

# GOOD PRACTICES FOR RESEARCH SOFTWARE DEVELOPMENT

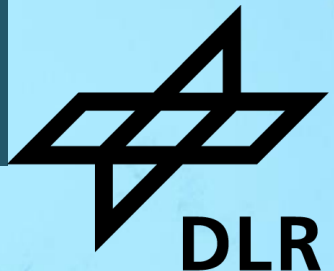
Workshop about Good Practices for Research Software Development,  
Hamburg, 19.02.2024

Tobias Schlauch <Tobias.Schlauch@DLR.de>

Institute for Software Technology

German Aerospace Center (DLR)

<http://www.dlr.de/sc>




# Good Practices for Research Software Development

## Why Should I Care?



- **It is for you!**

- To easily come back to your code 6 months later
- To enhance trust in your code and results produced with it
- To enhance chance to reuse (parts of) your code



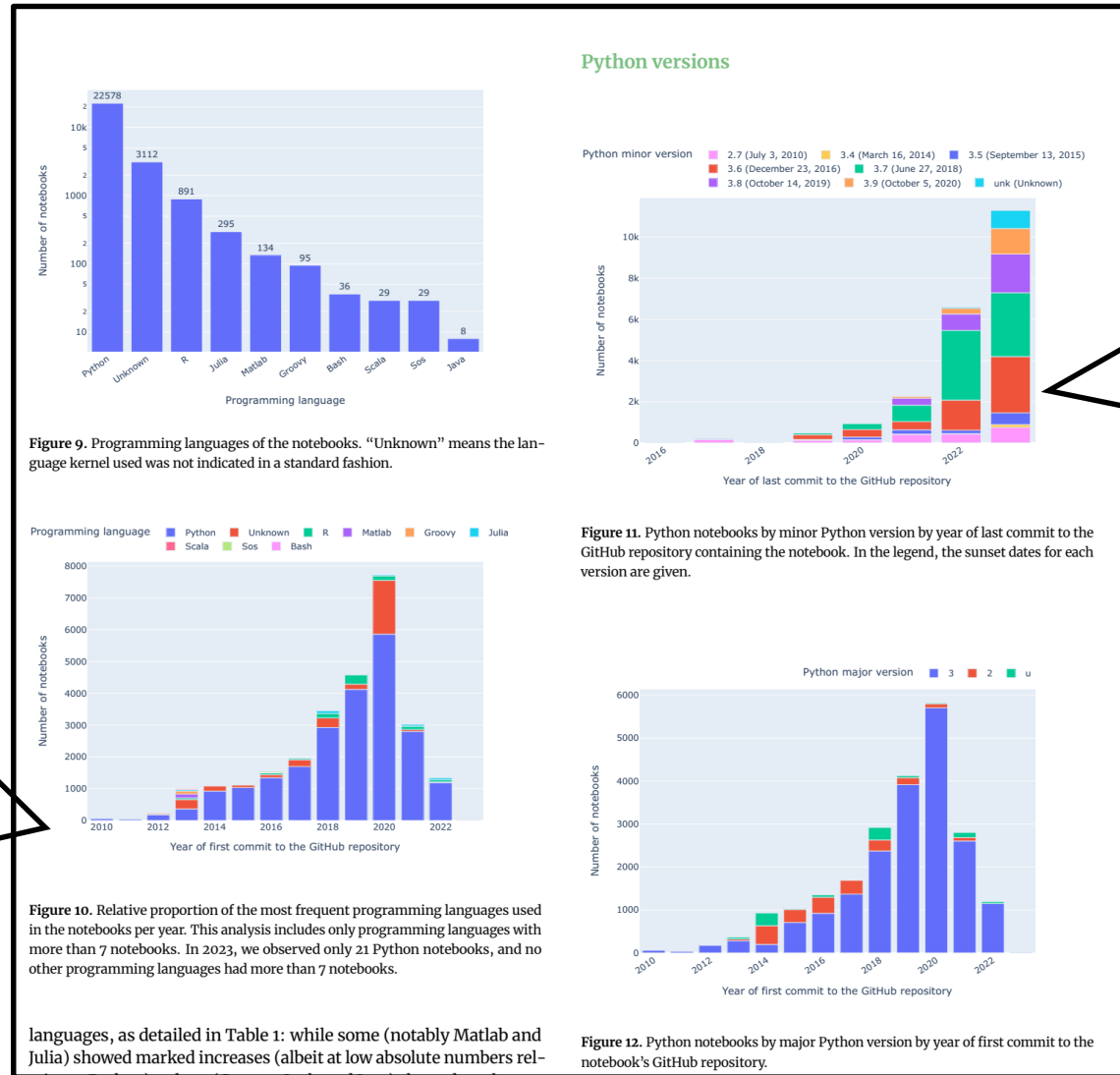
**For me:**  
Sometimes even  
after 1 week ;)

A speech bubble with a black outline and a tail pointing towards the left, containing text.

- **It helps others!**

- To get a better understanding of what you did
- To easier collaborate on your code
- To reproduce results based on your code

# Good Practices for Research Software Development Help to Enhance Reproducibility!



“How have they calculated this statistics?”

“Where is the analysis code and how should I run it?”

Source: Sheeba Samuel, Daniel Mietchen, “Computational reproducibility of Jupyter notebooks from biomedical publications”, <https://doi.org/10.48550/arXiv.2308.07333>, Licensed under: [CC-BY-4.0](https://creativecommons.org/licenses/by/4.0/)

# Good Practices for Research Software Development

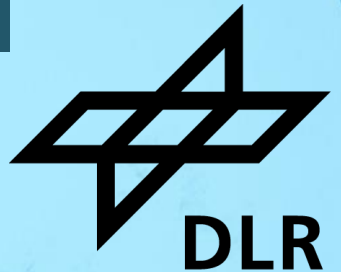
## What Recommendations Exist?



- [Model Policy on Sustainable Software at the Helmholtz Centers](#)
  - [DLR Software Engineering Guidelines](#)
  - Materials of the workshop [Foundations of Research Software Publication](#)
  - Your research domain, journal specific, ... recommendations?
- **There are many recommendations available! But how do I know exactly what to do ...?**

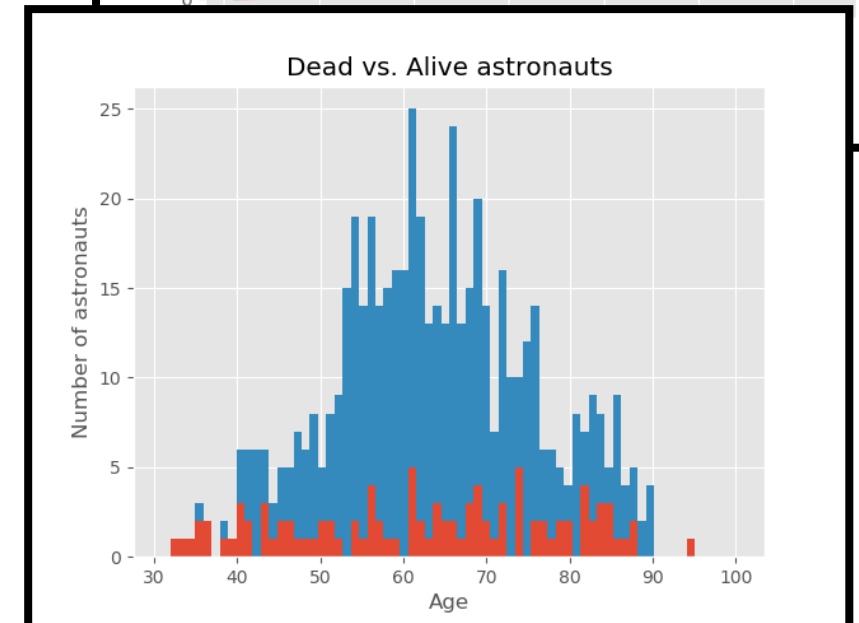
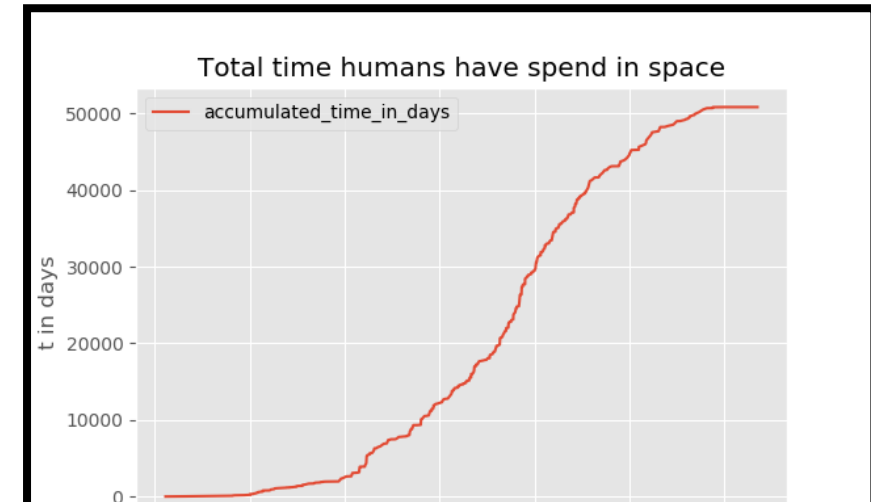


# OVERVIEW ABOUT TYPICAL GOOD PRACTICES



# Example: Astronaut Analysis

- [Astronauts Analysis](#) is a data publication consisting of:
  - Data set
  - Analysis script written in Python using [pandas](#) and [matplotlib](#)
  - Result plots
- **Scenario:**
  - I created it on my own as part of my job.
  - I want to publish it with my research paper.
  - I want to make its reuse as easy as possible and make it available under an open source license.



# Recommendations from the Workshop “Foundations of Research Software Publication”



- Step 1: Put your code under version control
- Step 2: Make sure that your code is in a sharable state
- Step 3: Add essential documentation
- Step 4: Add a license
- Step 5: Make your code citable
- Step 6: Release your code

Essential aspects  
which you should  
try to already  
address for  
“internal” software!

# Step 1: Put Your Code Under Version Control

## Where Should I Store My Code?



**Minimum:** Use a local Git repository + backup

**Recommended:** Use a code collaboration platform

astronauts.json

boxplot.png

combined\_histogram.png

female\_humans\_in\_space.png

humans\_in\_space.png

main.py

A local file explorer window showing a project directory. It contains a JSON file named 'astronauts.json', a Python script 'main.py', and several plot files: 'boxplot.png', 'combined\_histogram.png', and 'female\_humans\_in\_space.png'. The plots show age distributions and cumulative space time.



Repository

Ignore temporary Python files and only allow plot files in the results folder  
Schlauch, Tobias authored 1 year ago

1-put-into-git astronaut-analysis / +

Name	Last commit	Last update
results	Move plot files into a separate folder	1 year ago
.gitignore	Ignore temporary Python files and only...	1 year ago
astronauts.json	Add initial version	1 year ago
main.py	Add initial version	1 year ago

A screenshot of a GitHub repository page for 'astronaut-analysis'. The repository is under the 'Foundations of Research Software Publication' organization. It shows a commit history table with columns for Name, Last commit, and Last update. The files listed are 'results', '.gitignore', 'astronauts.json', and 'main.py', all with their last update dates.



# Step 1: Put Your Code Under Version Control

## What Belongs in the Repository?



- **Everything to make a usable version of your code** such as:
  - Source code, documentation, build scripts, test cases, configuration files, input data, ...
- **Avoid adding generated files** such as:
  - Third-party libraries, generated binaries, ...
- **How to handle large (data) files?**
  - Available could be [git-lfs](#), [git-annex](#), [Datalad](#) or your research data management publication repository
- **Please note:**
  - Details depend on the “product” that you manage in the Git repository
  - **.gitignore files** helps you to control what goes into your repository. See also <https://gitignore.io/> for templates.

# Step 1: Put Your Code Under Version Control

## Key Points



- Version control helps you to keep track of changes and is the basis for collaboration with others.
- Make sure to add all relevant files (or link them properly) to the source code repository.
- **.gitignore files** helps you to specify things that you do not want to share.
- Know your version control system properly.

# Step 2: Make Sure That Your Code Is in a Sharable State

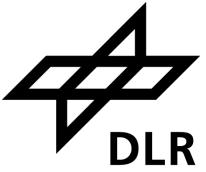
## General Hints



- Make sure others can run your code:
  - No dependencies on internal resources (servers, storage, licensed software, ...)
  - No absolute paths
  - Clearly state dependencies + provide required build / installation scripts (e.g.: [pip-tools](#), [poetry](#)) => crucial aspect of reproducibility
- Organize files in a suitable directory structure (e.g.: [Python Application Layouts](#), [Good Data Practices](#))
- Do not share sensitive data such as passwords, user accounts, SSH keys, internal IP addresses, etc. (e.g.: [gitleaks](#))
- Orientate on standards of your domain / community

# Step 2: Make Sure That Your Code Is in a Sharable State

## Improve Your Code Style and Structure



- Strive for understandable code:
  - Apply a code style – consistency is more important than convenience (e.g.: [PEP8](#))
  - Use a consistent and light code layout
  - Structure your code in suitable "building blocks" such as functions
  - Use specific and appropriate names for all artifacts
  - Provide sufficient level of code comments
- Read code of others for inspiration
- Try to do pair programming and reviews (even if it is [with your rubber duck](#))

# Step 2: Make Sure That Your Code Is in a Sharable State

## Think About Testing and Automation



- Small tests are done easily but already show effect:
  - Code linters and checkers help to find poor code snippets and help to enforce coding styles (e.g.: [flake8](#), [black](#))
  - Automated tests work as an executable documentation (e.g.: [pytest](#))
- Tests offer a good starting point for your automation efforts!



# Step 2: Make Sure That Your Code Is in a Sharable State

## Example After Step 2

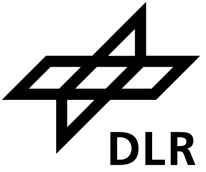


```
astronaut-analysis.py
1 """ This script analysis the astronaut data set and creates different plots as result. """
2
3
4 from datetime import date
5
6 import pandas as pd
7 import matplotlib.pyplot as plt
8
9
10 _ASTRONAUT_DATA_FILE = "../data/astronauts.json"
11
12
13 ##
14 # Data preparation functions
15 ##
16 def prepare_data_set(df):
17     df = rename_columns(df)
18     df = df.set_index("astronaut_id")
19
20     # Set pandas dtypes for columns with date or time
21     df = df.dropna(subset=["time_in_space"])
22     df["time_in_space"] = df["time_in_space"].astype(int)
23     df["time_in_space"] = pd.to_timedelta(df["time_in_space"], unit="m")
24     df["birthdate"] = pd.to_datetime(df["birthdate"])
25     df["date_of_death"] = pd.to_datetime(df["date_of_death"])
26     df.sort_values("birthdate", inplace=True)
27
28     # Calculate extra columns from the original data
29     df["time_in_space_D"] = df["time_in_space"].astype("timedelta64[D]")
30     df["alive"] = df["date_of_death"].apply(is_alive)
31     df["age"] = df["birthdate"].apply(calculate_age)
32     df["died_with_age"] = df.apply(died_with_age, axis=1)
33     return df
34
35
36 def rename_columns(df):
37     """
38     The original column naming in the data set is not useful
39     for programming with pandas. So we rename it.
40     """
```

- Applied PEP8 code style
- Cleaned up the code
- Added basic testing and more 😊

# Step 2: Make Sure That Your Code Is in a Sharable State

## Key Points

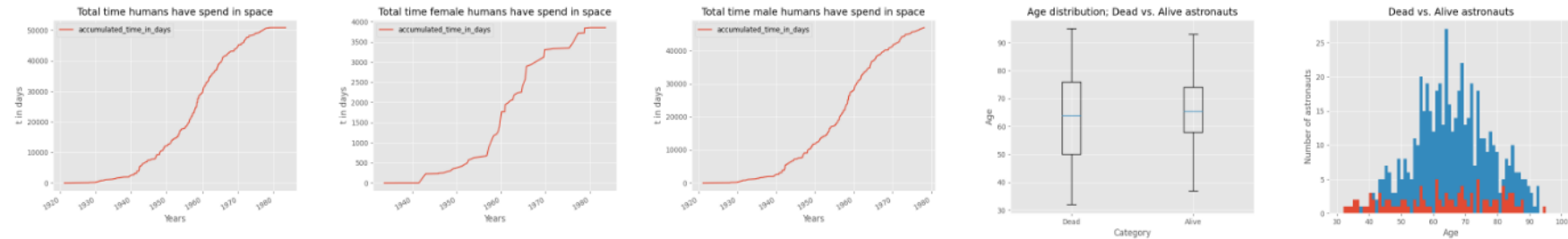


- Make sure that others can (re-)use your code
- Do not store secrets in your code repository
- Strive for understandable code
- Start introducing basic test automation

# Step 3: Add Documentation

## Astronaut Analysis

This analysis is based on publicly available astronauts data from [Wikidata](#). In this context, we investigated aspects such as time humans spent in space as well as the age distribution of the astronauts.



The repository is organized as follows:

- `data`: Contains the astronauts data set retrieved from Wikidata
- `code`: Contains the astronaut analysis script
- `results`: Contains the resulting analysis plots

## Astronaut Data

The data set has been generated using the following SPARQL query [1] (retrieval date: 2018-10-25).

You can also analyze a recent version of the astronaut data by replacing the data set and re-running the analysis script:

- Run the SPARQL query
- Download the resulting data formatted as JSON
- Replace the file `data/astronauts.json`
- Run the analysis script

## Astronaut Analysis Script

The script requires Python  $\geq 3.8$  and uses the libraries `pandas` as well as `matplotlib`.

The script has been successfully tested on Windows 10 and Linux with Python 3.8.

## Typical Structure:

- Software name
- Purpose
- Install
- Usage
- Contributing
- Citation Hint
- License

# Step 3: Add Documentation

## General Hints



- **Mind your target groups:**
  - **Typical perspectives:** Users, contributors
  - **Users:** Installation / usage instructions, tutorials, support channels, ...
  - **Contributors:** Contribution guidelines, technical overview, ...
- **Think about adding typical documentation files** such as:
  - README (project front page), CONTRIBUTING (contributions guidelines), CODE\_OF\_CONDUCT (communication rules), LICENSE (license information), CHANGELOG (major changes), CITATION (citation metadata)
- **Please note:**
  - [Markdown](#) or another markup language is quite often used to write documentation
  - Usually, you will need additional documentation, for example, in a `docs` directory (e.g.: [Sphinx](#), [MkDocs](#))

# Step 3: Add Documentation

## Key Points



- Provide documentation for relevant target groups
- Add a **README** file as a minimum documentation to your source code repository



# Astronaut Analysis Release 1.0.0



The screenshot shows the GitHub repository page for "Astronaut Analysis". At the top, it displays "13 Commits", "10 Branches", "1 Tag", "148 KiB Project Storage", and "1 Release". Below this, a message states "The repository contains the example code used in this workshop." A DOI "10.5281/zenodo.10001813" and a "Latest Release" tag "1.0.0" are visible. A commit by Tobias Schlauch is highlighted with a green checkmark and the hash "8a05544e". The repository structure includes folders for "LICENSES", "code", "data", and "results", and files for ".gitignore", ".gitlab-ci.yml", "CHANGELOG.md", and "LICENSE.md".

Name	Last commit	Last update
LICENSES	Add license and copyright information	2 years ago
code	Add license and copyright information	2 years ago
data	Add license and copyright information	2 years ago
results	Add license and copyright information	2 years ago
.gitignore	Add license and copyright information	2 years ago
.gitlab-ci.yml	Add license and copyright information	2 years ago
CHANGELOG.md	Add changelog and reference it	2 years ago
LICENSE.md	Add license and copyright information	2 years ago

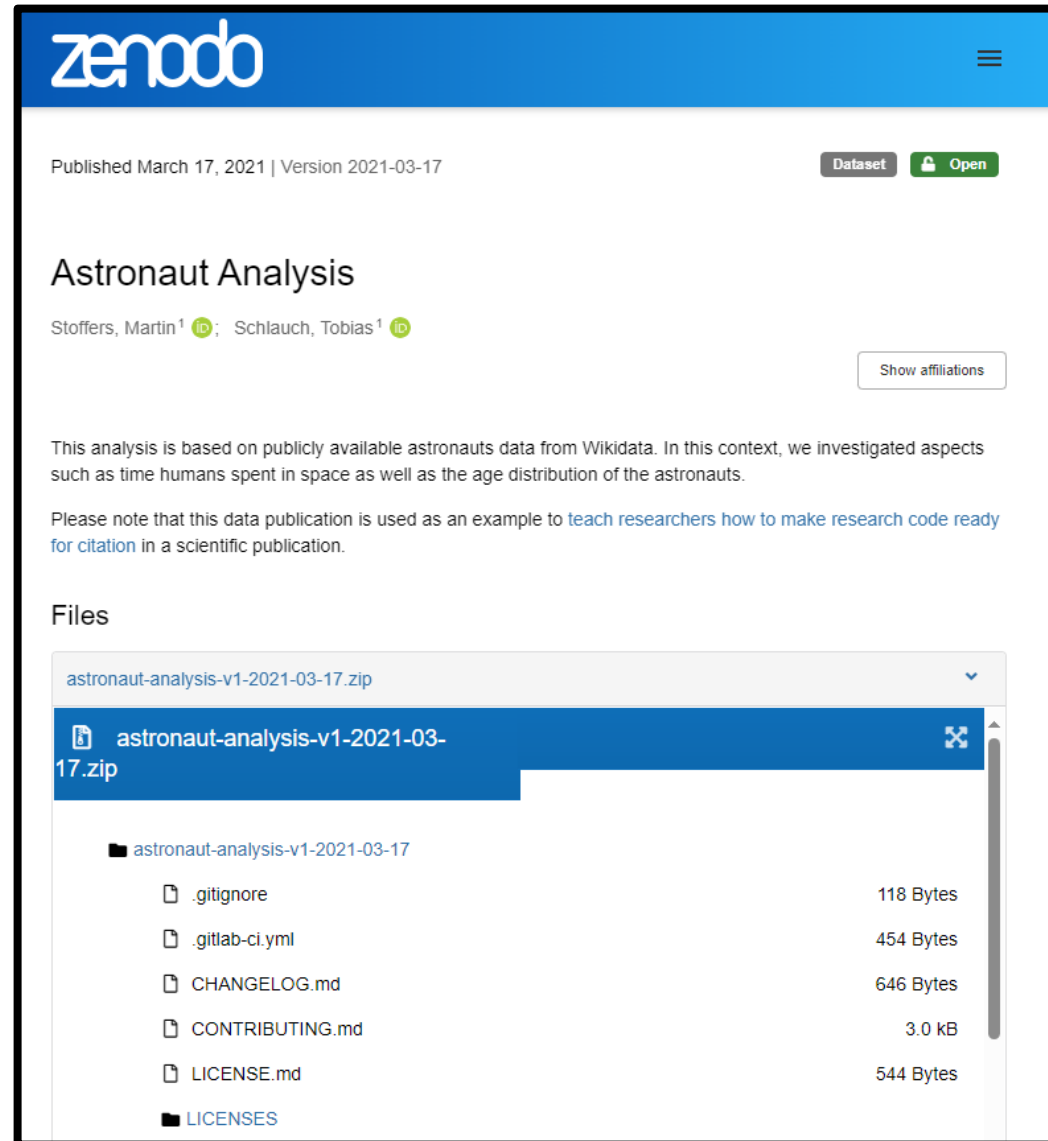
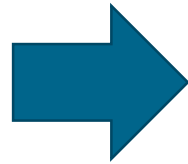
License information for code, data, results properly annotated via [REUSE](#)

Release 1.0.0 marked as **Git tag** in the repository

# Astronaut Analysis Release 1.0.0 (cont.)

DOI 10.5281/zenodo.10001813

```
cff-version: 1.2.0
title: Astronaut Analysis
message: >-
  If you use this dataset, please cite it using the
  metadata
  from this file.
type: dataset
authors:
  - given-names: Martin
    family-names: Stoffers
  - given-names: Tobias
    family-names: Schlauch
identifiers:
  - type: doi
    value: 10.5281/zenodo.10001813
```



The screenshot shows the Zenodo dataset page for "Astronaut Analysis". The page header includes the Zenodo logo and a menu icon. Below the header, it states "Published March 17, 2021 | Version 2021-03-17" and has buttons for "Dataset" and "Open". The title "Astronaut Analysis" is prominently displayed, followed by the authors "Stoffers, Martin" and "Schlauch, Tobias". A "Show affiliations" button is visible. The main text describes the analysis as being based on publicly available astronaut data from Wikidata and mentions its use as an example for teaching researchers. Below the text, a "Files" section shows a list of files for the dataset "astronaut-analysis-v1-2021-03-17.zip". The files listed are:

File Name	Size
.gitignore	118 Bytes
.gitlab-ci.yml	454 Bytes
CHANGELOG.md	646 Bytes
CONTRIBUTING.md	3.0 kB
LICENSE.md	544 Bytes
LICENSES	

## Citable Release:

- Citation metadata in [Citation File Format](#)
- DOI via [Zenodo](#)

# There Are Many Recommendations Available! But How Do I Know Exactly What to Do ...?



- Recommendations are typically made under certain assumptions. I.e., they leave out details and might not fit for your case directly ... ☹️
- Establishing detailed good practices on a **research group level** could help:
  - Similar tasks and projects make it easier to agree on relevant practices and details
  - Use generic recommendations as a starting point and leave out irrelevant aspects / add required details as needed
- “Executable” templates can help to get everyone better started:
  - Relevant tools: [Cookiecutter](#), [Cruft](#)
  - Example: [HCDC / Software Templates / Python Package · GitLab \(helmholtz.cloud\)](#)

# SUMMARY



# Summary



- **Good practices for research software development are important:**
  - Help you and others to work on code and have trust in results produced with it
  - Enhance chances for research to be reproducible
- **Existing recommendations are made under certain assumptions and need to be tailored to the right context:**
  - Existing guidelines might be too generic
  - Presented recommendations might be in some aspects too detailed or (currently) not relevant for your specific case
- **Research group could be the right level to establish effective good practice!**



# Thank you!

## What are your Questions?

Email: [Tobias.Schlauch@dlr.de](mailto:Tobias.Schlauch@dlr.de)

Mastodon: <https://norden.social/@schlauch>

HIFIS Mattermost: [@schlauch](#)



# Copyright and License Information



All content is © German Aerospace Center and licensed under [CC-BY-4.0](#) with the following exceptions:

- DLR logo, slide layout, © German Aerospace Center. All rights reserved.
- HIFIS logo, © HIFIS, [CC-BY-4.0](#).
- Philae landing on comet 67 P/Churyumov-Gerasimenko, slide 24, © German Aerospace Center. [CC-BY-3.0](#).