



Software sustainability of global impact models

Emmanuel Nyenah¹, Petra Döll^{1,2}, Daniel S. Katz³, and Robert Reinecke⁴

¹Institute of Physical Geography, Goethe-University Frankfurt, 60438 Frankfurt am Main, Germany

²Senckenberg Biodiversity and Climate Research Centre (SBiK-F), 60438 Frankfurt am Main, Germany

5 ³NCSA & CS & ECE & iSchool, University of Illinois Urbana-Champaign, Urbana, IL, 61801, USA

⁴Institute of Geography, Johannes Gutenberg-University Mainz, 55128 Mainz, Germany

Correspondence to: Emmanuel Nyenah (Nyenah@em.uni-frankfurt.de)

Abstract. Research software for simulating Earth processes enables estimating past, current, and future world states and guides policy. However, this modelling software is often developed by scientists with limited training, time, and funding, leading to software that is hard to understand, (re)use, modify, and maintain, and is, in this sense, non-sustainable. Here we evaluate the sustainability of global-scale impact models across ten research fields. We use nine sustainability indicators for our assessment. Five of these indicators – documentation, version control, open-source license, provision of software in containers, and the number of active developers – are related to best practices in software engineering and characterize overall software sustainability. The remaining four – comment density, modularity, automated testing, and adherence to coding standards – contribute to code quality, an important factor in software sustainability. We found that 29% (32 out of 112) of the global impact models (GIMs) participating in the Inter-Sectoral Impact Model Intercomparison Project were accessible without contacting the developers. Regarding best practices in software engineering, 75% of the 32 GIMs have some kind of documentation, 81% use version control, and 69% have open-source license. Only 16% provide the software in containerized form which can potentially limit result reproducibility. Four models had no active development after 2020. Regarding code quality, we found that models suffer from low code quality, which impedes model improvement, maintenance, reusability, and reliability. Key issues include a non-optimal comment density in 75%, insufficient modularity in 88%, and the absence of a testing suite in 72% of the GIMs. Furthermore, only 5 out of 10 models for which the source code, either in part or in its entirety, is written in Python show good compliance with PEP 8 coding standards, with the rest showing low compliance. To improve the sustainability of GIM and other research software, we recommend best practices for sustainable software development to the scientific community. As an example of implementing these best practices, we show how reprogramming a legacy model using best practices has improved software sustainability.



1 Introduction

30 Simulation models of the Earth system are essential tools for scientists and their outcomes are relevant for decision-makers (Prinn, 2013). They improve our understanding of complex subsystems of the Earth (Prinn, 2013; Warszawski et al., 2014) and enable us to perform numerical experiments that would otherwise be impossible in the real world, e.g., exploring future pathways (Wan et al., 2022; Satoh et al., 2022; Kemp et al., 2022). While so-called Earth System Models always include the simulation of atmospheric processes and thus compute climate variables and how they change due to greenhouse gas emissions, 35 so-called impact models enable us to quantitatively estimate the potential impacts of climate change on, e.g., floods (Sauer et al., 2021), droughts (Satoh et al., 2022), and food security (Schmidhuber and Tubiello, 2007). These impact models also quantify the historical development and current situation of, e.g., water stress, wildfire hazard, and fish population, thus providing crucial information for policymakers, scientists, and citizens. The central role of impact models can be seen in model intercomparison efforts of ISIMIP (Inter-Sectoral Impact Model Intercomparison Project) (Warszawski et al., 2014; ISIMIP, 40 2024) which encompasses more than 130 sectoral models (Frieler and Vega, 2019). ISIMIP uses bias-corrected climate forcings to assess the potential impacts of climate change in controlled experiments, and their outputs provide valuable contributions to the Intergovernmental Panel on Climate Change reports (Warszawski et al., 2014).

Impact models quantify physical processes related to specific components of the Earth system at various spatial and temporal 45 scales by using mathematical equations. The complexity of impact models is influenced by the complexity of the included physical processes, the choice of the perceptual and mathematical model, the computational effort needed for simulation, as well as their spatial-temporal resolution and spatial extent of the simulated domain (Azmi et al., 2021; Wagener et al., 2021). This complexity can result in models with very large source codes (Alexander and Easterbrook, 2015).

50 The software for these impact models is categorized as research software, which includes “source code files, algorithms, computational workflows, and executables developed during the research process or for a research objective” (Barker et al., 2022). Impact modelling research software is predominantly developed and maintained by scientists without formal training in software engineering (Hannay et al., 2009; Barton et al., 2022; Carver et al., 2022; Reinecke et al., 2022). Most of these researchers are self-taught software developers with little knowledge of software requirements (specifications and features of 55 software), industry-standard software design patterns (Gamma et al., 1994), good coding practices (e.g., using descriptive variable names), version control, software documentation, automated testing and project management practice (e.g. agile) (Carver et al., 2013, 2022; Hannay et al., 2009; Reinecke et al., 2022). We hypothesize that this leads to the creation of source code that is not well-structured, not easily (re)usable, difficult to modify and maintain, has scarce internal documentation (code comments) and external documentation (e.g. manuals, guides, and tutorials), and poorly documented workflows.

60



Research software that suffers from these shortcomings is likely difficult to sustain and has severe drawbacks for scientific research. For example, it can impede research progress, decrease research efficiency, and hinder scientific progress, as implementing new ideas or correcting mistakes in code that is not well-structured is more difficult and time-consuming. In addition, it increases the likelihood of erroneous results, thereby reducing reliability and hindering reproducibility (Reinecke et al., 2022). We argue that these harmful properties can be averted, to some extent, with sustainable research software.

There are various interpretations of the meaning of “sustainable research software”. Anzt et al. (2021) describe it as research software that is maintainable, extensible, flexible (adapts to user requirements), has a defined software architecture, is testable, has comprehensive in-code and external documentation, and is accessible (the software is licensed as Open Source with a digital object identifier (DOI) for proper attribution) (Anzt et al., 2021). Katz views research software sustainability as the process of developing and maintaining software that continues to meet its purpose over time (Katz, 2022). This includes adding new capabilities as needed by its users, responding to bugs and other problems that are discovered, and porting to work with new versions of the underlying layers, including software as well as new hardware (Katz, 2022). Both definitions share common aspects like the adaptation to user requirements but differ in scope and perspective. Katz’s definition is more user-oriented, focusing on the software’s ability to continue meeting its purpose over time. On the other hand, Anzt et al.’s definition is more developer-oriented, aiming to improve the quality and robustness of research software. We chose to adopt Anzt et al.’s definition in the following because it provides measurable qualities relevant to this study. In contrast, Katz’s definition is more challenging to measure and evaluate but is likely closer to the reality of software development. For example, one of the models in our analysis is more than 25 years old (Nyenah et al., 2023) and thus certainly was sustained during that period, while at the same time, it does not meet some sustainability requirements of Anzt et al.’s definition. It is possible that such software can be sustained but requires substantial additional resources.

Recent advances in developing sustainable research software have led to a set of community standard principles: FAIR (findable, accessible, interoperable, reusable) for research software (FAIR4RS), aimed towards increasing transparency, reproducibility, and reusability of research (Barker et al., 2022; Chue Hong et al., 2022). Software quality which impacts sustainability overlaps with the FAIR4RS principles, particularly reusability, but is not directly addressed by them (Chue Hong et al., 2022). Reusable software here means software can be understood, modified, built upon, or incorporated into other software (Chue Hong et al., 2022). A high degree of reusability is therefore important for efficient further development and improvement of research software, and thus for scientific progress. However, many models are not FAIR (Barton et al., 2022). To our knowledge, research software sustainability in Earth System Sciences has not been evaluated before. As an example of complex research software in the Earth System Sciences, in this study, we assess the sustainability of the software of global impact models (GIMs) that participate in the ISIMIP project to investigate factors that contribute to sustainable software development. The GIMs belong to the ten research fields (or impact sectors): agriculture, biomes, fire, fisheries, health, lakes, water (resources), water quality, Groundwater, and terrestrial biodiversity. In our assessment, we consider nine indicators of



95 research software sustainability, five of them related to best practices in software engineering and four related to source code
quality. We further provide first-order cost estimates required to develop these GIMs. We also demonstrate how
reprogramming legacy software using best practices can lead to significant improvements in code quality and thus
sustainability. Finally, we offer actionable recommendations for developing sustainable research software for the scientific
community.

100 **2 Methods**

2.1 Accessing GIM Source code

ISIMIP manages a comprehensive database of participating impact models (available in an Excel file at
<https://www.isimip.org/impactmodels/download/>), which provides essential information such as model ownership, name,
source code links, and simulation rounds. Initially, we identified 375 models across five simulation rounds (fast track, 2a, 2b,
105 3a, and 3b). As the focus of our analysis is on global impact models, we sorted the data by spatial domain and filtered out
models operating at local and regional scales, resulting in a subset of 264 GIMs. We then removed duplicate models,
prioritizing the most recent versions for inclusion, resulting in 112 unique models. For models with available source links, we
obtained their source code directly. In instances where source links were not readily available, we conducted manual searches
for source code by referring to code availability sections in reference papers. Additionally, we searched for source code using
110 model names along with keywords such as "GitHub" and "GitLab" using the Google search engine. As of April 2024, 32
model source codes out of the 112 unique model source codes were accessible in the described way. However, it's important
to note that our sample may suffer from a "survivor bias," as we are not investigating models that are no longer in use (GIMs
that couldn't be sustained over time). This bias could potentially skew our analysis towards models that have survived i.e.,
they are still in use and their source code is accessible. Due to time constraints, we refrained from contacting developers for
115 models that were not immediately accessible.

2.2 Research software sustainability indicators

We examine nine indicators of research software sustainability, distinguishing five indicators related to the best practice in
software engineering and four indicators of source code quality (Table 1).

120

125



Table 1: Indicators used for the assessment of research software sustainability

No.	Indicator	Description
1	Documentation	Enables software use and also makes software maintenance easier (Wilson et al., 2014)
2	Version control	Provides transparency and traceability throughout the software development lifecycle and enables collaboration between developers as well as user communities (Wilson et al., 2014)
3	Use of an open-source license	Allows code copying and reuse. This openness fosters a collaborative environment where the user community can provide valuable feedback and support. Users can potentially contribute to the software's development and maintenance, enhancing its overall quality (Jiménez et al., 2017),
4	Number of active developers	Prevent single points of failure in the development process and make software development as well as maintenance easier (Long, 2006)
5	Containerization	Makes the software easy to install and facilitates reproducibility (Nüst et al., 2020; Wilson et al., 2014)
6	Public availability of an (automated) testing suite ¹	Shows that software functionality can be or was tested
7	Compliance with coding standards (eg. PEP 8) ¹	Improves code quality, readability and makes maintenance easier (Capiluppi et al., 2009; Simmons et al., 2020; Wang et al., 2008)
8	Comment density ¹	Precursor to software maintainability and re-usability (Arafat and Riehle, 2009; Stamelos et al., 2002; He, 2019)
9	Modularity ¹	Necessary for extensible and flexible research software (Stamelos et al., 2002; Sarkar et al., 2008).

¹ Indicators that impact research software quality



In the following, we describe the indicators and their rationale and how we evaluated the GIMs with respect to each indicator.

135 *Documentation.* Documentation is crucial for understanding and effectively utilizing software (Wilson et al., 2014). This
includes various materials such as manuals, guides, tutorials that explain the usage and functionality of the software as well
reference model description papers. When assessing documentation availability, relying solely on a reference model
description paper may be insufficient, as it may not provide the level of detail necessary for the effective utilization and
maintenance of the research software. All GIMs used in this assessment have an associated description or reference paper (see
140 supplementary file ISIMIP_models.xlsx). Therefore, in addition to the reference model paper we checked for available
manuals, guides, readme files, and tutorials. We consider any of these resources, alongside the reference model paper, as
documentation for the model. These resources provide essential information such as user, contributor, and troubleshooting
guides, which are valuable for model usage and maintenance. In our assessment, we searched within the source code and
official websites (if available). We also utilized the Google search engine to find model documentation by inputting model
145 names along with keywords such as 'documentation,' 'manuals,' 'readme,' 'guides,' and 'tutorials'.

Version control. Version control systems such as Git and Mercurial facilitate track changes, and collaborative development,
and provide a history of software evolution. To assess whether GIMs use version control for development, we focused on
commonly used open-source version control hosting repositories such as GitLab, GitHub, BitBucket, Google Code, and Source
150 Forge. The hostname such as “github” or “gitlab” in the source link of models provides clear indications of version control
adoption in their development process. For other models, we searched within the Google search engine using model names
and keywords such as “Bitbucket”, “Google Code”, and “Source Forge”.

Use of an Open source license. Open-source licenses foster collaboration and transparency by enabling community
155 contributions and ensuring that software remains freely accessible. We determined the existence of open-source licenses by
checking license files within repositories or official websites against Open source initiative (OSI) approved licenses
(<https://opensource.org/licenses>).

Number of active developers. The presence of multiple active developers serves as a safeguard against halts within the
160 development process. In instances where a sole developer departs or transitions roles, the absence of additional contributors
could lead to disruptions or challenges in maintaining and advancing the software. We measured the number of active
developers by counting the individuals who made commits or contributions to the projects codebase within the period 2020-
2024. A higher number of developers indicates a greater capacity for bug review (enhancing source code quality) and code
maintenance. It can also lead to more frequent updates to the source code. On the other hand, the absence of active developers
165 suggests potential stagnation in software evolution, possibly impacting the relevance and usability of the software.



170 *Containerization.* Containerization provides convenient ways to package and distribute software, facilitating reproducibility and deployment. It encapsulates an application along with its environment, ensuring consistent operation across various platforms (Nüst et al., 2020). Despite these benefits, containerization in high-performance computing systems encounters challenges like performance, prompting the proposal of solutions (Zhou et al., 2023). Some popular containerization solutions include Docker (<https://www.docker.com/>) and Singularity (<https://sylabs.io/>). To evaluate the availability of container solutions, we conducted searches through reference papers, official websites, and software documentation for links to container images or image-building files such as “Dockerfiles”, and “singularity definition file (.def file)”. In addition, we also searched through source code repositories to identify the previous stated images or image-building files. Lastly, we utilize the Google search engine, inputting the name of the GIM, the sector, and keywords such as “containerization”, to ascertain if any other containerized solutions exist.

180 *Public availability of an (automated) testing suite.* Test coverage, which verifies the software’s functionality, is the property of actual interest. However, research software may have an automatic testing suite but not provide information on test coverage or test results. As a practical approach, we consider the availability of a testing suite as a proxy for the ability to test software functionality. By examining testing suites within repositories, we gain insights into the developers’ commitment to software testing, which contributes to enhancing software quality.

185 *Compliance with coding standards.* Coding standards are a set of industry-recognized best practices that provide guidelines for developing software code (Wang et al., 2008). Analysing the conformance to these standards can be complex, particularly when the source code is written in multiple languages. As an example analysis, we focused on GIMs containing Python in their source code as it is one of the most prevalent languages used in development. The tool used, known as Pylint, is designed to analyze Python code for potential errors and adherence to coding standards (Obermüller et al., 2021; Molnar et al., 2020). Pylint evaluates source files for their compliance with PEP8 conventions. To quantify adherence to this coding standard, it assigns a maximum score of 10 as perfect compliance but has no lower bound (Molnar et al., 2020). We consider scores below 6 as indicative of weak compliance as code contain several violations.

Comment density. Good commenting practice is valuable for code comprehension and debugging. Comment density is an indicator of maintainable software (Arafat and Riehle, 2009; He, 2019). Comment density is defined as

195
$$\text{Comment density} = \frac{\text{Number of lines of comment}}{\text{Total lines of code}} \quad (1)$$

Here, the total lines of code (TLOC) include both comments and source lines of code (SLOC) (SLOCCount, 2024). SLOC is defined as the physical non-blank, non-comment line in a source file. According to Arafat et al. (2009) and He (2019), the optimal comment density is 30-60% (Arafat and Riehle, 2009; He, 2019). For most programming languages, this range is



200 considered to represent a compromise between providing sufficient comments for code explanation and having too many comments that may distract from the code logic (Arafat and Riehle, 2009; He, 2019).

Modularity. Researchers typically pursue new knowledge by asking and then attempting to answer new research questions. When the questions can be answered via computation, this requires either building new software, adding new source code, or modifying existing source code. Addition and modification of source code are more easily achieved if the software has a modular structure that is implemented as extensible and flexible software (McConnell, 2004). Therefore, modularity is chosen as another indicator for research software sustainability. Modular programming is an approach where source codes are organised into smaller and well-manageable units (modules) that execute one aspect of the software functionality, such as the computation of evapotranspiration in a hydrological model (Sarkar et al., 2008; Trisovic et al., 2022). The aim is that each module can be easily understood, modified, and reused. Depending on the programming language, a module can be a single file (e.g. Python) or a set of files (e.g. C++).

To assess the modularity of research software, we use the TLOC per file as a metric. This metric reflects the organization of the source code into modules, each performing a specific function (Sarkar et al., 2008; Trisovic et al., 2022). We opted for this approach over measuring TLOC per function or subroutine due to variations in programming languages and the challenges associated with accurately measuring different functions using program-specific tools. For instance, in Python, a module that contains significantly more TLOC than usual (here over 1,000 TLOC) likely includes multiple functions. These functions may perform more than one aspect of the software's functionality, such as reading input files and computing other functions (e.g. evapotranspiration function), which contradicts the principle of modularity. Keeping the length of code in each file concise also enhances readability.

The ideal number of TLOC per file can vary with the language, paradigm (e.g., procedural or object-oriented), and coding style used in a software project (Fowler, 2019; McConnell, 2004). However, a common heuristic is to keep the code size per file under 1,000 lines to prevent potential performance issues such as crashes or slow program execution with some integrated development environments (IDEs) (Fowler, 2019; McConnell, 2004). IDEs are software applications that provide tools like code editors, debuggers, and build automation tools. As reported by Trisovic et al. (2022), based on interviews with top software engineers, a module with a single file should contain at least 10 lines of code, consisting of either functions or statements (Trisovic et al., 2022). We used this heuristic as a criterion for good modularity, assuming that 10-1,000 TLOC per file indicates adequate modularity. We also varied the upper bounds of the total lines of code to 5,000 and 500 to investigate how modularity changes across models and sectors.

2.3 Source code counter

To count SLOC, comment lines, and TLOC of computational models, the counting tool developed by Ben Boyter (<https://github.com/boyter50/scc>) was used (Sloc, Cloc, and Code, 2024). This tool builds on the industrial standard source code counter tool called SLOCCount (Source Lines of Code Count) (SLOCCount, 2024).



2.4 Software cost estimation

The cost of developing research software is mostly unknown and depends on many factors, such as project size, computing infrastructure, and developer experience (Boehm, 1981). A model that attempts to estimate the cost of software development is the widely used Constructive Cost Model (COCOMO) (Boehm, 1981; Sachan et al., 2016), which computes the cost of commercial software by deriving the person-months required for developing the code based on the lines of code. Sachan et al. (2016) used the TLOC and effort estimates of 18 very large NASA projects (Average TLOC = 35,000) to optimise the parameters of the COCOMO regression model (Sachan et al., 2016). Effort in person months is estimated following Eq. (2):

$$Effort = 2.022817(kTLOC)^{0.897183} \quad (2)$$

where total lines of code are expressed in 1,000 TLOC (kTLOC) (Sachan et al., 2016). We use this cost model to estimate the cost of GIMs.

3 Results and Discussion

3.1 GIM programming languages and access points

The source code of the 32 GIMs is written in 10 programming languages (Fig. 1a). Fortran and Python are the most widely used, with 11 and 10 models, respectively. The dominance of Fortran stems from its performance, and the fact that it is one of the oldest programming languages designed for scientific computing (Van Snyder, 2007), and was the main such language used at the time some of the GIMs were originally built. This specialization makes it particularly suitable for tasks involving numerical simulations and complex computations. On the other hand, Python enjoys popularity among model developers due to readability, large user community, and rich ecosystem of packages, including those supporting parallel computing. R, C++ and C follow with 5, 5, and 4 models respectively (Fig. 1a). GIMs may employ one or more programming languages to target specific benefits the programming languages offer, such as readability and performance. For example, one of the studied models, HydroPy, written in Python, enhances its runtime performance by integrating a routing scheme built in Fortran (Stacke and Hagemann, 2021).

24 (75%) of the readily accessible 32 GIMs were hosted on GitHub (Fig. 1b). The rest are made available on GitLab (2, or 6%), Zenodo (4, or 12%), or the official website of the model (2, or 6%) (Fig. 1b, see supplementary file ISIMIP_models.xlsx). We note that for one of the GIMs used for analysis, WaterGAP2.2e, only part of the complete model (the global hydrology model) was accessed (Müller Schmied et al., 2021). This might be the case for other models as well.

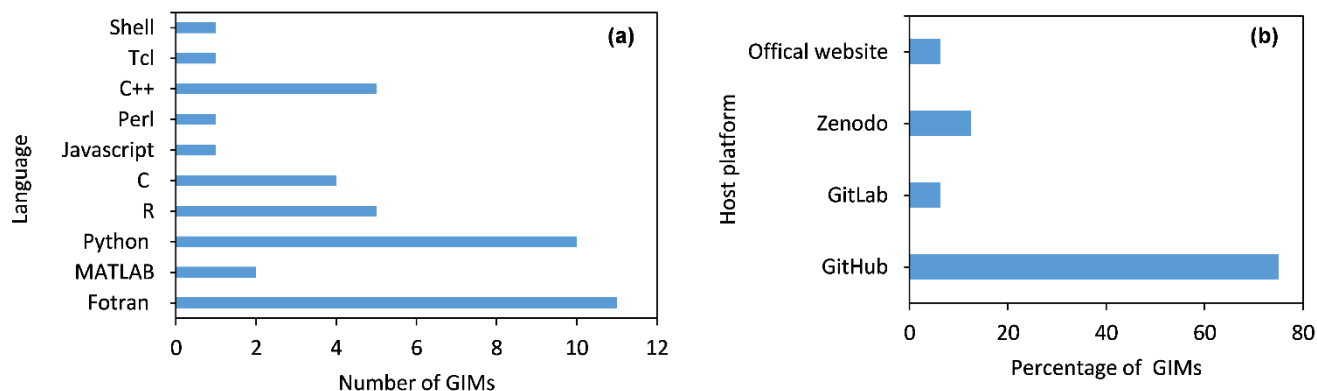


Figure 1: Programming languages for model development and model accessibility. (a) Bar plots showing programming languages used for developing 32 global impact models. (b) Bar plot showing open-source hosting platforms where 32 global impact models were accessed

265

3.2 Indicators of Software Sustainability

3.2.1 Software Engineering Practices

Documentation:

Our analysis reveals that 75% of the GIMs (24 out of 32) have publicly accessible documentation (Table 2). We observed a range of documentation formats across these GIMs. Specifically, 6 GIMs provided readme files, 13 had dedicated webpages for documentation, and 5 included comprehensive manuals (see supplementary file ISIMIP_models.xlsx). This prevalence of documentation practices among most models underscores the importance of documenting research software. However, a notable portion (25%) of the studied models either lack documentation or documentation has not been made publicly available (Table 2).

275

280



285 **Table 2:** Availability of Documentation, Version Control, Open-Source License, Test Suite, and Container for 32 Global Impact Models across 10 Sectors in Earth System Science. ‘x’, ‘-’, ‘not valid’ and ‘no info’ represent the availability, unavailability, not OSI-approved and absence of information, respectively.

No.	Sector	Model	Documentation	Version control	Open Source License	Test Suite	Container
1	Agriculture	CGMS-WOFOST	x	x	x	-	-
2	Agriculture	DSSAT-Pythia	x	x	no info	x	x
3	Agriculture	EPIC-TAMU	x	no info	x	-	-
4	Agriculture	LPJmL	x	x	x	-	-
5	Agriculture	ACEA	x	no info	not valid	-	-
6	Agriculture	LPJ-GUESS	x	no info	x	-	-
7	Biomes	CLASSIC	x	x	x	x	x
8	Biomes	MC2-USFS-r87g5c1	x	x	x	-	-
9	Fire	SSiB4/TRIFFID-Fire	-	x	no info	-	-
10	Fisheries	BOATS	-	x	no info	-	-
11	Fisheries	DBPM	-	x	no info	x	-
12	Fisheries	EcoTroph	x	x	no info	-	-
13	Fisheries	FEISTY	-	x	no info	-	-
14	Fisheries	ZooMSS	x	x	x	-	-
15	Groundwater	G ³ M	x	x	x	x	-
16	Groundwater	parflow	x	x	x	x	x
17	Lakes	ALBM	x	x	x	-	-
18	Lakes	GOTM	x	x	x	x	-
19	Lakes	SIMSTRAT-UoG	x	x	x	x	x
20	Terrestrial biodiversity	BioScen15-SDM-GAM/GBM	-	x	no info	-	-
21	Terrestrial biodiversity	BioScen1.5-MEM-GAM/GBM	-	x	x	-	-
22	Vector-borne diseases (health)	VECTRI	x	x	x	-	-
23	Water	CWatM	x	x	x	x	-
24	Water	DBH	x	no info	not valid	-	-
25	Water	HydroPy	x	no info	x	-	-
26	Water	PCR-GLOBWB	x	x	x	-	-
27	Water	WBM	x	x	x	-	-
28	Water	WaterGAP2.2e	-	no info	x	-	-
29	Water	VIC	x	x	x	x	x
30	Water	H08	x	x	x	-	-
31	Water	WAYS	-	x	x	-	-
32	Water quality	DynQual	x	x	no info	-	-
		Total	24	26	22	9	5



Version control:

290 We find that 81% (26 out of 32) of GIMs uses Git as their version control system reflecting the widespread acceptance of Git across the sectors (Table 2). In the remaining cases, GIMs were made available on Zenodo and the models' official websites (Table 2, Fig. 1b); information about the specific version control system used for these GIMs was unavailable. Developers' preference for Git highlights its user-friendly nature and effectiveness in supporting collaborative efforts.

Use of an open source license:

295 Most of the research software, 69% (22 out of 32), have open-source licenses (Table 2) with the “GNU General Public License” being the commonly used license (56%, 18 out of 32) (Fig. 2). However, the remaining 31% (10 out of 32) either have no information on the license even though the source code is made publicly available (8 or 25% of GIMs) or uses license which is not OSI-approved (1 GIM each with creative commons license and user agreement) (Fig. 2). This ambiguity or absence of licensing details can deter potential users and contributors, as it raises uncertainties about the permissions and restrictions associated with the software.
300

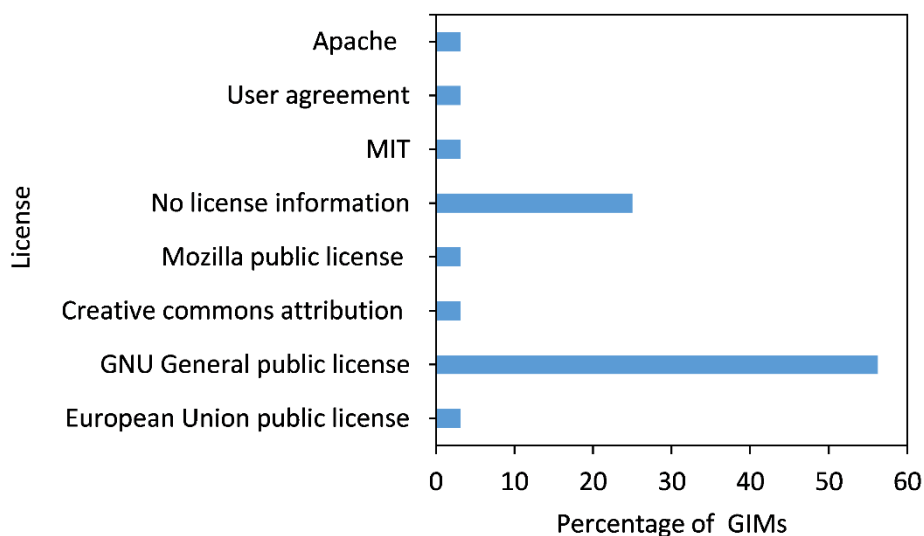


Figure 2: License distribution for 32 global impact models across 10 sectors. 8 (25%) GIMs lack license information, and two (6%) GIMs have licenses that are not OSI-approved.

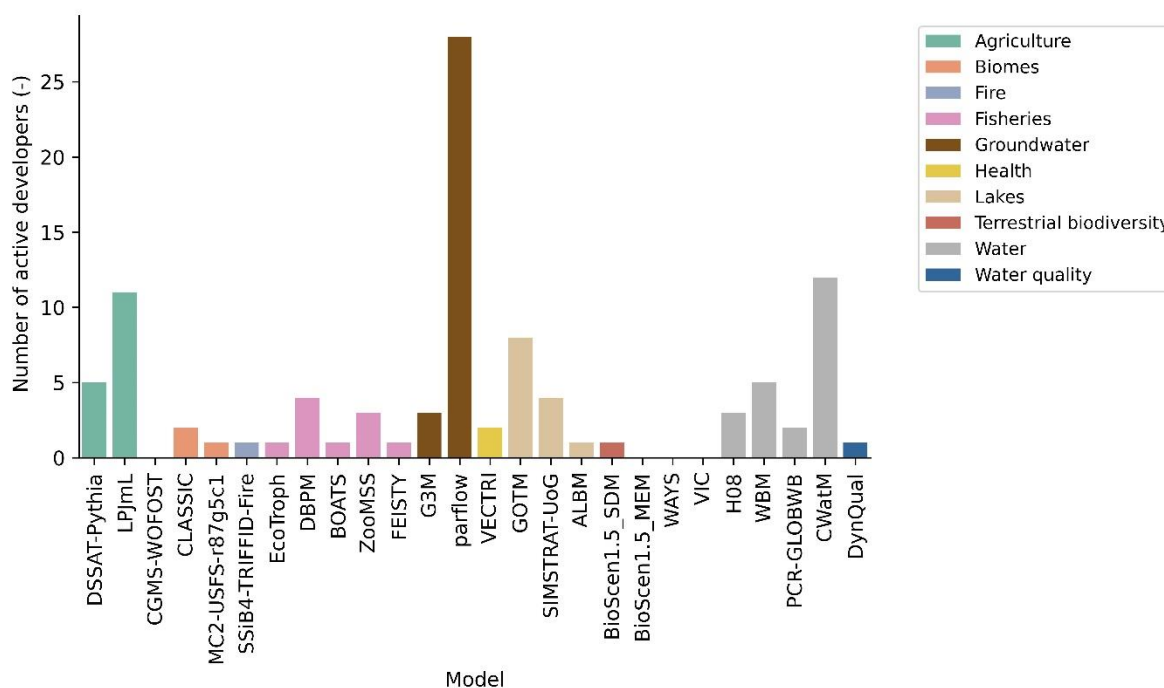
305

Number of active developers:

Our results reveal a diverse distribution of active developers across the GIMs. We have excluded GIMs without version control information from our results, as without could not be evaluated for this indicator, resulting in data for 26 GIMs. Notably, GIMs



such as parflow, CWatM, LPJmL, and GOTM have a significant number of active developers, with 28, 12, 11, and 8 developers respectively (Fig. 3). These values correlates with the size of GIMs source code, as evidenced by TLOC (282,722 for ParFlow, 33,286 for CWatM, 136,002 for LPJmL, and 29,477 for GOTM.). However, models such as WAYS, VIC, BioScen1.5-MEM, and CGMS-WOFOST had no active developers during the considered period of 2020 to 2024 (Fig. 3).



315 **Figure 3:** Number of active developers within 5 years (2020-2024) for 26 global impact models across 10 sectors. The results for the 6 remaining GIMs could not be measured since version control information could not be found. Zero value means no active developers within the 5 year period.

Containerization:

320 Only 5 (16%) of the GIMs have implemented containerized solutions (Table 2). Despite the recognized benefits of containerization in promoting reproducible research, provisioning of the software in containers is not yet a common practice in GIM development.



325 3.2.2 Code Quality Indicators

Public availability of an (automated) testing suite:

Our research indicates that 28% (9 out of 32) of the examined GIMs have a testing suite in place to test the software's functionality (Table 2). A typical test might involve ensuring that a global hydrological model such as CWatM runs without errors with different configuration file options (e.g., different resolutions and basins) (Burek et al., 2020). However, this practice is not widespread in the development of GIMs, with the majority (72%) lacking a testing suite (Table 2). This absence of testing suites in GIM development highlights a deficiency in the developers' dedication to software testing. The presence of a testing suite could lead to more frequent testing, thereby enhancing the overall quality of the software.

Compliance with coding standards:

335 We restricted our analysis to GIMs that include Python in their source code due to challenges described in section 2.2. Among the ten models we examined, we observed varying levels of adherence to the PEP8 style guide for Python. Five models (DSSAT-Pythia, parflow, HydroPy, VIC, and WAYS) demonstrated good compliance, each achieving a lint score above 6 out of a maximum of 10 (Fig. 4). Good compliance indicates minimal PEP8 code violations. However, the remaining five models showed lower compliance, with lint scores between 0 and 3 (Fig. 4). This suggests numerous violations leading to potential issues like poor code readability and an increased likelihood of bugs, which could hinder code maintenance.

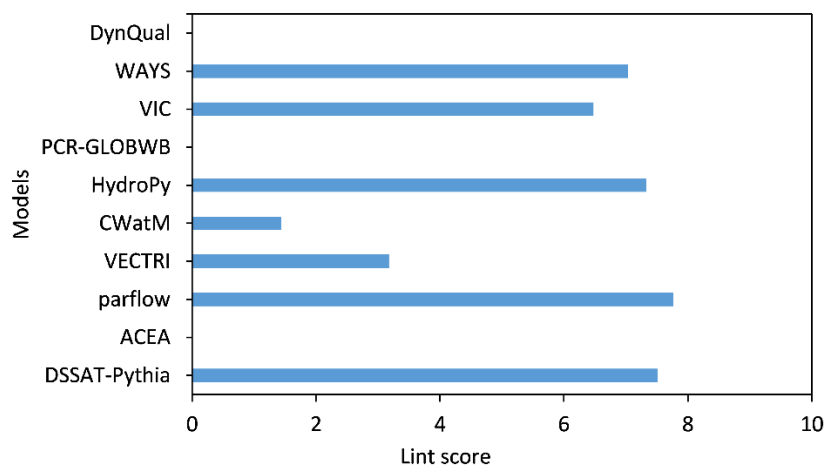


Figure 4: Lint scores of GIMs containing Python code.

345



Comment density:

25% (8 of 32) of the GIMs have well-commented source code, i.e. 30-60% of all source lines of code are comment lines (Fig. 5). The remaining 75% (24) of the GIMs have too few comments, which indicates that overall, commenting practice is low across the studied research fields.

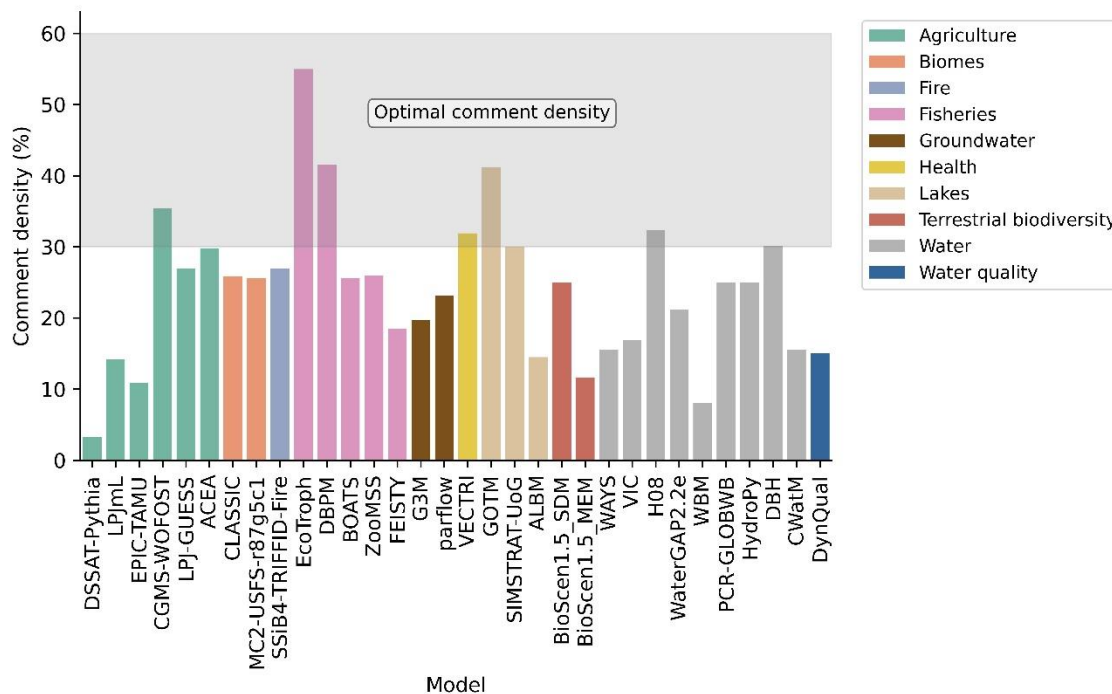


Figure 5: Comment density per model across 10 sectors. The grey zone denotes the optimal comment density (Arafat and Riehle, 2009; He, 2019).

355

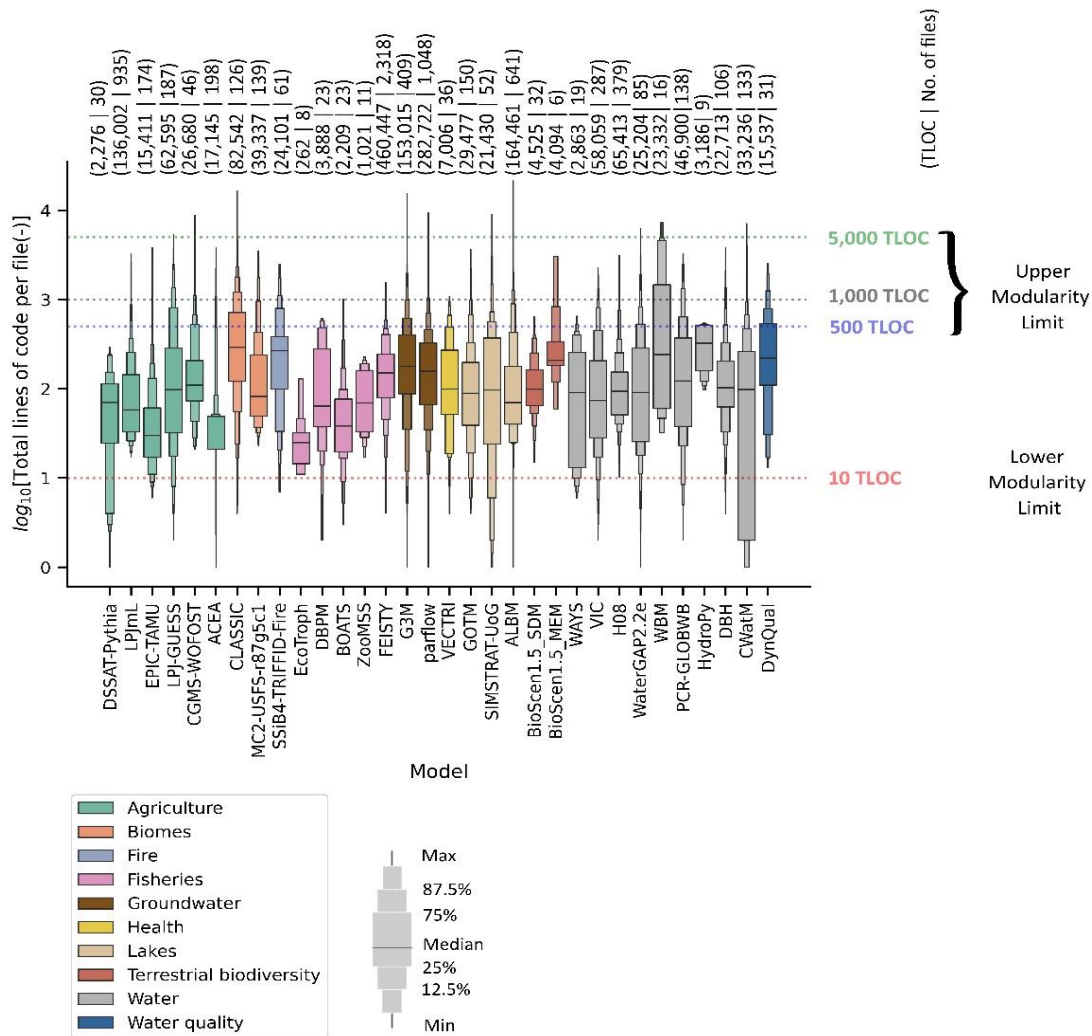
Modularity:

The investigated GIMs have TLOC values between 262 and 500,000, distributed over 6-2400 files (Fig. 6). Only 4 out of the 32 (12%) simulation models (EcoTroph, ZooMSS, HydroPy, and BioScen1.5_SDM) meet the criterion of having between 10 and 1,000 TLOC per file (Fig. 6). The remaining 28 GIMs either had at least one file exceeding 1,000 TLOC, which likely could be divided into smaller modules with distinct functionality or had at least one file less than 10 TLOC, which makes source code harder to navigate and understand, especially if the files are not well-named or documented. We also performed a sensitivity analysis by changing the criterion to 5,000 and 500 TLOC per file with the same lower limit of 10 TLOC. Nine simulation models (LPJmL, MC2-USFS-r87g5c1, EcoTroph, ZooMSS, BioScen1.5_SDM, BioScen1.5_MEM, H08, HydroPy, and DynQual) meet the 5,000-line criterion and two models (EcoTroph, ZooMSS) met the 500-line criterion (Fig. 6). Because code comments, which are included in TLOC, aid code comprehension, we also assessed modularity using the

365



criterion of 1,000 SLOC instead of 1,000 TLOC with 10 SLOC. Three GIMs (ZoomSS, BioScen1.5_SDM, and HydroPy) meet the 10-1,000 SLOC criterion (see supplementary Fig. S1).



370 **Figure 6:** Letter value plot (Hofmann et al., 2017) of total lines of code (TLOC) per file (logarithmic scale) of 32 global impact models across 10 sectors. The dotted blue, black, and green lines show upper modularity limits, the dotted red line the lower limit. The values (x|y) in the upper section of Fig. 6, show, for each GIM, TLOC | Number of files.



3.3 Cost of GIM software development

Research software is a valuable and complex research tool that often requires a lot of effort to develop and maintain (Carver et al., 2022; Reinecke et al., 2022). Here we use the cost estimate model from Sachan et al. (2016) (see section 2.4) in a scenario of “what if we would hire a commercial software company to develop the source code of the global impact models?” to provide a rough cost estimate for the software development of the 32 impact models. This cost estimate does not include developing the science (e.g., concepts, algorithms, and input data) nor costs of documenting, running, and maintaining the software, only the implementation of code. We assume that the COCOMO model is transferable to research software as the NASA projects used in cost model contain software that is similar to research software. As the TLOC of the impact model codes ranges from 262 to 500,000 TLOC (Fig. 7), the effort required to produce these models ranges from 1 to 495 person-months (Fig. 7).

The results suggest that these complex research software programs are expensive tools that require adequate funding for development and maintenance to make them sustainable. This is consistent with previous studies that have highlighted funding challenges for developing and maintaining sustainable research software in various domains (Carver et al., 2022; Reinecke et al., 2022; Carver et al., 2013; Merow et al., 2023; Eeuwijk et al., 2021).

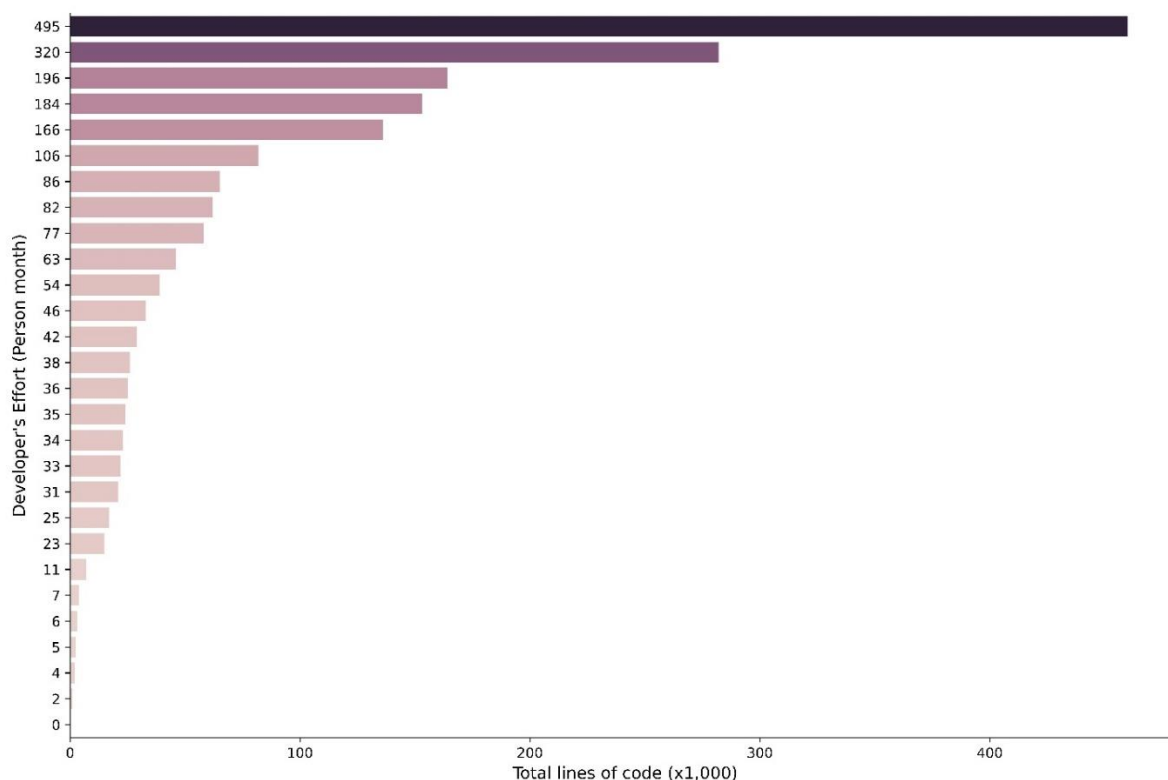


Figure 7 : Effort estimates of 32 global impact models across 10 sectors. Each bar represents one GIM. Darker colours represents large TLOC and effort values.



395 3.4 Case Study: Reprogramming legacy simulation models with best practices

Legacy codes often suffer from poor code readability and poor documentation, which hinder their maintenance, extension, and reuse. To overcome this problem, some of GIMs such as HydroPy (Stacke and Hagemann, 2021; Stacke, Tobias and Hagemann, Stefan, 2021) were reprogrammed, while others (e.g., WaterGAP, Nyenah et al., 2023) are in the process of being reprogrammed. We compared the legacy global hydrological model MPI-HM (in Fortran) and its reprogrammed version
400 HydroPy (in Python) in terms of the sustainability indicators. The reprogrammed model has improved modularity (Fig. 8a), which supports source code modification and extensibility. HydroPy has good compliance with the PEP8 coding standard, which improves readability and lower the likelihood of bugs in source code (Fig. 4). It has an open-source license and a persistent digital object identifier, which makes it easier to cite (Editorial, 2019). This research software refers to its associated publication for information and instructions on Zenodo to setup and run HydroPy. A software testing suite and container are
405 not yet available.

We find that HydroPy has a comment density of 25% (Fig. 8b), which is below the desired 30-60% range, but the developers argue that “the code is self-explanatory and comments are added only when necessary” (Stacke, 2023). MPI-HM has more comments (49%, Fig. 8b) because of its legacy Fortran code that limits variable names to a maximum length of 8 characters, so they have to be described in comments. Another reason is that the MPI-HM developers kept track of the file history in the
410 header, which adds to the comment lines in MPI-HM. This raises a question: *Is the comment density threshold metric still valid if a code is self-explanatory?* The need for comments can depend on the language’s readability (Python vs. Fortran), the complexity of the implemented algorithms and concepts, and the coder’s expertise. Nevertheless, comment density remains a valuable metric, especially for code written by novice developers.

Reprogramming legacy code not only allows developers to use more descriptive variable names, which increases code
415 readability and maintainability, but also enables them to share their code and documentation with the scientific community through open source platforms and tools. This practice enhances transparency and accountability, as the code can be inspected, verified, and reproduced by others. Reprogramming legacy code with best practices always improves code quality, which makes software more sustainable.

420

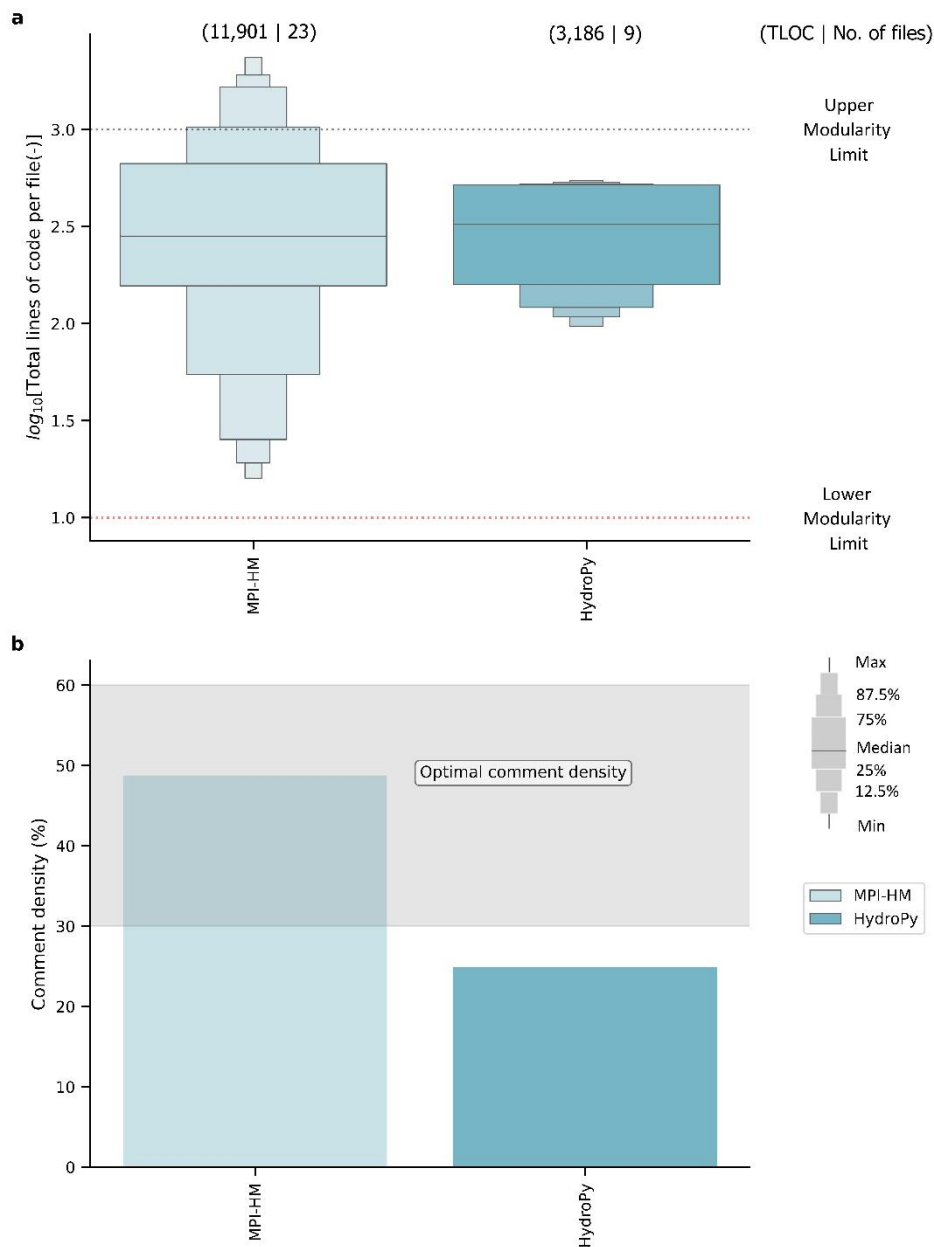


Figure 8: Modularity and commenting practice of a legacy (MPI-HM) and reprogrammed (HydroPy) global simulation model. (a) Letter value plot of total lines of code per file (logarithmic scale) of each model. The dotted black (red) line shows the upper (lower) modularity limit defined as the maximum of 1000 (minimum of 10) total lines of code per file. The values (x|y) shown in the upper section of Fig. 8a correspond to (TLOC | Number of files per model). (b) Comment density per model. The grey zone in Fig. 8b denotes the optimal comment density.



4 Limitations

430 Our study has limitations in the following regards. In the interest of timely analysis, developers were not contacted for models that were not readily available, causing a bias in the distribution of models. Specifically, the simulation model distribution does not favour certain sectors. For instance, only 2 out of the 18 global biomes impact models were readily available and therefore included in our assessment. This may affect the generalizability of our findings across different domains of Earth System Sciences.

435 Moreover, our sustainability indicators do not cover other relevant aspects of sustainable research software, such as user base size, code development activity (e.g. frequency of code contributions, date of last update or version), number of publications and citations, coupling and cohesion, information content of comments, software adaptability to user requirement and interoperability. A larger user base often results in more reported bugs, which ultimately enhances software reliability. However, determining the exact size of the user base presents challenges due to data reliability issues. Additionally, there is
440 the question of whether to include model output (data) users as part of the user base. Code development activity, such as the frequency of code contributions, indicates an ongoing commitment to improving and maintaining the software, but it does not necessarily reflect the quality of those contributions. In addition, the date of the last update or version is a useful metric, but it can be complex to interpret. For instance, research software might have an old last update date but still be widely used and reliable. Hence, these metrics were not evaluated here. The number of publications and citations referencing a model serves
445 as an indicator of its impact and relevance within the research community. Yet, collecting and analysing this data is a time-consuming and complex task. We further did not evaluate the interdependence of software modules (coupling) and how functions in a module work towards the purpose of the module (cohesion) (Sarkar et al., 2008), as language-specific tools are required to evaluate such properties.

In addition to the previously discussed limitations, the indicator analysed in this study are quantitative metrics that can be
450 measured. Factors such as information content of comments, software adaptability to user requirements and interoperability (Chue Hong et al., 2022) are examples of qualitative metrics that contribute to software sustainability. However, qualitative analysis is outside the scope of this study. We focus on measurable metrics that can be easily applied by the scientific community and by novice developers.

Also, we did not explore the analysis of code compliance to standards for other programming languages used for GIM
455 development. Specifically for Python, the Pylint tool provides a lint score for all source code analysed, making it easier to interpret results. However, the tools for other languages (e.g., linter for R) does not have this feature, which presents challenges in result interpretation.



5 Recommendations

460 Making our research software sustainable requires a combined effort of the modelling community, scientific publishers, funders, and academic and research organizations that employ modelling researchers (Barton et al., 2022; Barker et al., 2022; McKiernan et al., 2023; Research Software Alliance, 2023). Some scientific publishers, research organizations, funders and scientific communities adopted and proposed solutions to this challenge, such as 1) requiring that authors make source code and workflows available, 2) implementing FAIR standards, 3) providing training and certification programs in software engineering and reproducible computational research, 4) providing specific funding for sustainable software development, 5) 465 recognizing the scientific merit of sustainable research software, and 6) establishing the support of permanently employed research software engineers for disciplinary software developers (Carver et al., 2022; Eeuwijk et al., 2021; Editorial, 2018; Döll et al., 2023). In addition, we recommend the following actionable best practices for researchers developing software, based on literature and our own experience (summarized in Fig. 9):

- 470 • *Apply project management practices in software development (e.g., Agile):* This can help plan, organize, and monitor your software development process, as well as improve collaboration and communication within your team and with stakeholders. Project management practices can also help you identify and mitigate risks, handle changes, and deliver quality software on time and within budget (Anzt et al., 2021).
- 475 • *Consider software architecture (organisation of software components) and requirements (user needs):* This will help design your software in a way that meets the needs and expectations of your users. Considering software architecture (such as Model-Controller-View (Guaman et al., 2021)) and user requirements helps to design a software system that has a clear and coherent structure, well-defined functionality, and suitable quality (Jay and Haines, 2019).
- 480 • *Select an open-source license:* Choosing an open-source license will make your software accessible and open to the research community, enable collaborations with other developers and contributors, as well as protect your intellectual property rights (Carver et al., 2022; Anzt et al., 2021). Accessible software is crucial to reduce reliance on email requests (Barton et al., 2022).
- 485 • *Use version control:* Version control can help you track and manage changes to your source code and ensure your software is reproducible and traceable (Jiménez et al., 2017). Platforms like GitHub and GitLab are commonly used for this purpose. However, it's important to note that these platforms are not archival - the code can be removed by the developer at any time. A current best practice is to use both GitHub and GitLab for development, and to archive major releases on Zenodo or another archival repository.
- 490 • *Use coding standards (e.g., PEP8 for Python), good and consistent variable names, design principles, code quality metrics, peer code review, linters and software testing:* Coding standards help you write clear, consistent, and readable code that follows the best practices of your programming language and domain. Good variable names are descriptive and meaningful, reflecting the role and value of the variable. Design principles help adhere to the principles of



495

sustainable research software, such as modularity, reusability and interoperability. Code quality metrics can help measure and improve the quality of source code in terms of readability, maintainability, reliability, modularity and reusability. (Stamelos et al., 2002). Peer code review and linters (tools that analyse source code for potential errors) can help detect and fix errors, and vulnerabilities in your code, as well as improve your coding skills and knowledge (Jay and Haines, 2019). Software testing verifies if the research software performs as intended.

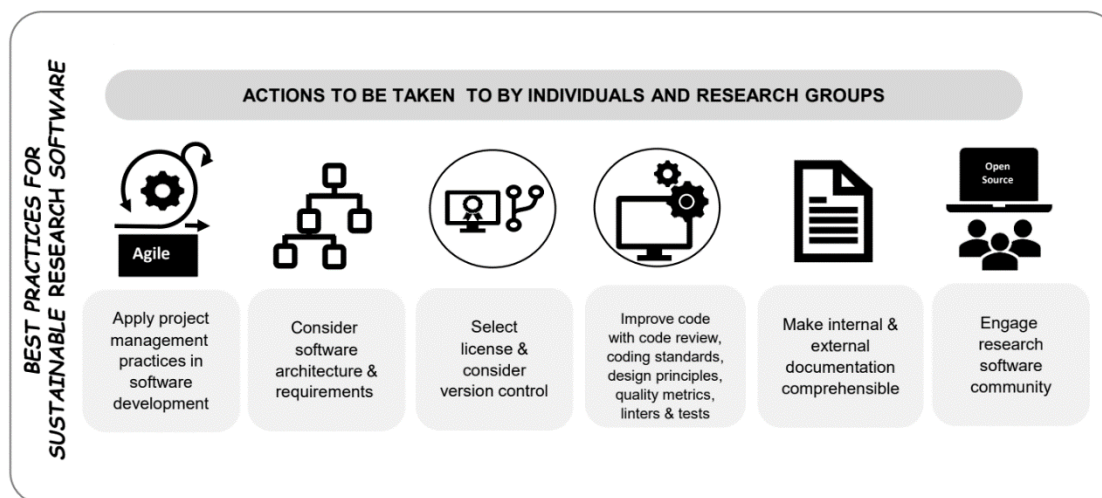
500

- *Make internal and external documentation comprehensible:* This can help you explain the purpose, functionality, structure, design, usage, installation, deployment, and maintenance of your software to yourself and others. Internal documentation refers to the comments and annotations within your code that describe what the code does and how it works. External documentation refers to manuals, guides, tutorials and any material that provide information about your software to users and developers. Comprehensible documentation can help you make your software more understandable, maintainable, and reusable. (Carver et al., 2022; Reinecke et al., 2022; Barker et al., 2022; Jay and Haines, 2019; Wilson et al., 2014)

505

- *Engage the research software community in the software development process.* This will help you get feedback, support, advice, collaboration, contribution and recognition from other researchers and developers who share your interests and goals. Engaging the research software community via conferences and workshops can also help you disseminate your software to a wider audience, increase its impact and visibility, and foster open science practices (Anzt et al., 2021). Additionally, consider utilizing containerization technologies, such as Docker, to simplify the installation and usage of your software. It helps eliminate the “it works on my machine” problem. This approach also facilitates easy sharing of your software with software users. Furthermore, implement continuous integration and automated testing to maintain the quality and reliability of your code. Continuous integration merges code changes from contributing developers frequently and automatically into a shared repository.

510



515 **Figure 9:** Actionable best practices for sustainable research software. The image summarizes the actions that modelling communities and individual developers should take, such as following project management practices, coding standards, reviews, documentation and community engagement strategies. These actions can help produce high-quality, robust, and reusable software that can be maintained.

6 Conclusion

520 The studied Earth system models are valuable and complex research tools that exhibit strengths and weaknesses in the use of certain software engineering practices (strengths, for example, in version control, open-source licensing, and documentation). However, notable areas remain for improvement, particularly in areas such as containerization and factors affecting code quality like comment density, modularity, and the availability of test suites. These shortcomings hinder the sustainability of such research software; they limit research reliability, reproducibility, collaboration, and scientific progress. To address this

525 challenge, we urge all stakeholders, such as scientific publishers, funders, as well as academic and research organizations, to facilitate the development and maintenance of sustainable research software. We also propose to use best practices for the developers of research software such as using project management and software design techniques, coding reviews, documentation, and community engagement strategies. We further suggest reprogramming the legacy code of well-established models. These practices can help achieve higher-quality code that is more understandable, reusable, and maintainable.

530

Efficient computational science requires high-quality software. While our study primarily focuses on Earth System Sciences, our assessment method and recommendations should be applicable to other scientific domains that employ complex research software. Future research could explore additional sustainability indicators, such as user base size, code development activity



(e.g. frequency of code contributions), software adaptability and interoperability, as well as code compliance standards for
535 various programming languages.

Code Availability

The Python scripts utilized for analysis can be accessed at <https://zenodo.org/doi/10.5281/zenodo.10245636>. Additionally, the line counting tool developed by Ben Boyter is available through the GitHub repository: <https://github.com/boyter/scc>.

Data Availability

540 The results obtained from the line count analysis are accessible at <https://zenodo.org/doi/10.5281/zenodo.10245636>. For convenient downloads of global impact models, links to the 32 global impact models, along with the respective dates of access, can be found in an Excel sheet named "ISIMIP_models.xlsx." present in the Zenodo repository.

Author contributions

EN and RR designed the study. EN performed the analysis and wrote the paper with significant contributions from PD, DK,
545 and RR. RR and PD supervised EN.

Competing interests

The authors declare no competing interests.

Acknowledgements

EN, RR, and PD acknowledge support from Deutsche Forschungsgemeinschaft (DFG) (Project number 443183317)

550 References

Alexander, K. and Easterbrook, S. M.: The software architecture of climate models: a graphical comparison of CMIP5 and EMICAR5 configurations, *Climate and Earth System Modeling*, <https://doi.org/10.5194/gmdd-8-351-2015>, 2015.

ISIMIP: <https://www.isimip.org/>, last access: 23 March 2024.

555 Anzt, H., Bach, F., Druskat, S., Löffler, F., Loewe, A., Renard, B., Seemann, G., Struck, A., Achhammer, E., Aggarwal, P., Appel, F., Bader, M., Brusch, L., Busse, C., Chourdakis, G., Dabrowski, P., Ebert, P., Flemisch, B., Friedl, S., Fritsch, B., Funk, M., Gast, V., Goth, F., Grad, J., Hegewald, J., Hermann, S., Hohmann, F., Janosch, S., Kutra, D., Linxweiler, J., Muth, T., Peters-Kottig, W., Rack, F., Raters, F., Rave, S., Reina, G., Reißig, M., Ropinski, T., Schaarschmidt, J., Seibold, H., Thiele,



- 560 J., Uekermann, B., Unger, S., and Weeber, R.: An environment for sustainable research software in Germany and beyond: current state, open challenges, and call for action [version 2; peer review: 2 approved], *F1000Research*, 9, <https://doi.org/10.12688/f1000research.23224.2>, 2021.
- Arafat, O. and Riehle, D.: The comment density of open source software code, in: 2009 31st International Conference on Software Engineering - Companion Volume, 195–198, <https://doi.org/10.1109/ICSE-COMPANION.2009.5070980>, 2009.
- 565 Azmi, E., Ehret, U., Weijs, S. V., Ruddell, B. L., and Perdigão, R. A. P.: Technical note: “Bit by bit”: a practical and general approach for evaluating model computational complexity vs. model performance, *Hydrology and Earth System Sciences*, 25, 1103–1115, <https://doi.org/10.5194/hess-25-1103-2021>, 2021.
- Barker, M., Chue Hong, N. P., Katz, D. S., Lamprecht, A.-L., Martinez-Ortiz, C., Psomopoulos, F., Harrow, J., Castro, L. J., Gruenpeter, M., Martinez, P. A., and Honeyman, T.: Introducing the FAIR Principles for research software, *Scientific Data*, 9, 622, <https://doi.org/10.1038/s41597-022-01710-x>, 2022.
- 570 Barton, C. M., Lee, A., Janssen, M. A., van der Leeuw, S., Tucker, G. E., Porter, C., Greenberg, J., Swantek, L., Frank, K., Chen, M., and Jagers, H. R. A.: How to make models more useful, *Proceedings of the National Academy of Sciences*, 119, e2202112119, <https://doi.org/10.1073/pnas.2202112119>, 2022.
- Boehm, B. W.: *Software engineering economics*, Prentice-Hall, Englewood Cliffs, N.J., 57–96, 1981.
- Sloc, Cloc, and Code: <https://github.com/boyter/scc.>, last access: 3 March 2024.
- 575 Burek, P., Satoh, Y., Kahil, T., Tang, T., Greve, P., Smilovic, M., Guillaumot, L., Zhao, F., and Wada, Y.: Development of the Community Water Model (CWatM v1.04) – a high-resolution hydrological model for global and regional assessment of integrated water resources management, *Geoscientific Model Development*, 13, 3267–3298, <https://doi.org/10.5194/gmd-13-3267-2020>, 2020.
- Capiluppi, A., Boldyreff, C., Beecher, K., and Adams, P. J.: Quality Factors and Coding Standards – a Comparison Between Open Source Forges, *Electronic Notes in Theoretical Computer Science*, 233, 89–103, <https://doi.org/10.1016/j.entcs.2009.02.063>, 2009.
- 580 Carver, J., Heaton, D., Hochstein, L., and Bartlett, R.: Self-Perceptions about Software Engineering: A Survey of Scientists and Engineers, *Comput. Sci. Eng.*, 15, 7–11, <https://doi.org/10.1109/MCSE.2013.12>, 2013.
- Carver, J. C., Weber, N., Ram, K., Gesing, S., and Katz, D. S.: A survey of the state of the practice for research software in the United States, *PeerJ Computer Science*, 8, e963, <https://doi.org/10.7717/peerj-cs.963>, 2022.
- 585 Chue Hong, N. P., Katz, D. S., Barker, M., Lamprecht, A.-L., Martinez, C., Psomopoulos, F. E., Harrow, J., Castro, L. J., Gruenpeter, M., Martinez, P. A., Honeyman, T., Struck, A., Lee, A., Loewe, A., van Werkhoven, B., Jones, C., Garijo, D., Plomp, E., Genova, F., Shanahan, H., Leng, J., Hellström, M., Sandström, M., Sinha, M., Kuzak, M., Herterich, P., Zhang, Q., Islam, S., Sansone, S.-A., Pollard, T., Atmojo, U. D., Williams, A., Czerniak, A., Niehues, A., Fouilloux, A. C., Desinghu, B., Goble, C., Richard, C., Gray, C., Erdmann, C., Nüst, D., Tartarini, D., Ranguelova, E., Anzt, H., Todorov, I., McNally, J., 590 Moldon, J., Burnett, J., Garrido-Sánchez, J., Belhajjame, K., Sesink, L., Hwang, L., Tovani-Palone, M. R., Wilkinson, M. D., Servillat, M., Liffers, M., Fox, M., Miljković, N., Lynch, N., Martinez Lavanchy, P., Gesing, S., Stevens, S., Martinez Cuesta, S., Peroni, S., Soiland-Reyes, S., Bakker, T., Rabemanantsoa, T., Sochat, V., Yehudi, Y., and WG, R. F.: FAIR Principles for Research Software (FAIR4RS Principles), <https://doi.org/10.15497/RDA00068>, 2022.
- 595 Döll, P., Sester, M., Feuerhake, U., Frahm, H., Fritzsche, B., Hezel, D. C., Kaus, B., Kolditz, O., Linxweiler, J., Müller Schmied, H., Nyenah, E., Risse, B., Schielein, U., Schlauch, T., Streck, T., and van den Oord, G.: Sustainable research software for high-



- quality computational research in the Earth System Sciences: Recommendations for universities, funders and the scientific community in Germany, <https://doi.org/10.23689/fidgeo-5805>, 2023.
- Editorial: Does your code stand up to scrutiny?, *Nature*, 555, 142–142, <https://doi.org/10.1038/d41586-018-02741-4>, 2018.
- Editorial: Giving software its due, *Nat Methods*, 16, 207–207, <https://doi.org/10.1038/s41592-019-0350-x>, 2019.
- 600 Eeuwijk, S. van, Bakker, T., Cruz, M., Sarkol, V., Vreede, B., Aben, B., Aerts, P., Coen, G., Dijk, B. van, Hinrich, P., Karvovskaya, L., Ruijter, M. K., Koster, J., Maassen, J., Roelofs, M., Rijnders, J., Schroten, A., Sesink, L., Togt, C. van der, Vinju, J., and Willigen, P. de: Research software sustainability in the Netherlands: Current practices and recommendations, Zenodo, <https://doi.org/10.5281/zenodo.4543569>, 2021.
- Fowler, M.: Refactoring, 2nd ed., Addison Wesley, Boston, MA, 2019.
- 605 Frieler, K. and Vega, I.: ISIMIP & ISIPedia - Inter-sectoral impact modeling and communication of national impact assessments, <https://unfccc.int/documents/197148>, 2019.
- Gamma, E., Helm, R., Johnson, R., and Vlissides, J.: Design patterns, Addison Wesley, Boston, MA, 1994.
- Guaman, D., Delgado, S., and Perez, J.: Classifying Model-View-Controller Software Applications Using Self-Organizing Maps, *IEEE Access*, 9, 45201–45229, <https://doi.org/10.1109/ACCESS.2021.3066348>, 2021.
- 610 Hannay, J. E., MacLeod, C., Singer, J., Langtangen, H. P., Pfahl, D., and Wilson, G.: How do scientists develop and use scientific software?, in: 2009 ICSE Workshop on Software Engineering for Computational Science and Engineering, 1–8, <https://doi.org/10.1109/SECSE.2009.5069155>, 2009.
- He, H.: Understanding Source Code Comments at Large-Scale, in: Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, New York, NY, USA, event-place: Tallinn, Estonia, 1217–1219, <https://doi.org/10.1145/3338906.3342494>, 2019.
- 615 Hofmann, H., Wickham, H., and Kafadar, K.: Letter-Value Plots: Boxplots for Large Data, *Journal of Computational and Graphical Statistics*, 26, 469–477, <https://doi.org/10.1080/10618600.2017.1305277>, 2017.
- Jay, C. and Haines, R.: Reproducible and Sustainable Research Software, in: *Web Accessibility: A Foundation for Research*, edited by: Yesilada, Y. and Harper, S., Springer, London, 211–221, https://doi.org/10.1007/978-1-4471-7440-0_12, 2019.
- 620 Jiménez, R. C., Kuzak, M., Alhamdoosh, M., Barker, M., Batut, B., Borg, M., Capella-Gutierrez, S., Hong, N. C., Cook, M., Corpas, M., Flannery, M., Garcia, L., Gelpí, J. L., Gladman, S., Goble, C., Ferreira, M. G., Gonzalez-Beltran, A., Griffin, P. C., Grüning, B., Hagberg, J., Holub, P., Hooft, R., Ison, J., Katz, D. S., Leskošek, B., Gómez, F. L., Oliveira, L. J., Mellor, D., Mosbergen, R., Mulder, N., Perez-Riverol, Y., Pergl, R., Pichler, H., Pope, B., Sanz, F., Schneider, M. V., Stodden, V., Suchecki, R., Vařeková, R. S., Talvik, H.-A., Todorov, I., Treloar, A., Tyagi, S., Gompel, M. van, Vaughan, D., Via, A., Wang, X., Watson-Haigh, N. S., and Crouch, S.: Four simple recommendations to encourage best practices in research software, <https://doi.org/10.12688/f1000research.11407.1>, 13 June 2017.
- 625 Katz, D. S.: Research Software: Challenges & Actions. The Future of Research Software: International Funders Workshop, Amsterdam, Netherlands., <https://doi.org/10.5281/zenodo.7295423>, 2022.
- 630 Kemp, L., Xu, C., Depledge, J., Ebi, K. L., Gibbins, G., Kohler, T. A., Rockström, J., Scheffer, M., Schellnhuber, H. J., Steffen, W., and Lenton, T. M.: Climate Endgame: Exploring catastrophic climate change scenarios, *Proceedings of the National Academy of Sciences*, 119, e2108146119, <https://doi.org/10.1073/pnas.2108146119>, 2022.



- Long, J.: Understanding the Role of Core Developers in Open Source Development, *Journal of Information, Information Technology, and Organizations (Years 1-3)*, 1, 075–085, 2006.
- McConnell, S.: in: *Code Complete, Second Edition*, Microsoft Press, USA, 565–596, 2004.
- 635 McKiernan, E. C., Barba, L., Bourne, P. E., Carter, C., Chandler, Z., Choudhury, S., Jacobs, S., Katz, D. S., Lieggi, S., Plale, B., and Tananbaum, G.: Policy recommendations to ensure that research software is openly accessible and reusable, *PLOS Biology*, 21, 1–4, <https://doi.org/10.1371/journal.pbio.3002204>, 2023.
- Merow, C., Boyle, B., Enquist, B. J., Feng, X., Kass, J. M., Maitner, B. S., McGill, B., Owens, H., Park, D. S., Paz, A., Pinilla-Buitrago, G. E., Urban, M. C., Varela, S., and Wilson, A. M.: Better incentives are needed to reward academic software development, *Nat Ecol Evol*, 7, 626–627, <https://doi.org/10.1038/s41559-023-02008-w>, 2023.
- 640 Molnar, A.-J., Motogna, S., and Vlad, C.: Using static analysis tools to assist student project evaluation, in: *Proceedings of the 2nd ACM SIGSOFT International Workshop on Education through Advanced Software Engineering and Artificial Intelligence, ESEC/FSE '20: 28th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, Virtual USA*, 7–12, <https://doi.org/10.1145/3412453.3423195>, 2020.
- 645 Nüst, D., Sochat, V., Marwick, B., Eglén, S. J., Head, T., Hirst, T., and Evans, B. D.: Ten simple rules for writing Dockerfiles for reproducible data science, *PLoS Comput Biol*, 16, e1008316, <https://doi.org/10.1371/journal.pcbi.1008316>, 2020.
- Nyenah, E., Reinecke, R., and Döll, P.: Towards a sustainable utilization of the global hydrological research software WaterGAP, *pico*, <https://doi.org/10.5194/egusphere-egu23-4453>, 2023.
- Obermüller, F., Bloch, L., Greifenstein, L., Heuer, U., and Fraser, G.: Code Perfumes: Reporting Good Code to Encourage Learners, in: *The 16th Workshop in Primary and Secondary Computing Education, WiPSCE '21: The 16th Workshop in Primary and Secondary Computing Education, Virtual Event Germany*, 1–10, <https://doi.org/10.1145/3481312.3481346>, 2021.
- 650 Prinn, R. G.: Development and application of earth system models, *Proceedings of the National Academy of Sciences*, 110, 3673–3680, <https://doi.org/10.1073/pnas.1107470109>, 2013.
- Reinecke, R., Trautmann, T., Wagener, T., and Schüler, K.: The critical need to foster computational reproducibility, *Environmental Research Letters*, 17, <https://doi.org/10.1088/1748-9326/ac5cf8>, 2022.
- 655 Research Software Alliance: *Amsterdam Declaration on Funding Research Software Sustainability*, , <https://doi.org/10.5281/ZENODO.8325436>, 2023.
- Sachan, R. K., Nigam, A., Singh, A., Singh, S., Choudhary, M., Tiwari, A., and Kushwaha, D. S.: Optimizing Basic COCOMO Model Using Simplified Genetic Algorithm, *Procedia Computer Science*, 89, 492–498, <https://doi.org/10.1016/j.procs.2016.06.107>, 2016.
- 660 Sarkar, S., Kak, A. C., and Rama, G. M.: Metrics for Measuring the Quality of Modularization of Large-Scale Object-Oriented Software, *IEEE Trans. Software Eng.*, 34, 700–720, <https://doi.org/10.1109/TSE.2008.43>, 2008.
- Satoh, Y., Yoshimura, K., Pokhrel, Y., Kim, H., Shiogama, H., Yokohata, T., Hanasaki, N., Wada, Y., Burek, P., Byers, E., Schmied, H. M., Gerten, D., Ostberg, S., Gosling, S. N., Boulange, J. E. S., and Oki, T.: The timing of unprecedented hydrological drought under climate change, *Nature Communications*, 13, <https://doi.org/10.1038/s41467-022-30729-2>, 2022.
- 665



- Sauer, I. J., Reese, R., Otto, C., Geiger, T., Willner, S. N., Guillod, B. P., Bresch, D. N., and Frieler, K.: Climate signals in river flood damages emerge under sound regional disaggregation, *Nat Commun*, 12, 2128, <https://doi.org/10.1038/s41467-021-22153-9>, 2021.
- 670 Schmidhuber, J. and Tubiello, F. N.: Global food security under climate change, *Proceedings of the National Academy of Sciences*, 104, 19703–19708, <https://doi.org/10.1073/pnas.0701976104>, 2007.
- Simmons, A. J., Barnett, S., Rivera-Villicana, J., Bajaj, A., and Vasa, R.: A large-scale comparative analysis of Coding Standard conformance in Open-Source Data Science projects, in: *Proceedings of the 14th ACM / IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM), ESEM '20: ACM / IEEE International Symposium on Empirical Software Engineering and Measurement, Bari Italy, 1–11*,
675 <https://doi.org/10.1145/3382494.3410680>, 2020.
- Stacke, T.: Personal communication, 2023.
- Stacke, T. and Hagemann, S.: HydroPy (v1.0): a new global hydrology model written in Python, *Geoscientific Model Development*, 14, 7795–7816, <https://doi.org/10.5194/gmd-14-7795-2021>, 2021.
- 680 Stacke, Tobias and Hagemann, Stefan: Source code for the global hydrological model HydroPy, <https://doi.org/10.5281/zenodo.4541381>, 2021.
- Stamelos, I., Angelis, L., Oikonomou, A., and Bleris, G. L.: Code quality analysis in open source software development, *Information Systems Journal*, 12, 43–60, <https://doi.org/10.1046/j.1365-2575.2002.00117.x>, 2002.
- Trisovic, A., Lau, M. K., Pasquier, T., and Crosas, M.: A large-scale study on research code quality and execution, *Sci Data*, 9, 60, <https://doi.org/10.1038/s41597-022-01143-6>, 2022.
- 685 Van Snyder, W.: Scientific Programming in Fortran, *Scientific Programming*, 15, 3–8, <https://doi.org/10.1155/2007/930816>, 2007.
- Wagener, T., Gleeson, T., Coxon, G., Hartmann, A., Howden, N., Pianosi, F., Rahman, M., Rosolem, R., Stein, L., and Woods, R.: On doing hydrology with dragons: Realizing the value of perceptual models and knowledge accumulation, *WIREs Water*, 8, e1550, <https://doi.org/10.1002/wat2.1550>, 2021.
- 690 Wan, W., Döll, P., and Zheng, H.: Risk of Climate Change for Hydroelectricity Production in China Is Small but Significant Reductions Cannot Be Precluded for More Than a Third of the Installed Capacity, *Water Resources Research*, 58, e2022WR032380, <https://doi.org/10.1029/2022WR032380>, 2022.
- Wang, Y., Zheng, B., and Huang, H.: Complying with Coding Standards or Retaining Programming Style: A Quality Outlook at Source Code Level, *Journal of Software Engineering and Applications*, 1, 88–91, <https://doi.org/10.4236/jsea.2008.11013>,
695 2008.
- Warszawski, L., Frieler, K., Huber, V., Piontek, F., Serdeczny, O., and Schewe, J.: The Inter-Sectoral Impact Model Intercomparison Project (ISI-MIP): Project framework, *Proceedings of the National Academy of Sciences*, 111, 3228–3232, <https://doi.org/10.1073/pnas.1312330110>, 2014.
- SLOCCount: <https://stuff.mit.edu/iap/debian/solutions/sloccount-2.26/sloccount.html>, last access: 4 March 2024.



700 Wilson, G., Aruliah, D. A., Brown, C. T., Hong, N. P. C., Davis, M., Guy, R. T., Haddock, S. H. D., Huff, K. D., Mitchell, I. M., Plumbley, M. D., Waugh, B., White, E. P., and Wilson, P.: Best Practices for Scientific Computing, *PLOS Biology*, 12, e1001745, <https://doi.org/10.1371/journal.pbio.1001745>, 2014.

Zhou, N., Zhou, H., and Hoppe, D.: Containerisation for High Performance Computing Systems: Survey and Prospects, *IEEE Trans. Software Eng.*, 49, 2722–2740, <https://doi.org/10.1109/TSE.2022.3229221>, 2023.

705