# ing.grid

## plotID - a toolkit for connecting research data and visualization

Martin Hock [iD] [1]

Hannes Mayr [iD] [1]

Manuela Richter [iD] [1]

Jan Lemmer [iD]

Peter F. Pelz [iD] [1]

1. Chair of Fluid Systems, Technische Universität Darmstadt, Darmstadt.

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

**Abstract.** While visualizations can carry a vast amount of information compared to text and are often used for validation, references to data and metadata resulting in these visualizations are not common. To provide such references, the software plotID provides two key modules that strive to seamlessly integrate into a generic, Python-based research workflow. The module *tagplot* generates or accepts a unique ID and anchors it (visibly) as a reference to a figure or picture. The module publish exports the figure along with the data, code and parameters used in its creation into folders named by the reference id for later reuse. The tools work to provide aid in research data management with simple base functionality as opposed to encompassing management frameworks. Later features and improvements will expand the scope and applicability to other programming environments.

## 1 Statement of need

Scientific results are published in the form of hypotheses, axioms and equations as well as text and diagrams. Likewise, research software is being published more and more frequently. The comprehensibility of scientific results is indispensable for scientific discourse and reproducibility. Hypotheses, axioms and equations are usually published in text form and can be referenced accordingly. Software can be made traceable and referencable through the use of version control software. But what about diagrams? A diagram published in a paper is difficult to trace because the (raw) data is usually not available. However, the traceability of diagrams and the data they contain is not only a challenge in publication but also in everyday research. Diagrams are used for visualization and are therefore often produced for interim results. While the researcher continues the research process with investigations, experiments or simulations, volatile but important information like metadata, background information and details of the data processing are lost.

To be able to reconstruct the complete path, a treasure map is needed, starting from a publication, marking important waypoints of the process back to the raw data and metadata. This map needs to be provided along with the product that will be reviewed the most - the created diagram. If diagrams - regardless of whether they are published later or only serve as interim results - are

**Figure 1:** Research workflow from left to right; Afterwards following the chain of references from right to left
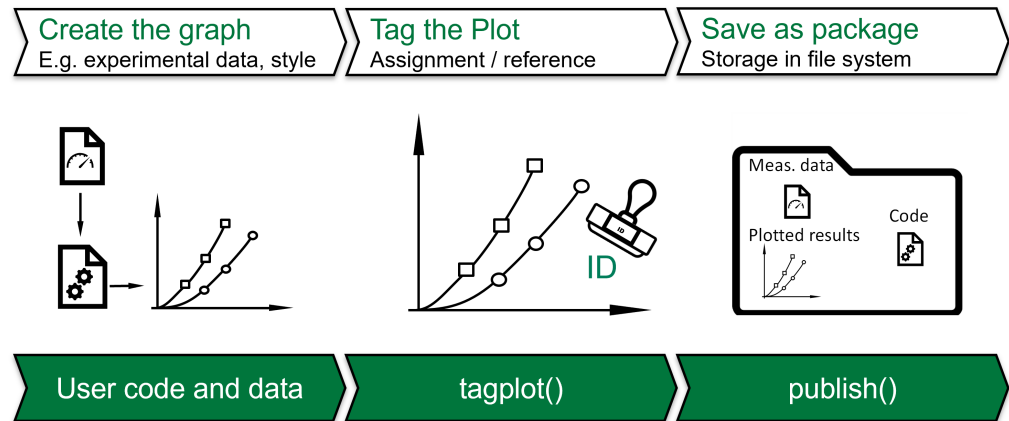'plotID-referencing' by Martin Hock, licensed under CC-BY-SA 4.0 ⓒⓘ◎

19 provided with an identifier, that connects to previous steps, traceability can be ensured. Figure
20 1 shows the order in which important elements are created and how the reference chain tracks
21 back.

22 A tool designed to meet these needs must meet the following requirements:

23 • Diagrams must have a unique identifier.

24 • The identifier must reference the raw data, relevant metadata and the code used to process
25 the data.

26 • The method must be easy to implement into the existing research workflow.

27 To reduce the effort of organizing figures along with all necessary data and metadata for later
28 review and reuse, the tool plotID was developed. plotID meets all the above-mentioned require-
29 ments and its implementation is described in this paper. The tool is limited to usage in an existing
30 python environment, but investigations on enabling independent installation and execution or
31 offering plotID as a web-based service are ongoing. The software depends on multiple python
32 libraries and is currently limited to visualizations from the [7] and general picture files such as
33 PNG and JPG.

34 Researchers often tend to keep an Excel table, noting down manually which data file corresponded
35 with which result along with input parameters. Sometimes an ID system is used (counting up or
36 using the date), but interim results like visualizations - used to verify results - are usually not
37 included. The authors could not find any other tool or software that aims at this specific task.
38 Few research data management frameworks strive to organize the full workflow via snapshots by
39 version controlling all software and data. Software execution commands have to be run through
40 the framework so the parameters can be captured and recorded in a snapshot. DataLad[6] works
41 like this. Labelling the figure with a corresponding ID, however, is unique to plotID. The scope
42 of plotID is also a different, limited one and integrating into frameworks with a broad scope like
43 DataLad is entirely possible as a future extension of functionality.

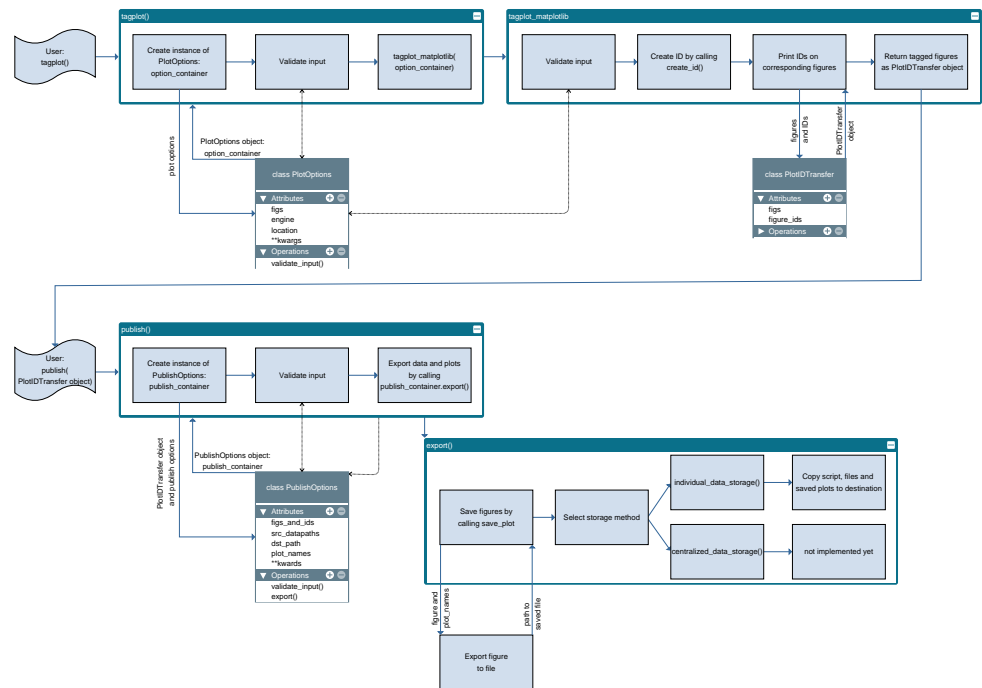**Figure 2:** Workflow integrating the plotID core functions
'plotID-workflow' by Martin Hock, licensed under CC-BY-SA 4.0 ⓒ ⓘ ⓢ

## 2   Methodology

The developed tool plotID is a software solution which covers the above-mentioned needs. The concepts and methods underlying the software are independent of the programming language. The aim of the software is to support the research workflow shown in figure 2 and to enable traceability. plotID aims to help during the early research process to decrease the work of making publications reproducible later on. To ensure ease of use, the tool has been designed to be integrated seamlessly into existing scripts. For this purpose, a GUI has been omitted. Instead, two main functions (building blocks) are provided, which can be inserted into existing user scripts as a one-liner. They are the core of plotID. The first module creates a (unique) ID and stamps this ID onto an object containing a visualization, while the second module helps organize all relevant code, software, and data that went into creating this graphic, into one complete package. Furthermore, connectivity to existing identifiers is ensured. If a specific visualization is later chosen to be included in a publication, the ID can be replaced by a permanent identifier like a DOI and the package of code, software and data can be published at the location referenced by the DOI. The ID on in the published paper will then directly reference the data, software and code used to create it, hence curating reproducibility. In the following, plotID is presented in more detail using the Python implementation.

## 3   Python - Implementation

The first version of plotID was implemented in Matlab since this is the most widely used programming language in the local working environment and the language the authors had the highest familiarity with. After reaching a usable state, the focus shifted to rewriting the tool in Python, the second most used programming language (locally). The goal was to make plotID accessible to a broader audience. Moreover, in contrast to Matlab, Python better fulfils the requirements for reusable software in the sense of the FAIR principles. In addition to being widely used in the engineering and research community, Python is non-proprietary, open source, easy to install or even shipped along many operating systems. Python also offers a package index (*PyPI*) and installer (*pip*) for easy distribution of software packages.

**Figure 3:** System architecture diagram
'plotID-system-architecture' by Martin Hock, licensed under CC-BY-SA 4.0 ⓒⓘⓢ

## 4   Core functions

The core functions of plotID are *tagplot()* and *publish()*. *tagplot()* generates an ID and adds
this ID to the figure object. *publish()* saves the figure object and image file, along with the
script file, plotID was called from – the essential elements necessary to recreate the visualization
from scratch. At this time, plotID does not yet export the python environment or the imported
modules. Refer to the subsection 'Script' for more information about the planned implementation.
Additional functionality might bring additional steps with the further development of plotID and
a widening of its scope.

### 4.1   tagplot()

The *tagplot()* function creates an ID and tags the figure object with this ID.

### 4.1.1   ID

*tagplot()* creates a unique ID (unique in a restricted system), that consists of a static part and
a generated part. The static part is handed over as a parameter and is meant to be used to
identify a project or organizational unit to which the figure is assigned. The generated part is by
default created from the UNIX-Time stamp in hexadecimal form. As an alternative option, a
random number generator can be used. The implementation of the ID is modular, enabling easy
implementation of individual needs or sources for IDs.

### 88 4.1.2 Tagging

89 In Python, there are multiple available packages that can produce visualizations from data.
90 Adding an ID needs to be implemented for many of these engines separately. For now plotID
91 supports figures created with *matplotlib* and raw image files. The ID is added as an attribute to
92 the object and the graphical, visible item.

### 93 4.1.3 Arguments

94 Necessary input arguments for *tagplot(figs, engine[, prefix, id_method, location])*:

95 • *figs*: the figure object or a list of objects, that is to be tagged

96 • *engine*: the plot/image engine to be used (currently only 'matplotlib' and 'image' (for plain
97   image files) are supported)

98 Optional input arguments are:

99 • *prefix*: to define a static part of each created ID (prefix='Ing.grid-'). Type: string.

100 • *id_method*: to define how the unique part of the ID is created ('random', 'time'). Type:
101   string.

102 • *location*: to define the position the ID is displayed in, relative to the full graphical ob-
103   ject (cardinal directions like 'west', custom inputs for rotation and position are to be
104   implemented). Type: string.

105 Output arguments are the tagged object and ID, if a list of objects was input, then the output is a
106 list as well.

107 At this point the figure object can still be modified, for example, to adjust colours or positioning
108 or recreate the full plot before exporting a final version.

### 109 4.2 publish()

110 This function starts the export process. The source files of the processed data, the visualization
111 (including the tagged ID), and the script hosting the call to the publish function are copied
112 together into a destination folder.

### 113 4.2.1 Script

114 A function in Python has access to the file path of the script which it was called from. With this,
115 the code for calculations can easily be gathered. For this reason, *publish()* cannot be called from
116 the command line or from within a script that has been started with the 'python -m' flag.

117 For dependent packages, the 'import' lines of the script can be compared with the output of 'pip
118 freeze' which returns all installed packages including their version. The overlap of these lists
119 can be written into a requirements.txt, which is added to the exported folder. Using the Python
120 package installer *pip* the original package versions can be reinstalled. This has not yet been
121 implemented. Furthermore, the user has to take care of including additional function files as
122 data paths, that have not been imported but are still accessed by the executed script.

**123**  **4.2.2   Data files**

**124**  Data files are handed over as a list of file or folder paths. Ideally, the script already manages
**125**  a list of all files that are read during the execution of the script. It is up to the user to control
**126**  this. By default, the data files are copied to each exported package. For large data files, the
**127**  *centralized* flag is intended.

**128**  The data files are copied to a central folder, relative to the export packages. For further exports,
**129**  the data files are compared to the ones already copied and only copied if new data files are
**130**  present. With this, a publication on a data repository could encompass the data files in addition
**131**  to multiple "satellite" folders containing the specific script, parameters and graphics. For HDF5
**132**  files, each package can contain an empty HDF5 file that only contains a link to the "real" central
**133**  data file. While this has proven to be helpful in the Matlab implementation, the Python version
**134**  aims to include the 'centralized' option in a future release.

**135**  **4.2.3   Arguments**

**136**  Necessary input arguments for *publish(src_datapath, dst_path, figure, plot_name[, ...])* are:

**137**  • *src_datapath*: This can be a single or a list of file or folder paths for source data and
**138**      additional function files. The type is a string or a list of strings.

**139**  • *dst_path*: This is the destination folder path. If it does not exist, the folder will be created.
**140**      The type is a string.

**141**  • *figure*: This is a figure object, the exact class depends on the plot engine used. This object
**142**      will be turned into an image file.

**143**  • *plot_name*: This is the name for the graphics objects. The type is a string or list of strings.
**144**      If a single name is passed for multiple objects, a raising number will be added.

**145**  Optional input arguments:

**146**  • *data_storage*: Currently only 'individual' and 'centralized' are available. 'Individual' will
**147**      store all data in each exported package, while 'centralized' stores the data files in a central
**148**      folder separate from the packages containing script and image files.

**149**  # **5   Example script**

**150**  The following script shows how plotID is used.

```python
10  # %% Import modules
11  import numpy as np
12  import matplotlib.pyplot as plt
13  from plotid.tagplot import tagplot
14  from plotid.publish import publish
15
16  # %% Set Project ID
17  PROJECT_ID = "MR04_"
18
```

```
19  # %% Create sample data
20  x = np.linspace(0, 10, 100)
21  y = np.random.rand(100) + 2
22  y_2 = np.sin(x) + 2
23
24  # %% Create sample figures
25
26  # 1. figure
27  IMG1 = 'image1.png'
28  FIG1 = plt.figure()
29  plt.plot(x, y, color='black')
30  plt.plot(x, y_2, color='yellow')
31  plt.savefig(IMG1)
32
33  # 2. figure
34  IMG2 = 'image2.png'
35  FIG2 = plt.figure()
36  plt.plot(x, y, color='blue')
37  plt.plot(x, y_2, color='red')
38  plt.savefig(IMG2)
```
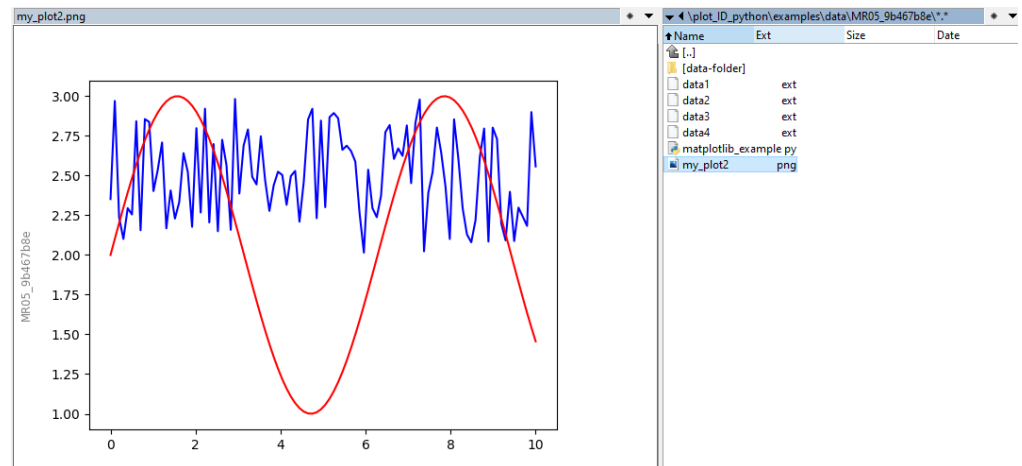
In this part, the plotID modules and those necessary to create figures and images are imported. The variable *PROJECT_ID* is set to provide a static part of the ID. Random data is used to create two figures with matplotlib and their image files.

```
42  # %% TagPlot
43
44  # If multiple figures should be tagged, they must be provided as list
        .
45  FIGS_AS_LIST = [FIG1, FIG2]
46  IMGS_AS_LIST = [IMG1, IMG2]
47
48  # Example for how to use tagplot with matplotlib figures
49  [TAGGED_FIGS, ID] = tagplot(FIGS_AS_LIST, 'matplotlib',
50          prefix=PROJECT_ID, id_method='time', location='west')
51
52  # Example for how to use tagplot with image files
53  # [TAGGED_FIGS, ID] = tagplot(IMGS_AS_LIST, 'image', prefix=
        PROJECT_ID,
54  #                                 id_method='random', location='west')
```

Both matplotlib objects are tagged with a generated ID in one line of code. Tagging the image files has been commented out in this case.

```
54  # %% Publish
```

**Figure 4:** Example export folder and tagged plot
'plotID-example-export' by Martin Hock, licensed under CC-BY-SA 4.0 ©①②

201  Files (README.md and LICENSE) and a folder from the code repository are used in place of
202  research data files. The folder ending with '-exports' is the destination, and 'testimage' is a freely
203  chosen name for the exported image files.

204  This also shows that the workflow does not depend on any kind of file format or pre-organized
205  structures. Any kind of data can be used, and even if the library creating the visualization is not
206  (yet) supported, the resulting image file can still be tagged.

207  Figure 4 shows the resulting export folder with (renamed) data files, the script file and the tagged
208  plot.

## 6  Distribution

210  Providing easy ways to acquire and use the software is important for adoption. The code is
211  open source under the Apache-v2.0 license. plotID requires a Python version $\geq 3.10$ and is
212  OS-independent. The current release version is v0.2.1 showing an alpha state.

213  At this time, the following distribution methods are available and described in the repository's[5]
214  README file.

### 6.1  Source Code

216  The plain source code is publicly available on a GitLab repository located under git.rwth-
217  aachen.de/plotID/plotID_python/[5] and can be directly downloaded or cloned with git.

```
1 git clone https://git.rwth-aachen.de/plotid/plotid_python.git
2 cd plotid_python
3 pip install -r requirements.txt
4 pip install .
```

**222**  ## 6.2 Python Package

**223**  plotID is listed in the official Python Package Index (PyPI)[10]. The installation is done with the
**224**  following command:

**225**  ```
pip install plotid
```

**226**  Distributing plotID independently from an existing Python installation is one of the aims of
**227**  later versions. Possible ways to achieve this are providing compiled executables or a central
**228**  web-hosted service.

**229**  ## 7 Ensuring good software quality

**230**  To ensure continuous good software quality, we adhere to best practices and the style guide
**231**  PEP-8[9]. This includes comments, docstrings and code formatting. To ensure adherence to
**232**  these guidelines, automated tests on the code are implemented.

**233**  ### 7.1 Unit tests

**234**  Python offers various libraries for unit testing. plotID is using the *unittest* module[11], which
**235**  is delivered with Python by default. Tests for each function are defined in the *tests* folder,
**236**  along with the *runner_test.py* script which organizes the execution of the tests, by discovering
**237**  the test files based on their location. The *unittest* module also measures the code covered by
**238**  the tests, and total coverage of less than 95% is considered as failed. The tests are executed
**239**  by a GitLab CI/CD pipeline[3] with every commit and merge request. Additional Jobs in the
**240**  pipeline execute Pylint[2] and Flake8[1] to check against coding style, programming errors and
**241**  cyclomatic complexity. Commits that fail the pipeline tests cannot be merged into the main
**242**  branch and will not make it into a release version. In the future, additional tests e.g. against
**243**  security risks introduced by dependencies and more detailed reports are planned.

**244**  ### 7.2 Documentation

**245**  To ensure easy access and understanding of the code, Python docstrings[8] have been implemented
**246**  in the source code from the beginning. The docstrings are compiled into HTML using the
**247**  Sphinx[12] python package and GitLab CI-CD[3] creating an automatically generated API
**248**  reference. The documents are hosted using GitLab Pages[4]. This documentation[13] will be
**249**  improved by adding the readme, example code, example use cases and an introductory text until
**250**  version 1.0.

**251**  ## 8 Conclusion

**252**  The idea of plotID is simple yet. As with most research data management operations, the benefit
**253**  for the additional work presents at a later time – although in this case, it presents itself for
**254**  the creator of data or visualizations and not only for later reuse. The code and open-source
**255**  implementation are still work-in-progress, but the core functionality is present. There are many
**256**  ideas to improve and add functionality present already. Bug reports, merge requests with code,
**257**  ideas for features and all feedback are welcome and best voiced in the GitLab repository.

## 9 Acknowledgements

## 10 Roles and contributions

**Martin Hock:** Conceptualization, Methodology, Coding, Tests, Writing – original draft

**Hannes Mayr:** Coding, Tests, Methodology

**Manuela Richter:** Conceptualization, Methodology, Coding

**Jan Lemmer:** Conceptualization, Methodology

**Peter F. Pelz:** Project administration, Supervision, Funding Acquisition

## References

[1] GitHub. *flake8/index.rst at main · PyCQA/flake8*. 2022. URL: https://github.com/PyCQA/flake8 (visited on 08/29/2022).

[2] GitHub. *PyCQA/pylint: It's not just a linter that annoys you!* 2022. URL: https://github.com/PyCQA/pylint (visited on 08/29/2022).

[3] *GitLab CI/CD | GitLab*. 2022. URL: https://docs.gitlab.com/ee/ci/ (visited on 08/19/2022).

[4] *GitLab Pages | GitLab*. 2022. URL: https://docs.gitlab.com/ee/user/project/pages/ (visited on 08/29/2022).

[5] GitLab RWTH Aachen. *PlotID / plotID_python · GitLab*. 2022. URL: https://git.rwth-aachen.de/plotid/plotid_python (visited on 08/19/2022).

[6] Yaroslav O. Halchenko et al. "DataLad: distributed system for joint management of code, data, and their relationship". In: *Journal of Open Source Software* 6.63 (2021), p. 3262. DOI: 10.21105/joss.03262. URL: https://doi.org/10.21105/joss.03262.

[7] *Matplotlib - Visualizations with python*. 2022. URL: https://matplotlib.org/ (visited on 10/10/2022).

[8] *PEP 257 – Docstring Conventions | peps.python.org*. 2022. URL: https://peps.python.org/pep-0257/ (visited on 08/29/2022).

[9] *PEP 8 – Style Guide for Python Code | peps.python.org*. 2022. URL: https://peps.python.org/pep-0008/ (visited on 10/11/2022).

[10] PyPI. *PyPI · The Python Package Index*. 2022. URL: https://pypi.org/ (visited on 08/19/2022).

293 [11] *unittest — Unit testing framework — Python 3.10.6 documentation*. 2022. URL: https:
294     //docs.python.org/3/library/unittest.html (visited on 08/29/2022).

295 [12] *Welcome — Sphinx documentation*. 2022. URL: https://www.sphinx-doc.org/en/m
296     aster/ (visited on 08/29/2022).

297 [13] *Welcome to PlotID's documentation! — plotID 0.1.2 documentation*. 2022. URL: https:
298     //plotid.pages.rwth-aachen.de/plotid_python/ (visited on 08/29/2022).