

Reproducible workflows with Jupyter: case study in materials simulation research

Hans Fangohr^{1,2,9}, Martin Lang^{1,9}, Samuel R Holt^{1,9}, Swapneel Pathak^{1,9}, Min Ragan-Kelley³, Marijan Beg^{4,2}, Juliette Belin⁵, Thomas Kluyver⁶, Olexandr Konovalov⁷, Nicolas M. Thiéry⁸

1) SSU Computational Science, Max Planck Institute for the Structure and Dynamics of Matter, Hamburg, Germany; 2) University of Southampton, Southampton, United Kingdom; 3) Simula Research Laboratory, Oslo, Norway; 4) Imperial College, London, United Kingdom; 5) Logilab, Paris, France; 6) European XFEL, Schenefeld, Germany; 7) University of St Andrews, St Andrews, United Kingdom; 8) Université Paris-Saclay, Orsay, France; 9) Center for Free-Electron Laser Science (CFEL)



Introduction

- Reproducibility is important as a principle of science
- Reproducibility makes science more efficient: *reproducibility enables re-usability*
- We focus on computational workflows (and assume experimental data is digitised)

Common barriers preventing reproducibility

- Source code of software used is not available
- Exact procedure, which led to the results reported in the publication, is not shared
- Information on the computing environment such as supporting libraries and (if required) code compilation details is not revealed

Proposal for reproducible publications [1]

- Publish software used in research as open source
 - for example in supplementary material of publication and/or
 - as GitHub repository with persistent archival and DOI from Zenodo
- Use Jupyter Notebooks to direct the study
 - notebooks encapsulate the exact procedure (sometimes called the “protocol”)
 - for the publication, provide notebooks that reproduce central results of the paper when executed: for example one notebook per figure
- Use a “binder-enabled” repository to define the required software environment
 - defines the software required to execute the notebooks, and enables automatic installation
 - mybinder.org provides zero-install cloud-based environments in which notebooks from binder-enabled repositories can be executed.

Example repository for reproducible publication [2]:

<https://github.com/lang-m/2022-paper-multiple-bloch-points>

What is a binder-enabled repository?

- (git) repository, typically hosting notebooks
 - containing some files that describe what software is required to execute the notebooks (such as `requirements.txt`, `apt.txt`, `environment.yml`, `Project.toml`, `install.R`, ...)
 - Binder uses existing conventions to describe software environments; thus there is little overhead to learn how to specify the requirements
 - A binder service (such as mybinder.org) can create a software environment (based on the specification) and start a Jupyter Notebook server within that environment, so that the notebooks in the repository can be executed.
 - A container with the same environment can be created locally using `repo2docker`
- Binder demo repository:
<https://github.com/fangohr/reproducibility-repository-example>

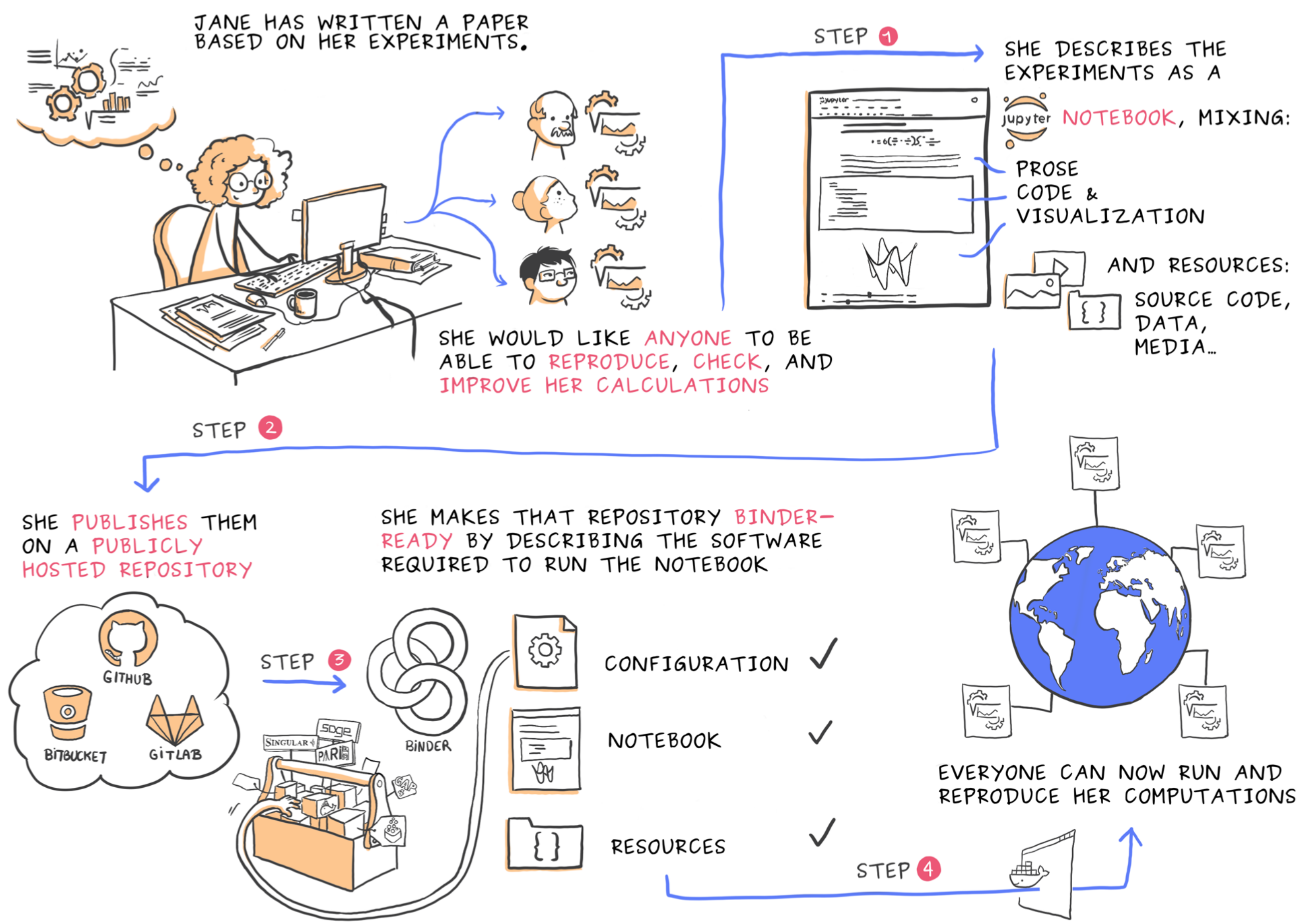
Issues

- Access and use of large data sets from Binder [3]
 - Sustainable funding of global mybinder.org service (e.g. <https://blog.jupyter.org/mybinder-org-reducing-capacity-c93ccfc6413f>)
- Some organisation address these issues by running their own Binder service (such as Max Planck Society: <https://notebooks.mpcdf.mpg.de/binder/>).

References

- M. Beg, J. Taka, T. Kluyver, A. Konovalov, M. Ragan-Kelley, N. M. Thiéry, and H. Fangohr, “Using Jupyter for reproducible scientific workflows,” *Computing in Science & Engineering*, vol. 23, pp. 36–46, Mar. 2021. <https://doi.org/10.1109/mcse.2021.3052101>.
- M. Lang, M. Beg, O. Hovorka, and H. Fangohr, “Bloch points in nanostrips,” *Scientific Reports*, vol. 13, Apr. 2023. <https://doi.org/10.1038/s41598-023-33998-z>.
- M. Ragan-Kelley, H. Fangohr, N. Thiéry, et al., “Supporting open, useful, and reproducible computational environments (SOURCE),” tech. rep., 2022. <https://github.com/minrk/horizon-widera-2022/blob/main/submitted-SOURCE-2022-04-20.pdf>.
- M. Beg, M. Lang, and H. Fangohr, “Ubermag: Toward more effective micromagnetic workflows,” *IEEE Transactions on Magnetics*, vol. 58, pp. 1–5, Feb. 2022. <https://doi.org/10.1109/tmag.2021.3078896>.

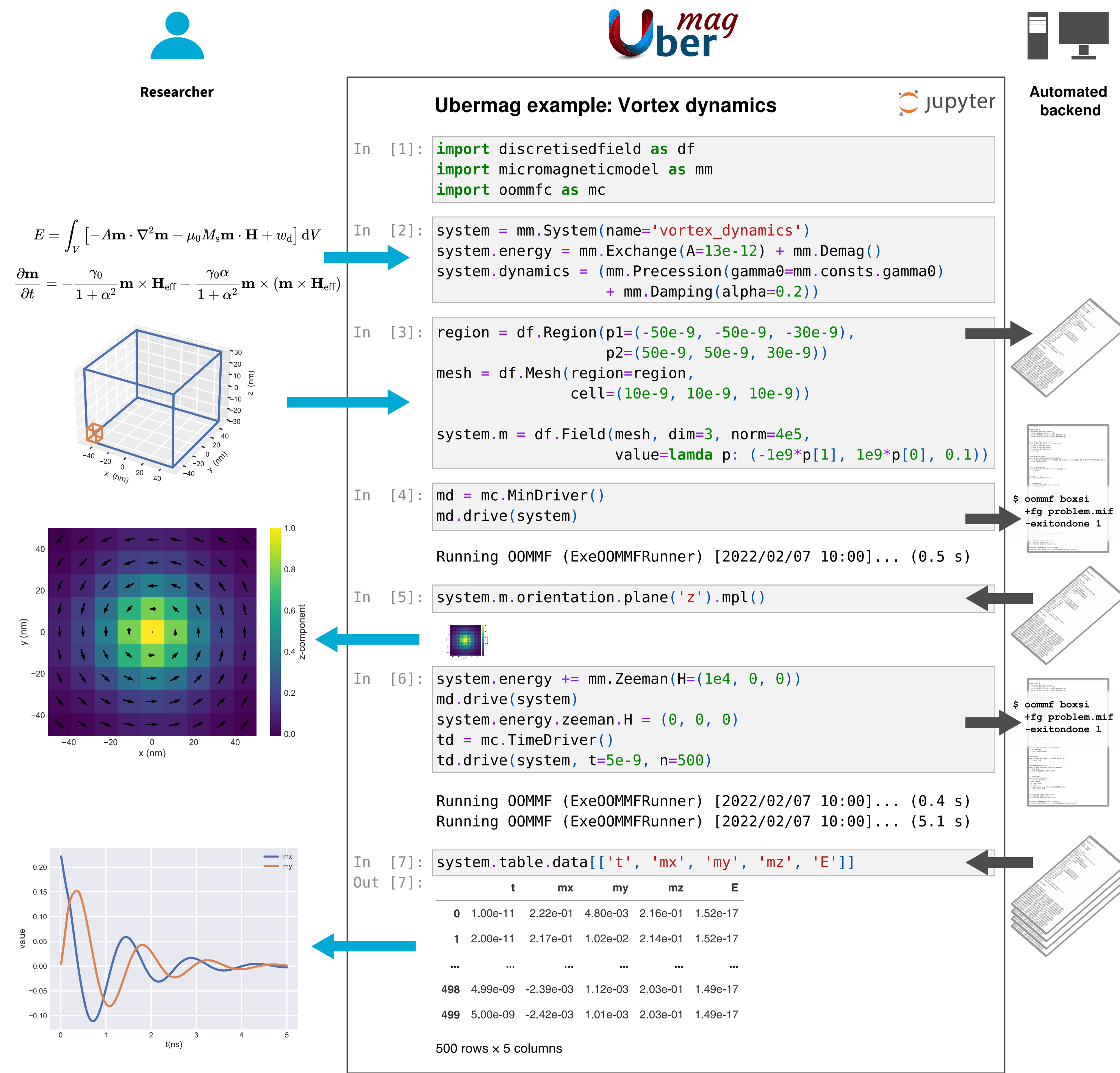
Reproducible publications with Jupyter — summary



A case study in which Jupyter Notebooks are used for reproducible data analysis in publications. The software environment is created through the Binder project when required [1].

Reproducible magnetic simulations with Ubermag

- “Ubermag” [4] (<https://ubermag.github.io>) provides a Python interface to (compiled) simulation software (referred to as “backends” below)
- Notebooks provide reproducible transcripts of simulation and analysis



Workflow and abstractions introduced by Ubermag for the computation of the dynamics of a magnetic vortex in a ferromagnetic nanostructure. The whole simulation can be driven from a single Jupyter Notebook as shown in the central column. The researcher defines the physical problem (cell 2) and the system to study (cell 3). Both can be visualised directly inside the notebook as indicated on the left-hand side. Ubermag internally communicates with a low-level simulation tool (cells 3, 4 and cell 6) shown as the automated backend in the right column, and returns the simulation results in a high-level format (cells 5 and 7). Plots of the simulation results and further data analysis are part of the notebook. The notebook shown here is a simplified and shortened version of an introductory notebook available at <https://ubermag.github.io/demo.html>.

