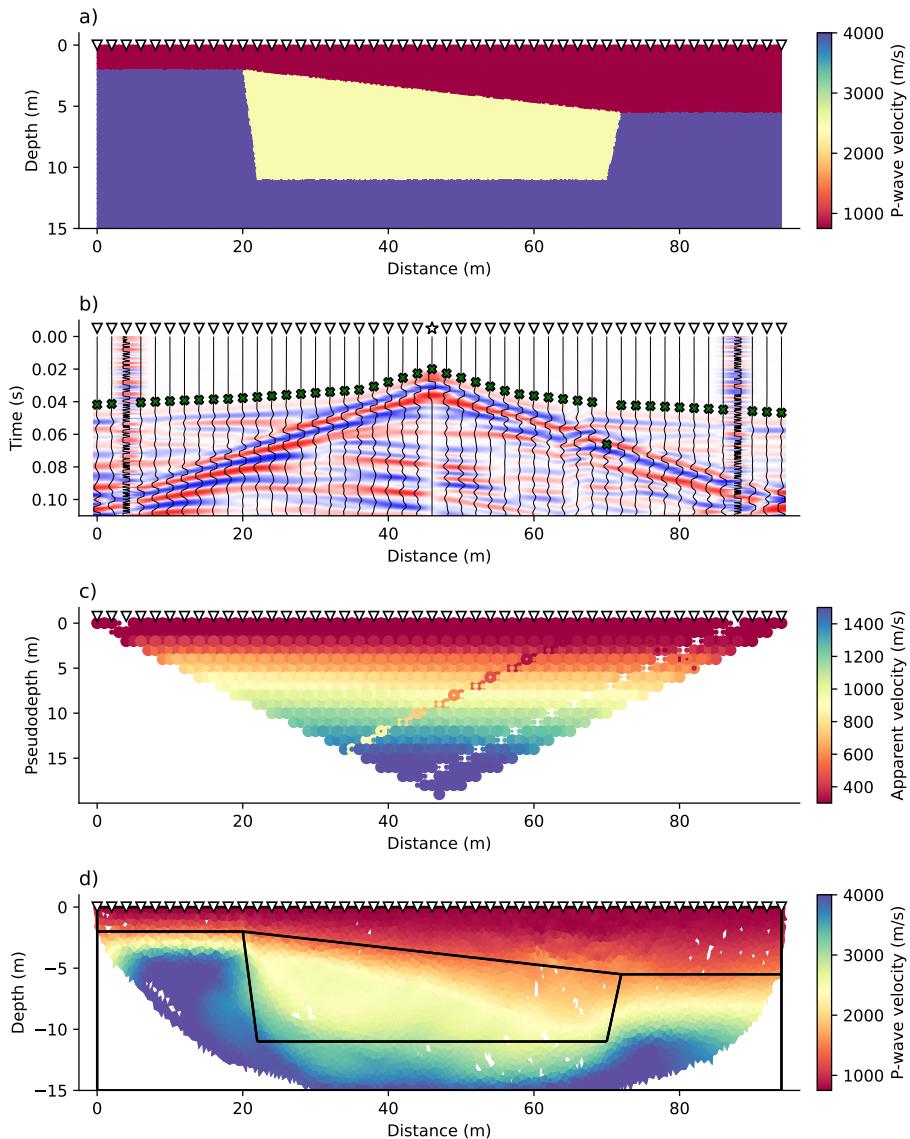


Figure 3: Typical formikoj workflow for the picking of first break traveltimes from seismic waveform data.



**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 traveltimes based on the survey geometry.
- (d) Subsurface model of the seismic P-wave velocity distribution obtained through the inversion of the synthetic P-wave traveltimes.

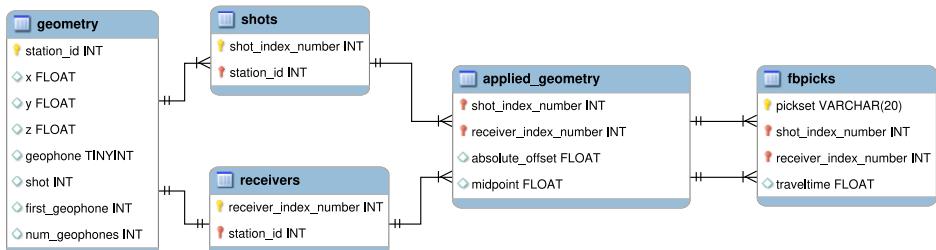
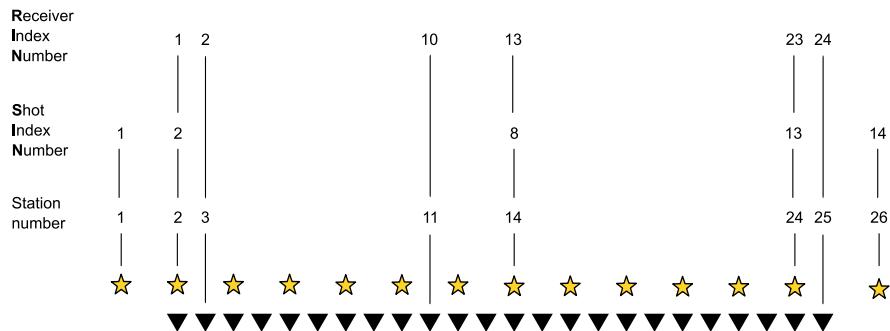
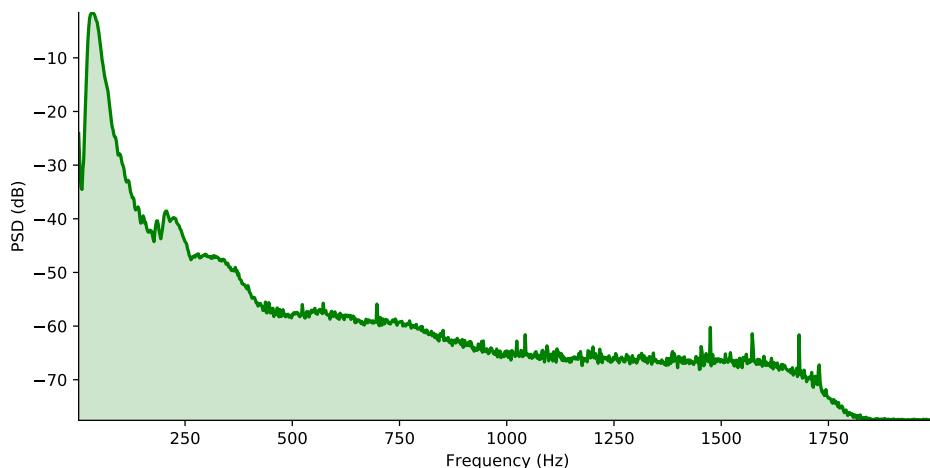


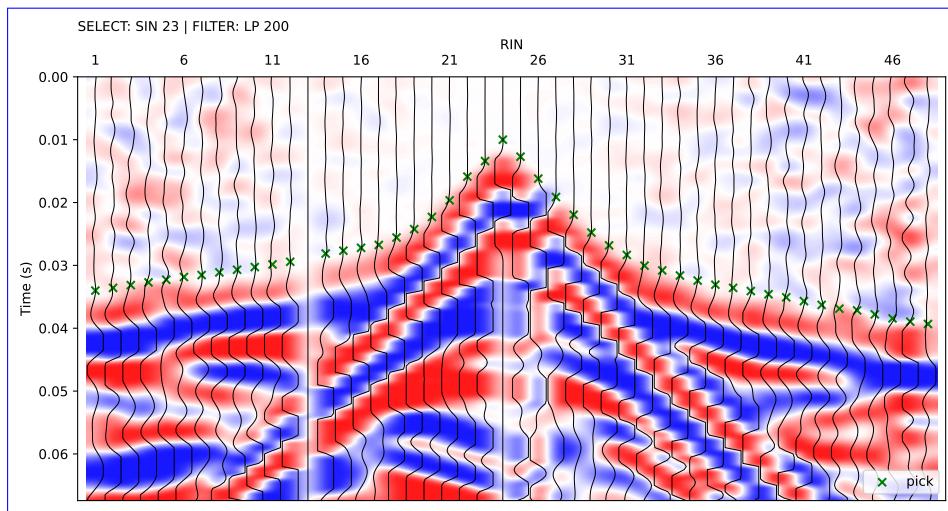
Figure 5: Entity-relationship diagram illustrating the SQLite database used in a SeismicRefractionManager project to store and manage processing related information.



**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 allowing for the definition of corresponding frequency filter settings.

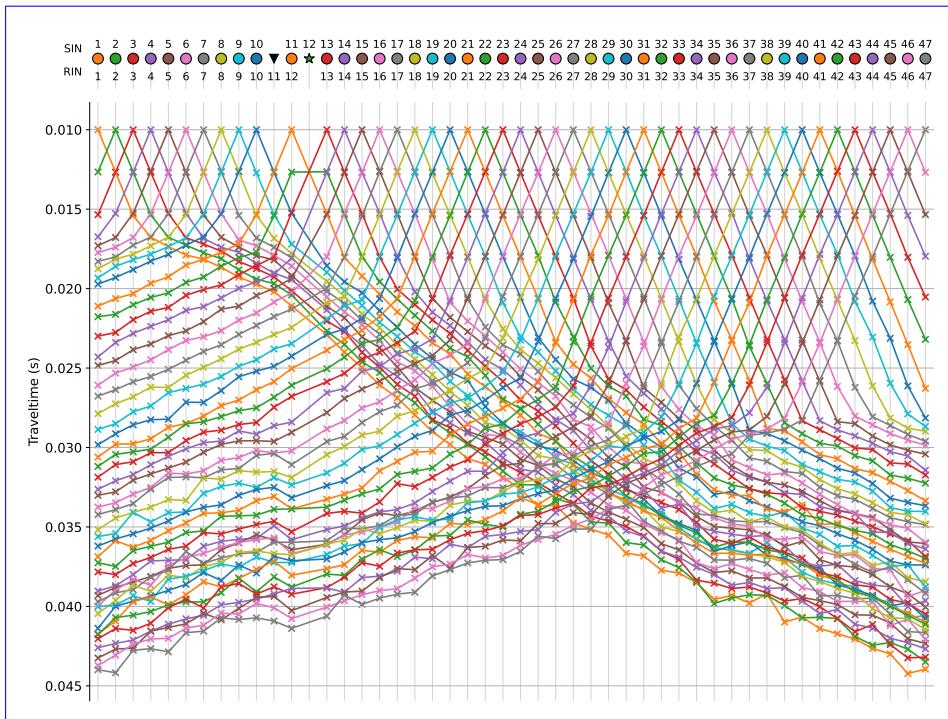


**Figure 8:** The seismogram plot presents the currently selected traces along the x axis, where the sort order is determined through the geometry. The corresponding trace data are illustrated as a function of time along the y axis (solid curves), with positive and negative amplitudes depicted in blue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

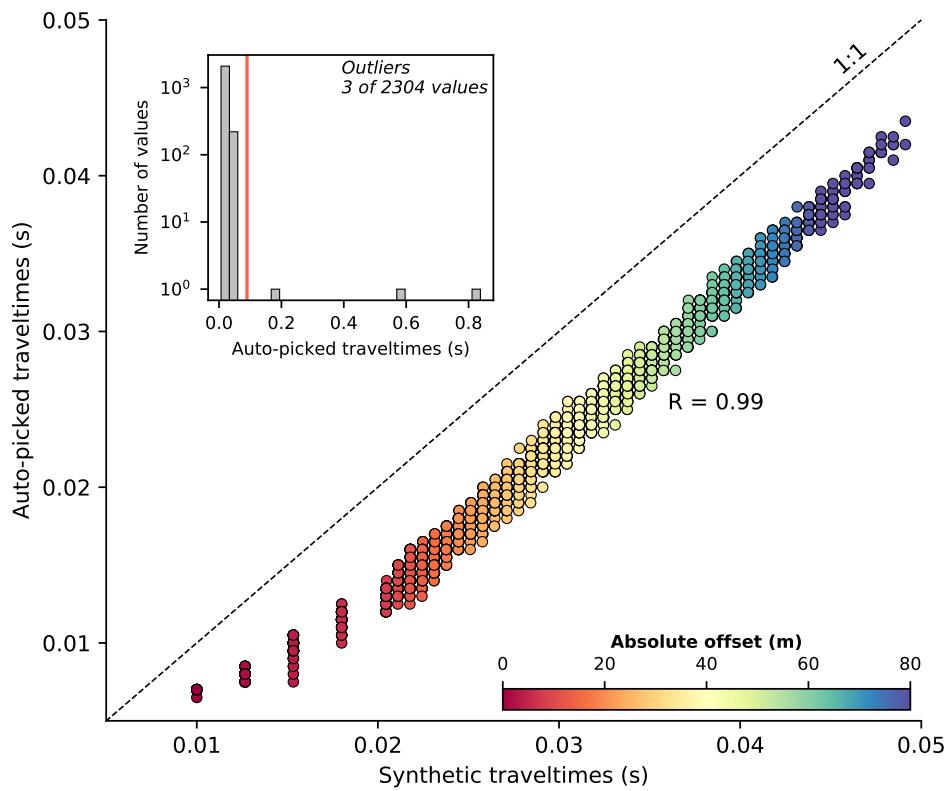
Proc: Fb pick | Scroll: Zoom | Time (s) = 0.000



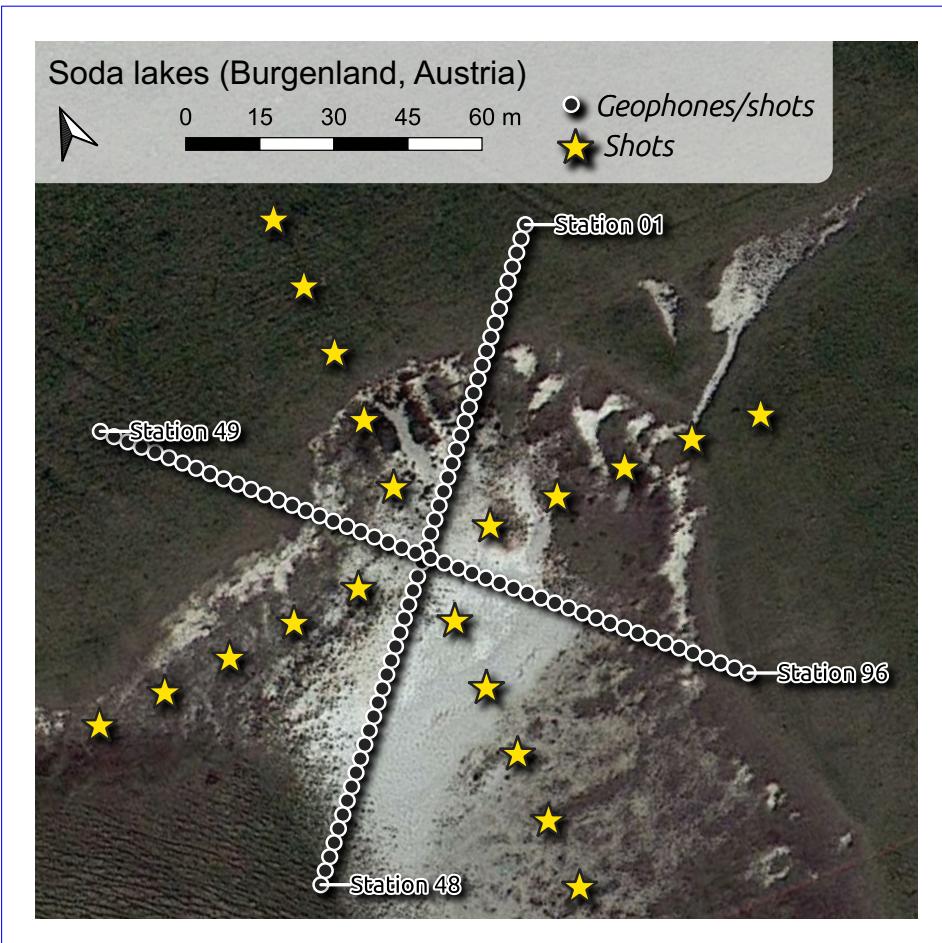
**Figure 9:** The status bar in the interactive seismogram plot displays the currently active processing and data scaling modes as well as the time (in seconds) at the current cursor position. By pressing the keys 'a', 'm', 'r', 'v' on the keyboard the different modes can be activated.



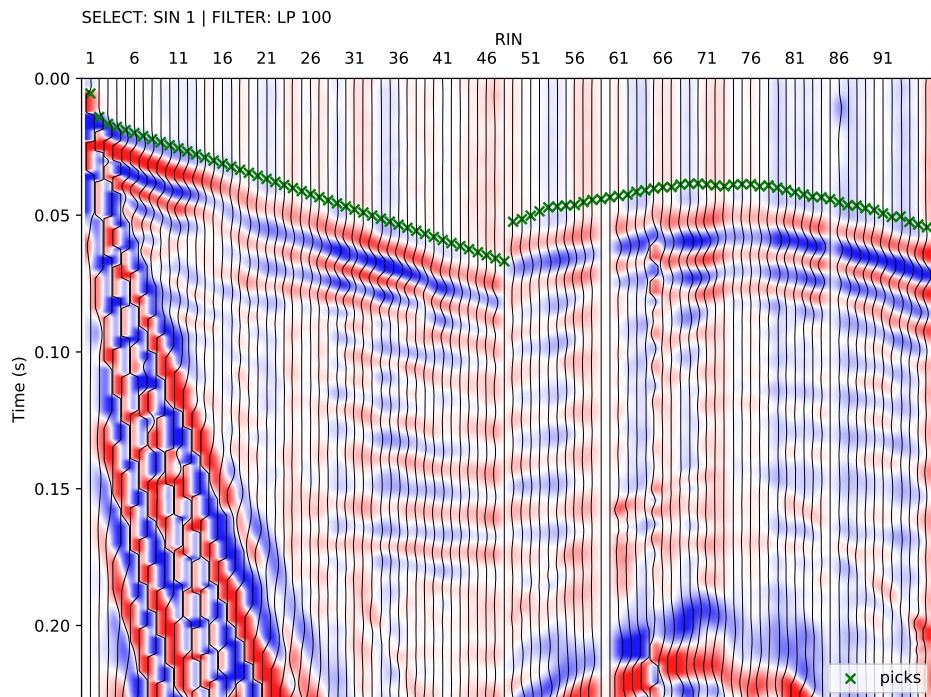
**Figure 10:** The traveltime diagram shows the first break traveltimes stored in the currently active pickset along the y axis (x symbols), where solid lines connect traveltimes assigned to a common shot. The sort order of the stations along the x axis reflects the geometry. Filled circles indicate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



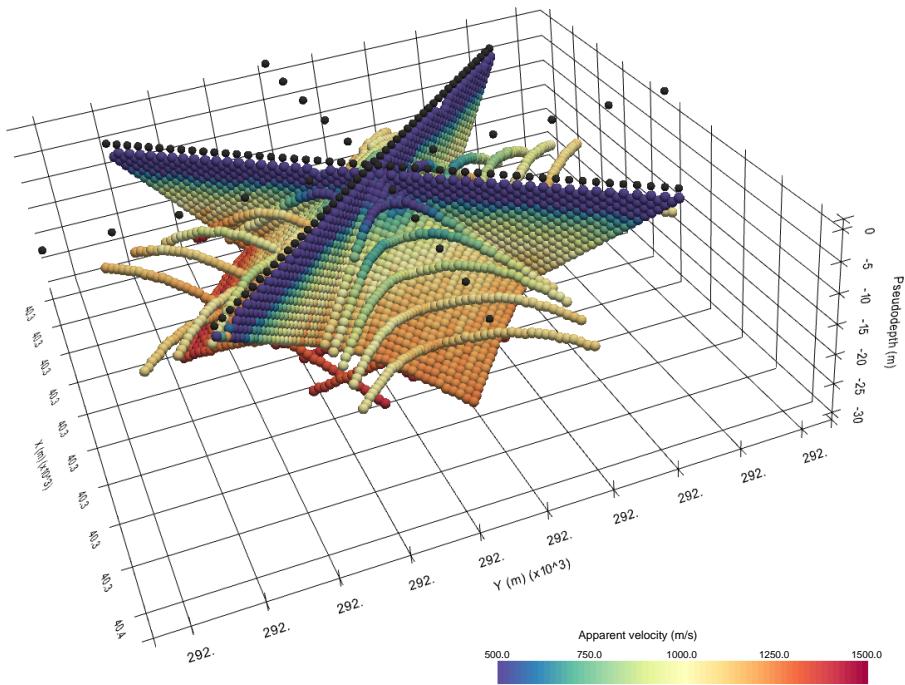
**Figure 11:** Comparison of forward modeled synthetic and automatically picked first break traveltimes for the synthetic seismic waveform data created in this study.



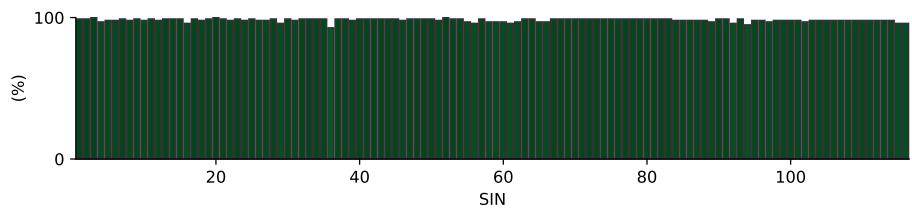
**Figure 12:** The Danube island soda lakes field dataset was collected along in a single line with a roll-along 3D survey layout; the with stations (co-located geophones and shots indicated by filled triangles indicate the direction dots) deployed in form of the measurements a cross. The five segments have an overlap. Additional shots (yellow stars) were conducted from of 50% a cross rotated by 45° to ensure an adequate data increase the coverage along of the entire profile dataset.



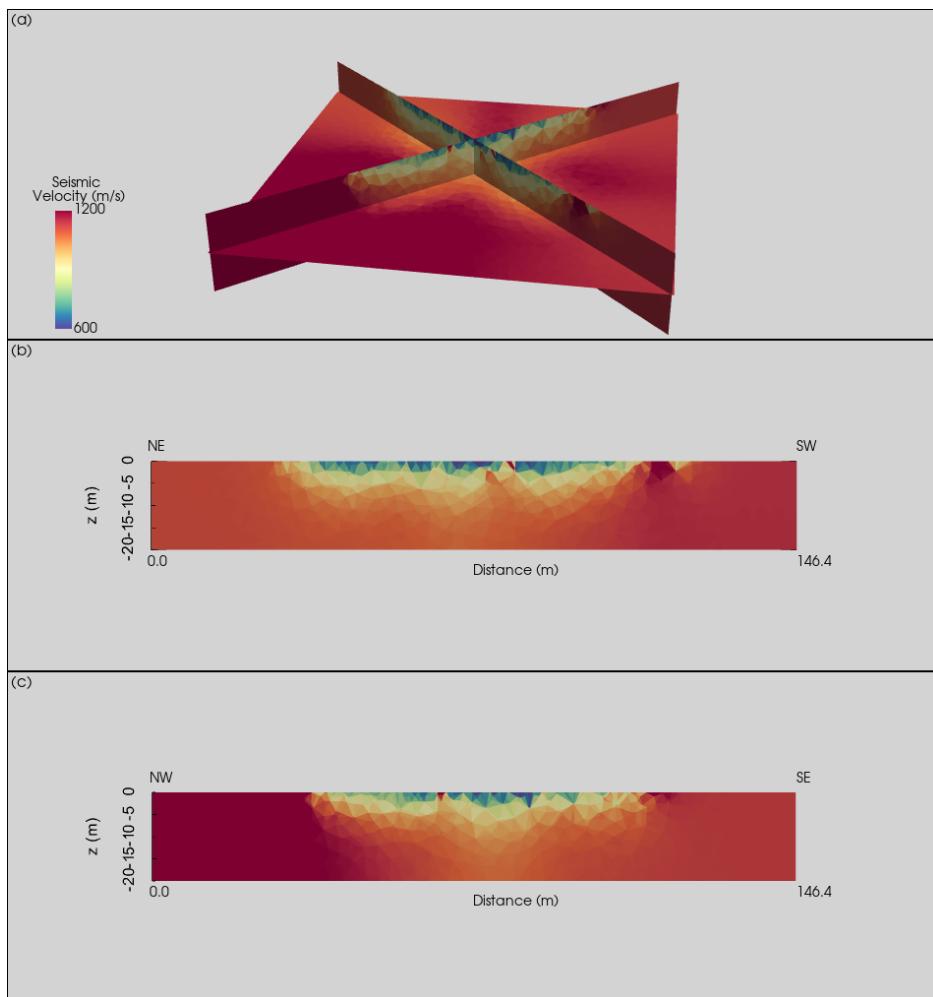
**Figure 13:** Exemplary seismic waveforms from the soda lakes dataset shown for shot index number (SIN) 1 with a 100 Hz lowpass filter applied to suppress high frequency noise. The recorded seismic waveforms clearly reflect the geometry of the geophones with RIN 1 to 48 deployed along the direction 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. (a) 3D representation of orthogonal slices through the resolved subsurface model. (b) 2D representation of the NE-SW slice. (c) 2D representation of the NW-SE slice.