



# Beyond ecosystem modeling: A roadmap to community cyberinfrastructure for ecological data-model integration

Istem Fer<sup>1</sup> | Anthony K. Gardella<sup>2,3</sup> | Alexey N. Shiklomanov<sup>4</sup> | Eleanor E. Campbell<sup>5</sup> | Elizabeth M. Cowdery<sup>2</sup> | Martin G. De Kauwe<sup>6,7,8</sup> | Ankur Desai<sup>9</sup> | Matthew J. Duveneck<sup>10</sup> | Joshua B. Fisher<sup>11</sup> | Katherine D. Haynes<sup>12</sup> | Forrest M. Hoffman<sup>13,14</sup> | Miriam R. Johnston<sup>15</sup> | Rob Kooper<sup>16</sup> | David S. LeBauer<sup>17</sup> | Joshua Mantooth<sup>18</sup> | William J. Parton<sup>19</sup> | Benjamin Poulter<sup>4</sup> | Tristan Quaife<sup>20</sup> | Ann Raiho<sup>21</sup> | Kevin Schaefer<sup>22</sup> | Shawn P. Serbin<sup>23</sup> | James Simkins<sup>24</sup> | Kevin R. Wilcox<sup>25</sup> | Toni Viskari<sup>1</sup> | Michael C. Dietze<sup>2</sup>

<sup>1</sup>Finnish Meteorological Institute, Helsinki, Finland

<sup>2</sup>Department of Earth and Environment, Boston University, Boston, MA, USA

<sup>3</sup>School for Environment and Sustainability, University of Michigan, Ann Arbor, MI, USA

<sup>4</sup>Biospheric Sciences Laboratory (618), NASA Goddard Space Flight Center, Greenbelt, MD, USA

<sup>5</sup>Earth Systems Research Center, University of New Hampshire, Durham, NH, USA

<sup>6</sup>ARC Centre of Excellence for Climate Extremes, Sydney, NSW, Australia

<sup>7</sup>Climate Change Research Centre, University of New South Wales, Sydney, NSW, Australia

<sup>8</sup>Evolution & Ecology Research Centre, University of New South Wales, Sydney, NSW, Australia

<sup>9</sup>Department of Atmospheric and Oceanic Sciences, University of Wisconsin-Madison, Madison, WI, USA

<sup>10</sup>Harvard Forest, Harvard University, Petersham, MA, USA

<sup>11</sup>Jet Propulsion Laboratory, California Institute of Technology, Pasadena, CA, USA

<sup>12</sup>Department of Atmospheric Science, Colorado State University, Fort Collins, CO, USA

<sup>13</sup>Computational Earth Sciences Group and Climate Change Science Institute, Oak Ridge National Laboratory, Oak Ridge, TN, USA

<sup>14</sup>Department of Civil and Environmental Engineering, University of Tennessee, Knoxville, TN, USA

<sup>15</sup>Department of Organismic and Evolutionary Biology, Harvard University, Cambridge, MA, USA

<sup>16</sup>NCSA (National Center for Supercomputing Applications), University of Illinois at Urbana Champaign, Urbana, IL, USA

<sup>17</sup>College of Agriculture and Life Sciences, University of Arizona, Tucson, AZ, USA

<sup>18</sup>The Fulton School at St. Albans, St. Albans, MO, USA

<sup>19</sup>Natural Resource Ecology Laboratory, Colorado State University, Fort Collins, CO, USA

<sup>20</sup>UK National Centre for Earth Observation and Department of Meteorology, University of Reading, Reading, UK

<sup>21</sup>Fish, Wildlife, and Conservation Biology Department, Colorado State University, Fort Collins, CO, USA

<sup>22</sup>National Snow and Ice Data Center, Cooperative Institute for Research in Environmental Sciences, University of Colorado, Boulder, CO, USA

<sup>23</sup>Brookhaven National Laboratory, Environmental and Climate Sciences Department, Upton, NY, USA

<sup>24</sup>University of Delaware, Newark, DE, USA

<sup>25</sup>Ecosystem Science and Management, University of Wyoming, Laramie, WY, USA

This is an open access article under the terms of the Creative Commons Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2020 The Authors. *Global Change Biology* published by John Wiley & Sons Ltd

**Correspondence**

Istem Fer, Finnish Meteorological Institute,  
P.O. Box 503, 00101 Helsinki, Finland.  
Email: istem.fer@fmi.fi

**Funding information**

Australian Research Council, Grant/Award Number: CE170100023 and DP190101823; Academy of Finland, Grant/Award Number: 297350 and 327214; National Science Foundation, Grant/Award Number: 1062204, 1062547, 1457897, 1458021 and 1261582; National Aeronautics and Space Administration, Grant/Award Number: 80NSSC17K0711; Energy Biosciences Institute; Amazon AWS; Business Finland, Grant/Award Number: 6905/31/2018; UK NERC; California Institute of Technology; Climate and Environmental Sciences Division; Department of Energy, Grant/Award Number: DE-SC0012704; NSW; Oak Ridge National Laboratory, Grant/Award Number: DE-AC05-00OR22725

**Abstract**

In an era of rapid global change, our ability to understand and predict Earth's natural systems is lagging behind our ability to monitor and measure changes in the biosphere. Bottlenecks to informing models with observations have reduced our capacity to fully exploit the growing volume and variety of available data. Here, we take a critical look at the information infrastructure that connects ecosystem modeling and measurement efforts, and propose a roadmap to community cyberinfrastructure development that can reduce the divisions between empirical research and modeling and accelerate the pace of discovery. A new era of data-model integration requires investment in accessible, scalable, and transparent tools that integrate the expertise of the whole community, including both modelers and empiricists. This roadmap focuses on five key opportunities for community tools: the underlying foundations of community cyberinfrastructure; data ingest; calibration of models to data; model-data benchmarking; and data assimilation and ecological forecasting. This community-driven approach is a key to meeting the pressing needs of science and society in the 21st century.

**KEYWORDS**

accessibility, benchmarking, community cyberinfrastructure, data, data assimilation, ecosystem models, interoperability, reproducibility

**1 | INTRODUCTION**

Kindled by rapid environmental change, the scientific community is deeply invested in understanding and predicting nature's dynamics (Dietze et al., 2018; Hanson & Walker, 2020; Rineau et al., 2019). Thankfully, recent decades have seen an explosion of environmental data globally that is being delivered to us faster than ever before (Farley et al., 2018; LaDeau et al., 2017; Reichstein et al., 2019; Schimel et al., 2019). Process-based ecosystem models play a critical role in translating data into mechanistic understanding, as they provide us with the ability to synthesize and reformulate knowledge across organizational, spatial, and temporal scales, and to generate testable predictions from alternative hypotheses (Fisher et al., 2014; Hanson & Walker, 2020; Medlyn et al., 2015). Despite having more data than ever before, we have not seen comparable progress in our capacity to forecast natural systems with process-based models (Bonan & Doney, 2018; Dietze et al., 2018; Lovenduski & Bonan, 2017). For example, model projections out to the year 2100 do not agree on whether terrestrial ecosystems will be a carbon sink or source in response to climate change, and these discrepancies have not changed despite years of apparent model improvement (Arora et al., 2020; Friedlingstein et al., 2006, 2014). Perhaps this is not unexpected: adding model complexity without being informed by data does not equate to improved predictions, new processes (e.g., nutrients) may increase realism but may undo previous calibrated performance unless calibration is renewed easily. Overall, it is not a simple task to evaluate multiple model ensembles, making conclusions about forecast capacity complicated (Herger et al., 2019; Lovenduski & Bonan, 2017). A new strategy is needed to approach challenges in

advancing our ecological understanding, reducing uncertainties, and integrating the disparate science communities of global change biology (Bonan & Doney, 2018; Dietze et al., 2018). The goal of this paper is to better characterize the bottlenecks that have obstructed the rates at which new information has been integrated into ecosystem models, and to lay out a roadmap to overcome these bottlenecks. While many of the examples here are focused on terrestrial ecosystem models, the principles highlighted are general across different systems and processes.

A more predictive global change science needs to be based on ecosystem models that capture important processes rather than merely reproducing patterns (Bonan & Doney, 2018; Lovenduski & Bonan, 2017; Medlyn et al., 2015). Modeling efforts should be geared toward generating hypotheses that are testable against data (Hanson & Walker, 2020). Most current modeling activities, however, are more likely to be informed by high-volume high-level observational data (e.g., landscape level biogeochemical fluxes) than experimental manipulations (Wieder et al., 2019) or studies focused on low-level process details (e.g., interactions between non-structural carbohydrate reserves, drought, and mortality; Keenan et al., 2013). This is in direct contrast with the incredibly diverse range of data generated by ecology as a discipline (Hanson & Walker, 2020). Until modeling tools become more accessible, new communities of model users who can expand model-based interpretation and hypothesis testing beyond its limited scope will be curbed by informatics bottlenecks that impede wider representation.

More importantly, current approaches in confronting models with data frequently fail to actively engage the non-modeler community, who often possess a more detailed understanding of processes and

study systems (Jeltsch et al., 2013; Seidl, 2017). This bottleneck not only impacts the pace and the quantity but also the quality of modeling efforts. The division between empirical and modeling research is further exacerbated by the current “uniqueness of models”; that is, each model comes with an idiosyncratic learning curve due to the lack of standards around model interfaces and operation. To restore the balance, we need to concurrently increase modeling literacy and lower the technical barrier for modeling activities (Seidl, 2017). This barrier, overall, hinders efforts to replicate findings, extend analyses to other models and locations, and routinely confront model-based hypotheses with data (Gil et al., 2016).

We argue that a major step toward reducing these model-data bottlenecks lies in the development and support of community-wide cyberinfrastructure: a computational environment where we can effortlessly operate on data, simulate natural phenomena, perform model evaluation, and interpret results (Dietze et al., 2013; Eyring et al., 2019; Gil et al., 2016; also see Appendix A for a glossary of terms). While the general idea is not new, their application has been limited in ecology. However, there are several converging initiatives that make it timely to reinvigorate efforts (see Appendices C and D for example initiatives and their overview, and Box 1).

### BOX 1 How to support and sustain community cyberinfrastructure?

The ongoing maintenance and development of common cyberinfrastructure tools are essentially conditioned upon uptake and support by the community. This effort typically starts with building a bottom-up community (Boettiger et al., 2015) involving:

- Support widely adopted languages by the domain scientists (e.g., R and Python) so that:
  - experienced users can get off to a running start,
  - inexperienced users would be motivated to invest efforts with the co-benefit of learning a popular language,
  - larger communities of these languages can bring further support.
- Initiate strong ties with the demographic that can highly benefit from community solutions such as early career researchers.
- Establish codes of conduct for inclusion and diversity, and encourage participation regardless of experience level.
- Always adhere to open software best practices to build a reputation that can in return attract human resources and funding.

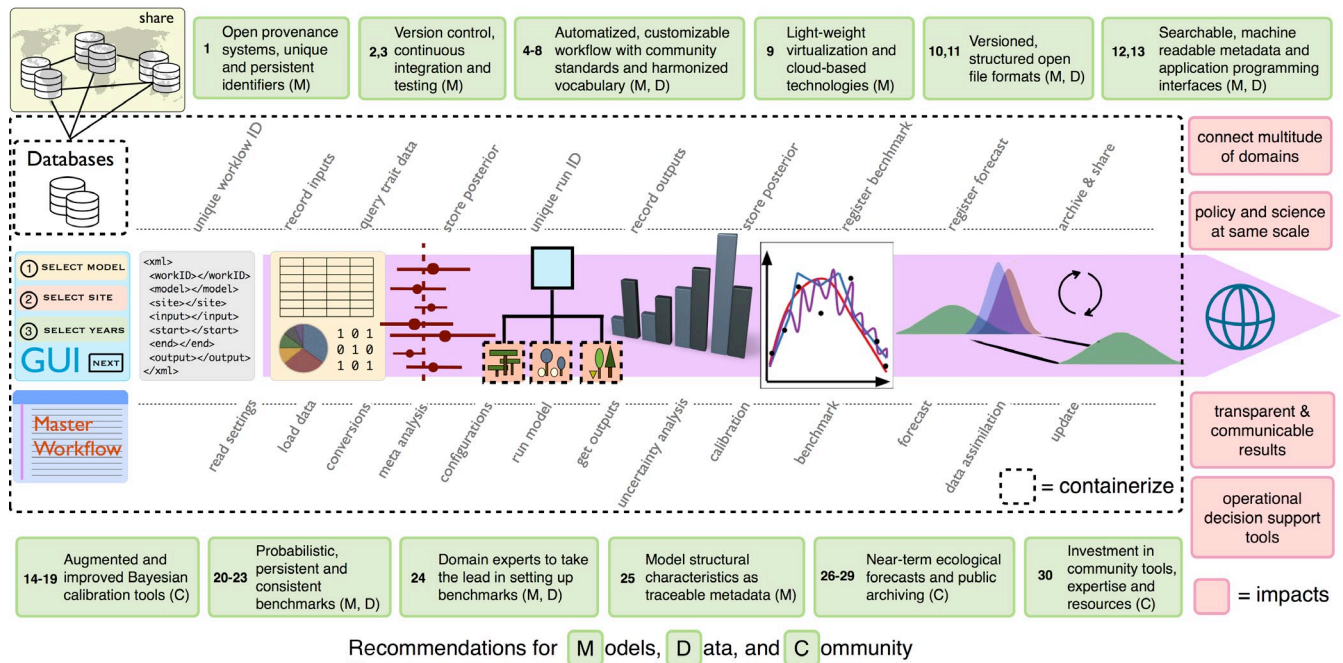
Luckily, these efforts do not need to start from scratch: the community can adopt and build upon existing systems (Appendix C). While we acknowledge that getting involved with community development requires upfront investment of time and resources of individuals, the benefits from participation are significant overall:

- Contributions to community tools perpetuate and increase their value, elevate recognition of their contributors (Dai et al., 2018; Lowndes et al., 2017).
- Community involvement provides larger support and career networks (McKiernan et al., 2016).
- In a research landscape that is ever diversifying, community cyberinfrastructure will be an active learning platform where ecologists gain advanced capability (Dietze et al., 2013).

As the community grows, successful strategies could be taken as an example, such as the WRF (The Weather Research and Forecasting Model) community (Powers et al., 2017):

- Financial and personnel burdens are spread out among the community, while the main support and steering responsibility could remain centralized.
- A help service that is responsible for user assistance is fundamental.
- Building committees in charge of coordination and direction is effective, e.g.:
  - Developers committee, to maintain code design, testing and upkeep.
  - Release committee, to oversee and time major releases.
  - Review committee, for scientific evaluation of major module/package contributions.

Open software and data management plans are increasingly becoming an important requirement by funding agencies (Powers & Hampton, 2019) for which use of community cyberinfrastructure could be fittingly proposed. Thus, we suggest such proposals to include a budget item or person hours for the support of community tools when possible. While projects without funding should also be welcome, short-term funding opportunities for open research (McKiernan et al., 2016; Powers & Hampton, 2019) will help bottom-up community building. However, viability over the long-term requires sustainable funding structures and top-down support from funding agencies, networks, and the private sector. There are currently several appropriate venues for cyberinfrastructure projects (e.g., NSF Cyberinfrastructure for Sustained Scientific Innovation), but as communities make their cyberinfrastructure needs better known (e.g., through communication with funding agencies and uptake), we expect such opportunities to increase in number and variety. Ultimately, [R30] it is important that community and funding agencies support the sustainability of these tools as critical components of the collective scientific infrastructure in a similar way they do with the physical infrastructure (field stations, sensor networks, satellites) and data repositories.



**FIGURE 1** Schematic of a community cyberinfrastructure example and summary of recommendations (numbers in the green boxes refer to our recommendations in the main text). Users start with a high-level Graphical User Interface (GUI) to provide their setup for a modeling activity. These selections are translated into a human and machine-readable markup language and read in by the master workflow which then executes a sequence of modularized tasks. At this stage, a unique identifier is assigned to the workflow to be executed. This ID, which points to the full workflow output and access to the metadata required to repeat it, can be shared among collaborators and published in papers. Next, the selections of the user are queried with the database, and actions are decided depending on whether requested items are already processed in an earlier modeling activity and ready to use or need to be retrieved and processed. Then, each module performs a well-defined task in the specified order. Crucial information for provenance of the whole workflow is recorded in the database during associated steps. Key outputs from analyses, such as calibration posteriors, are stored in a way that enables their exchange and re-use between different workflows. An important feature of this cyberinfrastructure is that both its parts and itself as a whole are virtualized (containerized) to add an additional layer of abstraction and automation, and to ensure interoperability

In the following sections, we present a roadmap to the key features of a community cyberinfrastructure, and discuss specific challenges and solutions for model-data activities. These activities include but are not limited to (a) obtaining and processing data (data ingest); (b) estimating model parameters through statistical comparisons between models and real-world observations (calibration); (c) evaluating and comparing performance skills through standardized and repeatable multi-model tests (evaluation and benchmarking); and (d) combining model predictions with multiple observations to update our understanding of the state of the system (data assimilation). We provide specific recommendations for the measurement community, the modeler and developer community, and the broader community throughout each section (Figure 1; Appendix B).

## 2 | FAIR CYBERINFRASTRUCTURE ESSENTIALS

There should be few things more repeatable in science than running a deterministic model. In practice, running a process-based simulation model is often fraught with roadblocks to any new user or developer (Dietze et al., 2013). Tackling this at the individual model level leads to redundant efforts across models and inhibits

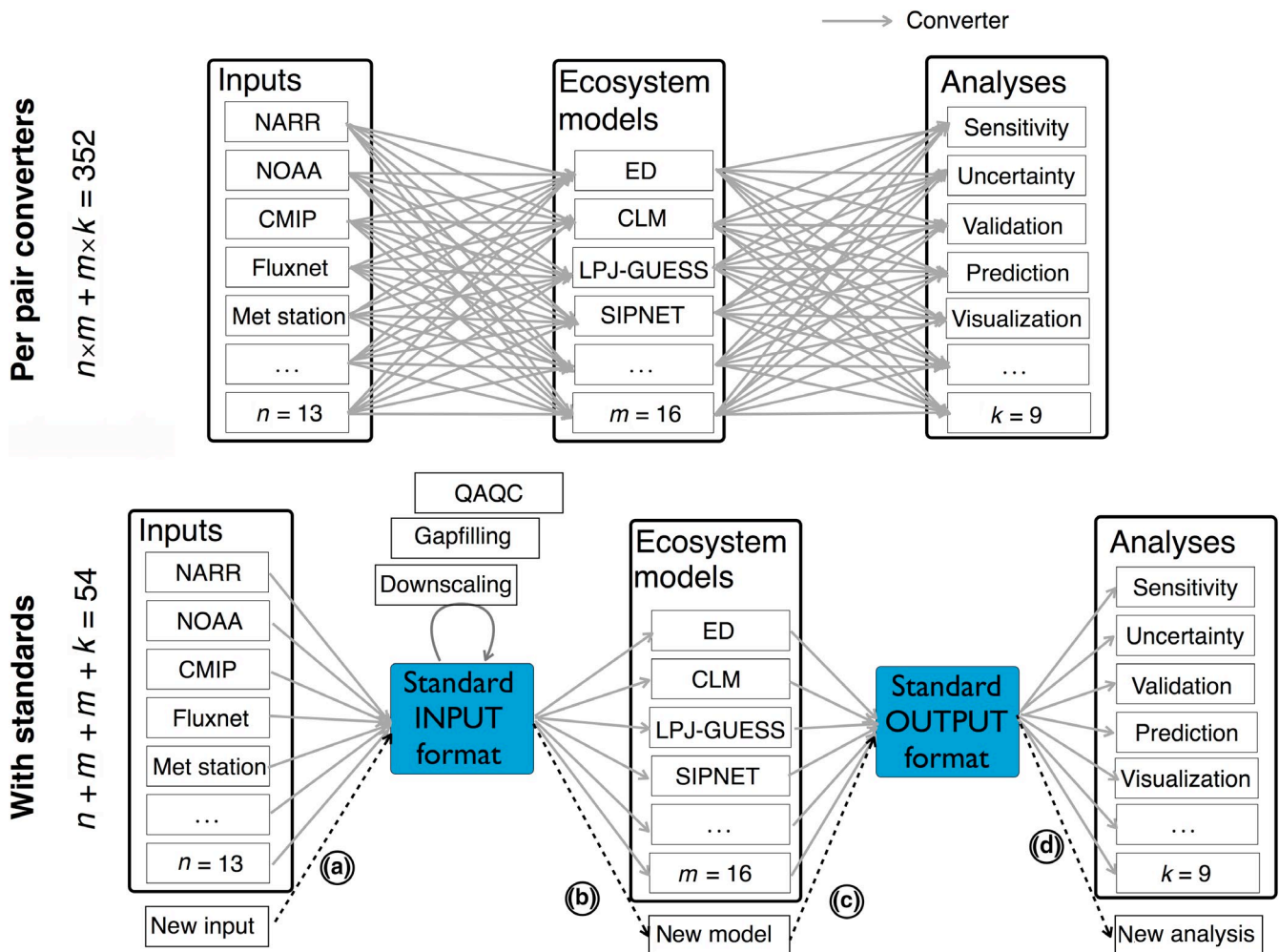
economies of scale that could be gained by sharing informatics tools across communities (for examples of shared ecological informatics infrastructure please see Appendix C). Besides, the larger community of users associated with common infrastructure will foster innovation and create an incentive for developers to make better, more sophisticated algorithms that have gone through more extensive testing (Gil et al., 2016). The revolutionary success of the open source and free programming language R (R Core Team, 2020) aptly exemplifies the importance of community involvement in developing and sharing standard tools for a massive reduction in redundant efforts, as well as having access to a much larger community support (Boettiger et al., 2015; Lai et al., 2019).

Here we briefly highlight the FAIR (*findable, accessible, interoperable, and reusable*) cyberinfrastructure essentials to facilitate a catalog of model-data activities (for more details on FAIR principles for research software and data, please see Culina et al., 2018; Gil et al., 2016; Hasselbring et al., 2020 and the references therein):

- **Findability** refers to the ease with which permanent records of the key metadata about each model-data activity and computational output can be found (Hasselbring et al., 2020). Recording the full, transparent history of an analysis to enable findability is

known as provenance. For large model-data workflows executing multiple models or experiments, we recommend **[R1; R for recommendation]** model developers utilize open community provenance databases, which assign unique and persistent identifiers to each model-data activity (Gil et al., 2016; LeBauer et al., 2013). Such identifiers could be used in publications, pointing readers to the full computational output and the metadata required to repeat a model run (Fer et al., 2018). **[R2]** The workflow and provenance system themselves should also be version controlled (e.g., using GitHub) to ensure a fully reproducible record (Piccolo & Frampton, 2016). **[R3]** Then, any changes to their code need to be automatically tested to ensure expected behavior by tools for continuous integration (e.g., Travis CI, [travis-ci.com](https://travis-ci.com); Github Actions, [github.com/features/actions](https://github.com/features/actions)).

- **Accessibility** in modeling goes beyond obtaining the model code. A broader technical barrier exists in terms of the abilities required to effectively deploy simulation models and perform complex analyses. **[R4]** A well-defined automated workflow that coordinates individual tasks (Figure 1) should be set up by the developers to (a) reduce barriers to entry; (b) ensure replication is possible; and (c) reduce costs of manual operation. The process of focusing on the design of this workflow, which is also known as abstraction, requires standardizing and generalizing the important tasks involved, and devising how they are related to one another. Leveraging systemized approaches (e.g., tidyverse in R, or pandas in Python) throughout the workflow design promotes consistency, creates predictable expectations, and fosters knowledge transfer across projects. Abstraction further facilitates presenting



**FIGURE 2** Reduction in redundant work when adopting common formats. There are " $n$ " data types that must be linked to " $m$ " simulation models and " $k$ " post-simulation analyses. In the top panel, the conventional approach where modeling teams work independently requires implementing  $n \times m$  different input and  $m \times k$  different output conversions. As data, models, and analyses are added, and effort scales quadratically. On the other hand, the bottom panel shows that by working as a community, and adopting common formats and shared analytical tools, the number of converters necessary to link models, data, and analyses reduces to an  $m + n$  and  $m + k$  problem, and scales linearly. When a new input source or a new analysis is added to the system, it can immediately get access to  $m$  models by writing only one converter, (a) and (d) respectively. Likewise, when a new model is added, it can get access to  $n$  inputs and  $k$  analyses by writing one converter for each, (b) and (c) respectively. This scaling also extends beyond data conversions to the development of tools and analyses. For example, if input data need to be extracted, downscaled, debiased, gap-filled, or have their uncertainties estimated, each of these steps does not need  $m \times n$  variants but rather just one tool that can be applied to the standard



the user with a [R5] more intuitive and accessible interface that handles everything from running ecosystem models in place to submitting complex analyses to remote high-performance computing resources under the hood.

- *Interoperability* is critical to building cyberinfrastructure that works seamlessly across many models, but this requires predictable file formats for model inputs, outputs, and data constraints used by the community. While reducing the proliferation of both data and model formats would alleviate this in the long term, in the short term [R6] using standard data pipelines can remedy the redundant efforts put into building custom tools. For example, consider the common problem of managing the data streams in and out of the models with two cases where (a) every developer team works independently (Figure 2, top panel); and (b) a common pipeline with internal standards is used (Figure 2, bottom panel). Not only is the latter approach much more scalable, but these tools can be made more reliable and sophisticated as less code will be written and tested by more people. [R7] We recommend the ecological community leverage existing standard formats as the internal standards, such as the Climate and Forecast convention (Eaton et al., 2017), and the use of ontologies to provide harmonized vocabularies and semantic frameworks (e.g., Stucky et al., 2018).
- *Reusability* of community models and tools builds on interoperability but also requires [R8] individual tasks involved be isolated and modularized in the workflow (Figure 1). Modularity would allow (a) internal modifications to their implementation without altering the overall behavior of the system; (b) independent reuse of tools outside of specific systems; and (c) users to swap in/out alternative algorithms/tools and customize their workflow. Community cyberinfrastructure should further be available to users without having to deal with obscure system requirements and dependencies. Similar to what programming language R has achieved, more standardized installation procedures and fewer configuration steps significantly reduce user time for setup and increase adoption, reusability, and reproducibility. Fortunately, modern virtualization technologies offer a number of tools that allow users to run packaged software, called containers, complete with all its dependencies (Piccolo & Frampton, 2016). [R9] We recommend developer communities adopt recent lightweight containerization systems (such as e.g., Docker—[www.docker.com](http://www.docker.com); Singularity—[singularity.lbl.gov](http://singularity.lbl.gov)) that are easy to install, set up, upgrade, and scale up with new locations to run the models. Containerization allows existing infrastructures to be run reliably across a variety of computing resources, including cloud-based virtual services (Farley et al., 2018; Hasselbring et al., 2020).

### 3 | DATA INGEST OPPORTUNITIES

Data play a critical role in modeling activities; however, due to their sheer volume and diversity, they can be difficult to locate and obtain as sifting through deluge of data manually is impractical (Reichstein

et al., 2019; Waide et al., 2017). [R10] To make data FAIR, we recommend data producers use consistent naming structures (e.g., Assistance for Land-surface Modelling activities convention, also please see Appendix A for more details) and open file formats (e.g., comma-separated values, netCDF; Hart et al., 2016). [R11] Next, data should be stored in data repositories where datasets are versioned, data citations are provided, and that support [R12] standard, searchable metadata, and machine-readable Application Programming Interfaces (e.g., the Oak Ridge National Laboratory Distributed Active Archive Center, Cook et al., 2016; Environmental Data Initiative, Gries et al., 2019; Open Science Framework, Sullivan et al., 2019). When those repositories are part of jointly searchable networks (e.g., DataONE—[www.dataone.org](http://www.dataone.org)), it could further allow developers to leverage one set of tools for many sources.

Admittedly, data providers may have to invest significant time and resources to follow these recommendations. These costs include the following: preparing descriptive metadata to prevent misuse, choosing the right repository with appropriate licensing and without isolating data from relevant disciplines, and finding means (funding and expertise) to manage data especially for small projects (Culina et al., 2018; Gil et al., 2016; Waide et al., 2017). Furthermore, other valid concerns such as data leakage and insufficient recognition are frequently raised (Bond-Lamberty et al., 2016). While these issues are not specific to the roadmap discussion here, community cyberinfrastructure tools can alleviate them to a certain extent. For example, investments in optimizing standardized protocols, terminologies and file formats for community tools during data collection and processing will help with metadata preparation and repository selection. By getting involved with community cyberinfrastructure, small projects can gain access to larger community expertise and support. Cyberinfrastructure data ingest pipelines can automatically query licenses as chosen by the data provider (Culina et al., 2018) and streamline citations to credit researchers seamlessly. Community tools (such as Brown Dog, [browndog.ncsa.illinois.edu](http://browndog.ncsa.illinois.edu)) can access and index data collections, in particular small uncured and/or unstructured data collections, thereby preventing data loss, increasing discovery, and further securing recognition.

On the big data side, approaches for scientifically and computationally interacting with high-volume, high-velocity data become increasingly available (Reichstein et al., 2019). While it is important to generalize these cutting-edge tools and share with the community, modeling activities frequently involve a subset of data (e.g., a specific region or period) for which time to transfer data often exceeds the time to process it. Thus, we endorse the recent paradigm of [R13] cloud computing and online services (e.g., Google Earth Engine) that allow users to select, subset, transform, or perform other operations on the data without having to download and expand (see Gomes et al., 2020 for more examples). Within this set up, community cyberinfrastructure also provides a medium where a diverse array of data delivered by Internet of Things techniques can be integrated into models in a sensible manner (Fang et al., 2014). As developers combine cloud-based cyberinfrastructure tools with cutting-edge data platforms, this would free the users from their local constraints

altogether. Empowering more groups to interact with large datasets brings its own push toward progress in terms of scientific proficiency and diversity (Nagaraj et al., 2020).

## 4 | WAY FORWARD IN CALIBRATION

After data ingest, another persistent challenge in process-based ecosystem modeling is calibration: the process of using data to constrain model parameters (Dietze et al., 2013; Seidel et al., 2018; van Oijen, 2017). Some model parameters may be directly informed by ecological trait data (e.g., turnover rates). In this case, meta-analysis tools can pull data together from open-access, machine-readable, curated databases (LeBauer et al., 2013, 2018; Shiklomanov et al., 2020). A non-negligible portion of model parameters, however, are often not directly measurable; therefore, there is a need to estimate parameters indirectly using inverse methods that infer what parameter combinations produce model predictions compatible with observations (Hartig et al., 2012). **[R14]** When doing this, we recommend the community take the Bayesian approach to transfer the information from data to probability distributions about models and parameters (Hartig et al., 2012; LeBauer et al., 2013). Bayesian approach allows combining information from multiple sources and scales, iteratively updating our understanding as new data become available, propagating uncertainty into model predictions to inform decision making, and it is becoming more effective in dealing with complex systems with the increase in computing power and numerical methods (van Oijen, 2017).

Most off-the-shelf Bayesian tools (e.g., JAGS—[mcmc-jags.sourceforge.net](http://mcmc-jags.sourceforge.net); STAN—[mc-stan.org](http://mc-stan.org)), however, are not designed to work with external “black box” models. Process-based models cannot simply be “plugged-into” these tools and are often too complicated to be re-implemented in the specific syntax of these software. In addition, **[R15]** these tools need to support re-reading their own outputs (posteriors) as new inputs (priors), which is critical for iterative updating of the analyses. Due to lack of available tools, models are frequently used uncalibrated (or hand-tuned; Seidel et al., 2018). Assessment of uncalibrated (or naively calibrated) models can cause poor calibration to be mistaken for inadequate model structure or mask real problems with the model structure, hindering overall progress in model development (van Oijen, 2017). **[R16]** Using multiple data constraints can be critical to ensuring that a model is getting the right answer for the right reason (Medlyn et al., 2015). Even when a model is calibrated for one setting (e.g., site or period), it does not guarantee reliable performance at another setting because there is variability and heterogeneity in natural systems. More flexible techniques, such as hierarchical Bayesian calibration, can formally quantify the scales of unexplained system variability and inform directions for model development (van Oijen, 2017), but there are even fewer available tools for their standard implementation with external models.

Within a community cyberinfrastructure, the challenge of developing advanced calibration tools only needs to be faced by

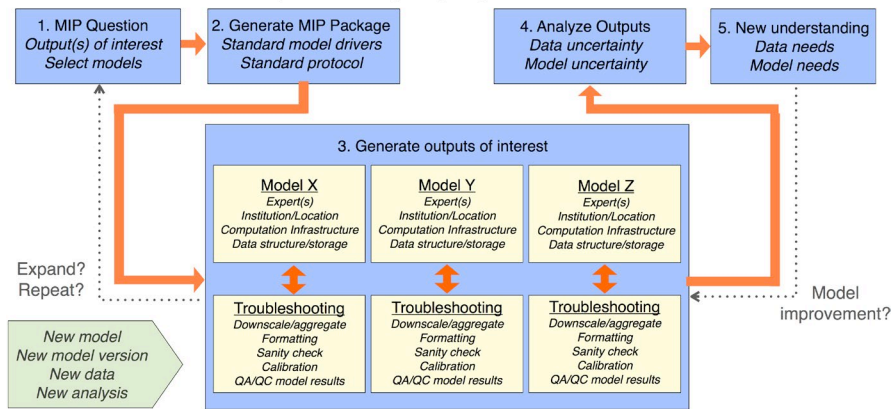
statistics experts. Software alternatives for calibrating “black-box” models are becoming increasingly available (Fer et al., 2018; Hartig et al., 2019; Huang et al., 2019). **[R17]** Community cyberinfrastructure will be most successful if hierarchical calibration tools are able to account for all kinds of ecological variability and heterogeneity (Farley et al., 2018), and if coupling to a calibration workflow is part of model development. When calibration tools are implemented in community cyberinfrastructure, they can seamlessly link multiple data constraints with multiple models. As such workflows are tracked by provenance systems, **[R18]** results from one analysis (e.g., posteriors) can readily be used by a subsequent analysis elsewhere, accelerating our ability to confront models with data. Investing in such standardization and generalization will not only allow a wider audience to adopt these methods as common practices but also foster progress on **[R19]** developing novel, more advanced calibration techniques (e.g., with emulators, Fer et al., 2018; deep learning, Tao et al., 2020).

## 5 | MODEL INTERCOMPARISON AND BENCHMARKING

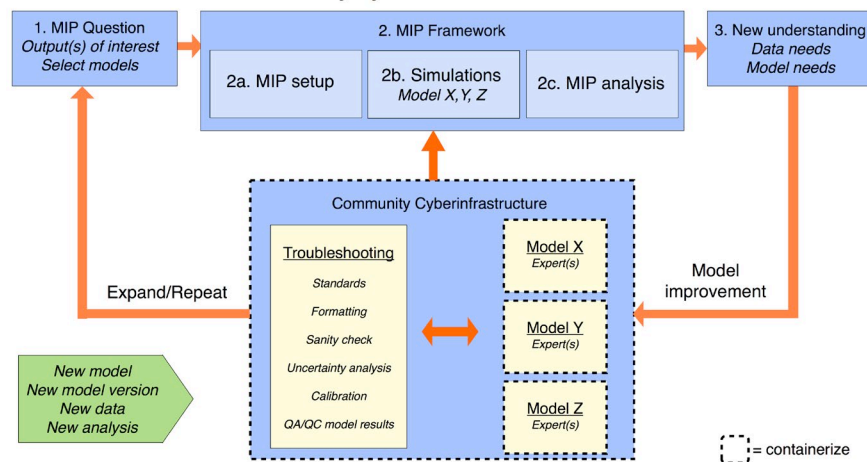
Comparing models to data is at the heart of hypothesis testing and model evaluation (Best et al., 2015; Fisher et al., 2014). While process-based models are frequently compared to multiple datasets across their lifespan, it is remarkably rare to put an ecosystem model through all its past assessment exercises every time it is updated unless a workflow has been automated (Best et al., 2015; Collier et al., 2018). **[R20]** To verify progress, and assess the tradeoffs between model parsimony and complexity, key datasets need to be set as “benchmarks” to track and compare performance through time (Best et al., 2015; Luo et al., 2012). Benchmark data can also be used to compare across models as part of model intercomparison projects (MIPs). However, the lack of automated and shared workflows also makes traditional MIPs logistically challenging to coordinate and repeat (Figure 3, top panel). Modeling groups could face incompatibilities in their results due to differences in their model configurations (e.g., calibrated vs. uncalibrated). Furthermore, due to the cost of performing a MIP, model output requests and experimental designs are typically kept simple. For example, MIPs largely focus on single model realizations which can lead to biased or overprecise decisions about model performances.

Many of the utilities that are particularly valuable for MIPs and benchmarking are already included in embedding each individual model in the community cyberinfrastructure (Figure 3, bottom panel). The use of a cyberinfrastructure also opens up the possibility of more advanced MIP benchmark activities, such as running ensembles to propagate input uncertainty to model output uncertainty. Generating multi-model ensembles with uncertainties is also practical for studying model structural errors (Bonan & Doney, 2018) and for model averaging which could potentially reduce prediction errors (Dormann et al., 2018). **[R21]** We recommend the community move toward benchmarks that account for model and data uncertainty, and leverage

### Traditional Model Intercomparison Project (MIP) Framework



### MIP Framework with a Community Cyberinfrastructure



**FIGURE 3** Traditional multi-model intercomparison project (MIP) workflow versus community cyberinfrastructure. Historically, each model and associated experts/infrastructure individually engage with MIPs (top). While stimulating model improvement is intended, it is not inherently nor readily available in traditional MIPs. In a community cyberinfrastructure, by contrast, both standardization of inputs and outputs and troubleshooting are included in embedding each individual model in the system (bottom) where MIP analyses are a use case. MIP conclusions relevant for model or cyberinfrastructure development can be fed directly back into this framework

this information when computing model performance scores (e.g., benchmarking that takes into account the uncertainty bounds in models and observations to calculate a score based on overlap probability).

Once a model is integrated into community cyberinfrastructure, it becomes trivial to add its alternative versions, benchmark against existing MIPs and seamlessly feed back to future model developments (Collier et al., 2018; Kelley et al., 2013; Wieder et al., 2019). For example, advancing model versions would benefit from being continually tested against the Free-Air CO<sub>2</sub> Experiments (FACE-MIP; De Kauwe et al., 2014; Hoffman et al., 2017) and the Arctic-Boreal Vulnerability Experiment (Fisher et al., 2018). Within or in addition to existing frameworks, interactive environments (e.g., Rstudio/Jupyter) would allow users to perform more extensive analyses with pre-loaded and aligned models and data. However, a number of challenges remain, including how to deal with datasets and metrics that are incomplete or inconsistent with each other (Collier et al., 2018; Hoffman et al., 2017). [R22] Thus, we further recommend model developers enable direct comparison to observations when possible. For example, instead of relying on modeled data products (e.g., leaf area index) whose uncertainties are harder to determine, models can be augmented to predict observations (e.g., reflected spectral radiance) as measured by the instruments. In other words, bringing models to data, rather than the other way around, may eventually reduce artificial inconsistencies between datasets that stem from additional

manipulations for making data and models match. Concomitantly, community cyberinfrastructure would facilitate [R23] interaction with a compilation of standard datasets that models need to be able to reproduce repeatedly (Anderson-Teixeira et al., 2018; Kraemer et al., 2020; Reyer et al., 2020).

## 5.1 | Who sets up benchmarks?

To address the bottleneck that only a small fraction of the data collected by ecologists (often with the aim of improving projections) ever makes its way into ecosystem models and scale up, data generators and disciplinary experts need also be equipped with tools for data-model comparison, not only the “modeler” minority (Seidl, 2017). Through community cyberinfrastructure, [R24] domain experts will more easily be able to compare multiple models to their data and set up persistent benchmarks. For example, with input/output standardization and data harmonization, the person leading the MIP no longer needs to be concerned with multiple file formats and model-specific terminology while assessing the underlying processes and mechanisms represented in the models. As cyberinfrastructure automates tedious activities associated with a MIP, experts can focus on their analysis rather than the logistics, making modeling activities more relevant for their science.



Yet, even before the challenges of running a model or a MIP, it is nearly impossible for non-modelers to keep abreast of which models exist, their most updated version, and their respective strengths and weaknesses (Jeltsch et al., 2013; Schwalm et al., 2019). [R25] Therefore, we further recommend developers encode model structural characteristics as traceable metadata. Although there are preliminary examples of this (e.g., MSTMIP encoding presence and absence of process representations; Huntzinger et al., 2016), standards need to be developed by the community to provide information about key structural characteristics of models. As a result, process representations that repeatedly perform below average across multiple MIPs can be considered rejected hypotheses (Schwalm et al., 2019), which community cyberinfrastructure could track and in return inform the development of the next generation of models as advancing new hypotheses can regain focus. In time, by centralizing these comparisons into databases, community cyberinfrastructure allows new users to discover new models and to evaluate their updated process representations with minimal technical barriers while allowing the modeling minority to focus on learning from their colleagues and improving models, rather than the status quo where the majority of their time is spent on mundane informatics issues.

## 6 | DATA ASSIMILATION AND ECOLOGICAL FORECASTING

For ecology to respond to the pace of global change, and better inform environmental decisions, the nature of the relationship between ecological models and data must be reconsidered. While most ecological analyses tend to be non-specific and a posteriori (e.g., ANOVA models), and most ecological forecasts are long term (e.g., 2100 projections), there is much to be learned from [R26] making near-term ecological forecasts that can be tested and updated as new observations become available (Dietze et al., 2018; Fox et al., 2009). Adopting an iterative forecasting approach will not only make ecology more relevant to the society, by providing information on fast, decision-relevant timescales, but will also transform basic ecological science and theory (Dietze et al., 2018), by accelerating the pace at which specific, quantitative, and falsifiable predictions are confronted with data.

Like calibration, the data assimilation methods that drive forecasting, through a formal fusion of data and modeled states (or both states and parameters), also require advanced statistical and computational expertise. Ecological models and data frequently violate the statistical assumptions embedded in assimilation algorithms developed in other disciplines (e.g., normality, homoscedasticity, independence); hence, [R27] many existing tools need to be reassessed and generalized by experts within community tools to appropriately meet the ecological model-data characteristics (Raiho et al., 2020). Making a forecast operational also requires [R28] a higher level of repeatability and efficient scheduling of cyclic workflows, where a large number of jobs are executed at regular intervals and each

forecast cycle depends on previous ones (Oliver et al., 2019). Overall, the breadth of expertise and investment of resources needed to set up a forecasting pipeline using state-of-the-art data assimilation methods often exceeds the limits of individualistic efforts (White et al., 2019).

Community-level development of automated pipelines provides a key economy of scale in data assimilation and forecasting and builds upon many of the features already discussed (Dietze et al., 2018): informatics tasks of gathering, processing, and standardizing new data will maximize data use and diversity of contributions. Managing the execution of analytical workflows will refine analyses and make them applicable to new problems. [R29] By publicly archiving and reporting results community cyberinfrastructure enables comparisons of different forecasting approaches, future syntheses, and assessment of improvement over time. These features are integral to the vision for such an infrastructure and could then be coupled to, and build upon, existing community tools for workflow scheduling (Oliver et al., 2019) and data assimilation (Fox et al., 2018; Raiho et al., 2020; Pinnington et al., 2020).

## 7 | CONCLUSIONS

Scientists, managers, and policymakers increasingly rely on models to understand the impact of decisions on ecological processes (Arneth et al., 2014; Bonan & Doney, 2018; Smith et al., 2019). As the barriers to entry for using the latest models and data are lowered, decisions will be made with better information, and scientific problems will be solved more quickly. Community cyberinfrastructure is the engine to bring time frames associated with model-data integration in line with the pressing needs of managers, policymakers, and society more broadly. We summarize our major recommendations for promptly meeting the dispersed and variable model-data synthesis needs of the ecological community as follows.

### 7.1 | Integrated community principles and practices

Modeling needs to be open, verifiable, and credible. Three key concepts in modeling cyberinfrastructure—abstraction, automation, and provenance—open up the possibility for realistic replication, community-wide transparency, and model-based ecological analysis. Adopting common cyberinfrastructure tools that are accessible, reproducible, interoperable, scalable, and community driven will play a critical role in reshaping how ecologists interact with models.

### 7.2 | Reusable data and software

Data processing remains a bottleneck to model improvement. To foster effective discovery and reuse of both data and software, we

recommend human- and machine-friendly community-scale approaches. Developing reusable tools based on community standards and involving the measurement community more deeply in data-model integration are both essential for scaling up modeling efforts.

### 7.3 | More advanced calibration techniques

Testing hypotheses should be done with properly calibrated models. Inconsistencies in model comparison due to different calibration procedures will be reduced by employing shared Bayesian calibration tools that are set up to work with process-based models. Hierarchical Bayesian calibration solutions and novel algorithms, developed and generalized under community cyberinfrastructure, will help us better capture the inherent variability and heterogeneity in ecological systems.

### 7.4 | Persistent benchmarks

Model benchmarking and intercomparison are dynamic activities that need to continually inform model improvement. We recommend a more streamlined, easily repeated, and modified process for benchmarking a suite of models with varying levels of process complexity and scale. Community cyberinfrastructure will allow domain experts to determine and more directly influence the most salient datasets that models need to replicate to demonstrate that they are capturing processes correctly, and then take the lead in setting up and performing these benchmarks.

### 7.5 | Near-term ecological forecasts

Automated data assimilation and forecasting pipelines are a necessity for ecology to support decision-making in an increasingly non-equilibrium world that has moved outside of historical norms. Building these forecasting systems requires complex automated systems, and community cyberinfrastructure is well-positioned for putting the parts of operational forecasts together.

Process-based models, though imperfect, are our window into the future functioning of ecosystems under global change. The next generation of ecological models will need to ingest increasingly diverse and expansive data to inform and test new process representations and scaling approaches, allow rapid detection and explanation of global change patterns, and even possibly allow them to be prevented. This need is now more pressing than ever. To achieve ecological model-data integration in a way that is transparent, easily communicable, and scales up to the size and diversity of the ecological community, we must invest in community cyberinfrastructure.

### ACKNOWLEDGEMENTS

We are grateful to Kristina Anderson-Teixeira and Mingkai Jiang for their reviews which improved this manuscript to a great extent. We

further thank Gab Abramowitz, Veronika Eyring, Michael Fienen, Andy Fox, Yuan Gao, Birgit Hassler, Xin Huang, Randall Hunt, Lifan Jiang, and Jeremy White for their helpful overview on example cyberinfrastructure tools. The PEcAn project which organized the workshop where the authors of this paper came together is supported by the NSF (ABI no. 1062547, ABI no. 1458021, ABI no. 1457897, ABI no. 1062204, DIBBS no. 1261582), NASA Terrestrial Ecosystems, the Energy Biosciences Institute, and an Amazon AWS education grant. We would also like to thank Boston University for providing the venue for the workshop that inspired this article. I.F. and T.V. acknowledge funding from the Strategic Research Council at the Academy of Finland (decision 327214), the Academy of Finland (decision 297350), and Business Finland (decision 6905/31/2018) to the Finnish Meteorological Institute. T.Q. is funded by the UK NERC National Centre for Earth Observation. J.B.F. contributed to this work from the Jet Propulsion Laboratory, California Institute of Technology, under a contract with the National Aeronautics and Space Administration. California Institute of Technology. J.B.F. was supported in part by NASA programs: CARBON and CMS. S.P.S. was partially supported by NASA CMS (grant #80NSSC17K0711), and through the DOE Reducing Uncertainties in Biogeochemical Interactions through Synthesis and Computation Science Focus Area (RUBISCO SFA), which is sponsored by the Earth & Environmental Systems Modeling (EESM) Program in the Climate and Environmental Sciences Division (CESD), and the Next-Generation Ecosystem Experiments (NGEE-Arctic and NGEE-Tropics) supported by the Office of Biological and Environmental Research in the Department of Energy, Office of Science, as well as through the United States Department of Energy contract no. DE-SC0012704 to Brookhaven National Laboratory. M.D.K. acknowledges funding from the Australian Research Council (ARC) Centre of Excellence for Climate Extremes (CE170100023), the ARC Discovery Grant (DP190101823) and support from the NSW Research Attraction and Acceleration Program. F.M.H. was partially supported by the Laboratory Directed Research and Development Program of Oak Ridge National Laboratory, which is managed by UT-Battelle, LLC, for the U.S. Department of Energy under contract DE-AC05-00OR22725. Additional support was provided by the Data Program, by the Reducing Uncertainties in Biogeochemical Interactions through Synthesis and Computation Science Focus Area (RUBISCO SFA) in the Earth & Environmental Systems Modeling (EESM) Program, and by the Next-Generation Ecosystem Experiments (NGEE-Arctic and NGEE-Tropics) Projects in the Terrestrial Ecosystem Science (TES) Program. The Data, EESM, and TES Programs are part of the Climate and Environmental Sciences Division (CESD) of the Office of Biological and Environmental Research (BER) in the U.S. Department of Energy Office of Science.

### CONFLICT OF INTEREST

The authors declare no competing interests.

### AUTHOR CONTRIBUTION

All authors were present in the workshop where these ideas were discussed. I.F. and A.K.G. lead the writing with extensive feedback

from M.C.D. and with contributions from all authors. All authors have read and approved the manuscript.

## CODE AVAILABILITY

Code availability not applicable to this article. However, we note for the interested reader that all example community tools mentioned in Appendix C are open source and available on online code repositories.

## DATA AVAILABILITY STATEMENT

Data sharing not applicable to this article as no datasets were generated or analyzed during the particular study.

## ORCID

Istem Fer  <https://orcid.org/0000-0001-8236-303X>  
 Alexey N. Shiklomanov  <https://orcid.org/0000-0003-4022-5979>  
 Eleanor E. Campbell  <https://orcid.org/0000-0002-9272-6276>  
 Elizabeth M. Cowdery  <https://orcid.org/0000-0002-6538-6296>  
 Martin G. De Kauwe  <https://orcid.org/0000-0002-3399-9098>  
 Ankur Desai  <https://orcid.org/0000-0002-5226-6041>  
 Matthew J. Duveneck  <https://orcid.org/0000-0003-1264-3081>  
 Joshua B. Fisher  <https://orcid.org/0000-0003-4734-9085>  
 Katherine D. Haynes  <https://orcid.org/0000-0001-9637-7766>  
 Forrest M. Hoffman  <https://orcid.org/0000-0001-5802-4134>  
 Miriam R. Johnston  <https://orcid.org/0000-0001-7481-8794>  
 Rob Kooper  <https://orcid.org/0000-0002-5781-7287>  
 David S. LeBauer  <https://orcid.org/0000-0001-7228-053X>  
 Benjamin Poulter  <https://orcid.org/0000-0002-9493-8600>  
 Tristan Quaife  <https://orcid.org/0000-0001-6896-4613>  
 Ann Raiho  <https://orcid.org/0000-0002-2552-3399>  
 Kevin Schaefer  <https://orcid.org/0000-0002-5444-9917>  
 Shawn P. Serbin  <https://orcid.org/0000-0003-4136-8971>  
 Kevin R. Wilcox  <https://orcid.org/0000-0001-6829-1148>  
 Toni Viskari  <https://orcid.org/0000-0002-3357-1374>  
 Michael C. Dietze  <https://orcid.org/0000-0002-2324-2518>

## REFERENCES

- Anderson-Teixeira, K. J., Wang, M. M. H., McGarvey, J. C., Herrmann, V., Tepley, A. J., Bond-Lamberty, B., & LeBauer, D. S. (2018). ForC: A global database of forest carbon stocks and fluxes. *Ecology*, 99, 1507. <https://doi.org/10.1002/ecy.2229>
- Arneth, A., Brown, C., & Rounsevell, M. D. A. (2014). Global models of human decision-making for land-based mitigation and adaptation assessment. *Nature Climate Change*, 4, 550–557. <https://doi.org/10.1038/nclimate2250>
- Arora, V. K., Katavouta, A., Williams, R. G., Jones, C. D., Brovkin, V., Friedlingstein, P., Schwinger, J., Bopp, L., Boucher, O., Cadule, P., Chamberlain, M. A., Christian, J. R., Delire, C., Fisher, R. A., Hajima, T., Ilyina, T., Joetzer, E., Kawamiya, M., Koven, C. D., ... Ziehn, T. (2020). Carbon-concentration and carbon-climate feedbacks in CMIP6 models and their comparison to CMIP5 models. *Biogeosciences*, 17, 4173–4222. <https://doi.org/10.5194/bg-17-4173-2020>
- Best, M. J., Abramowitz, G., Johnson, H. R., Pitman, A. J., Balsamo, G., Boone, A., Cuntz, M., Decharme, B., Dirmeyer, P. A., Dong, J., Ek, M., Guo, Z., Haverd, V., van den Hurk, B. J. J., Nearing, G. S., Pak, B., Peters-Lidard, C., Santanello, J. A., Stevens, L., & Vuichard, N. (2015). The plumbing of land surface models: Benchmarking model performance. *Journal of Hydrometeorology*, 16, 1425–1442. <https://doi.org/10.1175/JHM-D-14-0158.1>
- Boettiger, C., Chamberlain, S., Hart, E., & Ram, K. (2015). Building software, building community: Lessons from the rOpenSci project. *Journal of Open Research Software*, 3(1), e8. <https://doi.org/10.5334/jors.bu>
- Bonan, G. B., & Doney, S. C. (2018). Climate, ecosystems, and planetary futures: The challenge to predict life in Earth system models. *Science*, 359, eaam8328. <https://doi.org/10.1126/science.aam8328>
- Bond-Lamberty, B., Smith, A. P., & Bailey, V. (2016). Running an open experiment: Transparency and reproducibility in soil and ecosystem science. *Environmental Research Letters*, 11, 084004. <https://doi.org/10.1088/1748-9326/11/8/084004>
- Collier, N., Hoffman, F. M., Lawrence, D. M., Keppel-Aleks, G., Koven, C. D., Riley, W. J., Mu, M., & Randerson, J. T. (2018). The International Land Model Benchmarking (ILAMB) system: Design, theory, and implementation. *Journal of Advances in Modeling Earth Systems*, 10, 2731–2754. <https://doi.org/10.1029/2018MS001354>
- Cook, R. B., Vannan, S. K. S., McMurry, B. F., Wright, D. M., Wei, Y., Boyer, A. G., & Kidder, J. H. (2016). Implementation of data citations and persistent identifiers at the ORNL DAAC. *Ecological Informatics*, 33, 10–16. <https://doi.org/10.1016/j.ecoinf.2016.03.003>
- Culina, A., Baglioni, M., Crowther, T. W., Visser, M. E., Woutersen-Windhauer, S., & Manghi, P. (2018). Navigating the unfolding open data landscape in ecology and evolution. *Nature Ecology & Evolution*, 2, 420–426. <https://doi.org/10.1038/s41559-017-0458-2>
- Dai, S.-Q., Li, H., Xiong, J., Ma, J., Guo, H.-Q., Xiao, X., & Zhao, B. (2018). Assessing the extent and impact of online data sharing in eddy covariance flux research. *Journal of Geophysical Research: Biogeosciences*, 123, 129–137. <https://doi.org/10.1002/2017JG004277>
- De Kauwe, M. G., Medlyn, B. E., Zaehle, S., Walker, A. P., Dietze, M. C., Wang, Y. P., Luo, Y., Jain, A. K., El-Masri, B., Hickler, T., & Wärlind, D. (2014). Where does the carbon go? A model-data intercomparison of vegetation carbon allocation and turnover processes at two temperate forest free-air CO<sub>2</sub> enrichment sites. *New Phytologist*, 203, 883–899. <https://doi.org/10.1111/nph.12847>
- Dietze, M. C., Fox, A., Beck-Johnson, L. M., Betancourt, J. L., Hooten, M. B., Jarnevich, C. S., Keitt, T. H., Kenney, M. A., Laney, C. M., Larsen, L. G., Loescher, H. W., Lunch, C. K., Pijanowski, B. C., Randerson, J. T., Read, E. K., Tredennick, A. T., Vargas, R., Weathers, K. C., & White, E. P. (2018). Iterative near-term ecological forecasting: Needs, opportunities, and challenges. *Proceedings of the National Academy of Sciences of the United States of America*, 115, 1424–1432. <https://doi.org/10.1073/pnas.1710231115>
- Dietze, M. C., LeBauer, D., & Kooper, R. (2013). On improving the communication between models and data. *Plant, Cell & Environment*, 36, 1575–1585. <https://doi.org/10.1111/pce.12043>
- Dormann, C. F., Calabrese, J. M., Guillera-Aroita, G., Matechou, E., Bahn, V., Bartoň, K., Beale, C. M., Ciuti, S., Elith, J., Gerstner, K., Guelat, J., Keil, P., Lahoz-Monfort, J. J., Pollock, L. J., Reineking, B., Roberts, D. R., Schröder, B., Thuiller, W., Warton, D. I., ... Hartig, F. (2018). Model averaging in ecology: A review of Bayesian, information-theoretic, and tactical approaches for predictive inference. *Ecological Monographs*, 88, 485–504. <https://doi.org/10.1002/ecm.1309>
- Eaton, B., Gregory, J., Drach, B., Taylor, K., Hankin, S., Caron, J., Signell, R., Bentley, P., Rappa, G., Höck, H., & Pamment, A. (2017). *Netcdf Climate and Forecast (CF) metadata conventions*. <http://cfconventions.org/>
- Eyring, V., Cox, P. M., Flato, G. M., Gleckler, P. J., Abramowitz, G., Caldwell, P., Collins, W. D., Gier, B. K., Hall, A. D., Hoffman, F. M., Hurtt, G. C., Jahn, A., Jones, C. D., Klein, S. A., Krasting, J. P., Kwiatkowski, L., Lorenz, R., Maloney, E., Meehl, G. A., ... Williamson, M. S. (2019). Taking climate model evaluation to the next level.

- Nature Climate Change*, 9, 102–110. <https://doi.org/10.1038/s41558-018-0355-y>
- Fang, S., Da Li, X., Zhu, Y., Ahati, J., Pei, H., Yan, J., & Liu, Z. (2014). An integrated system for regional environmental monitoring and management based on internet of things. *IEEE Transactions on Industrial Informatics*, 10(2), 1596–1605. <https://doi.org/10.1109/TII.2014.2302638>
- Farley, S. S., Dawson, A., Goring, S. J., & Williams, J. W. (2018). Situating ecology as a big-data science: Current advances, challenges, and solutions. *BioScience*, 68, 563–576. <https://doi.org/10.1093/biosci/biy068>
- Fer, I., Kelly, R., Moorcroft, P. R., Richardson, A. D., Cowdery, E. M., & Dietze, M. C. (2018). Linking big models to big data: Efficient ecosystem model calibration through Bayesian model emulation. *Biogeosciences*, 15, 5801–5830. <https://doi.org/10.5194/bg-15-5801-2018>
- Fisher, J. B., Hayes, D. J., Schwalm, C. R., Huntzinger, D. N., Stofferahn, E., Schaefer, K., Luo, Y., Wulschleger, S. D., Goetz, S., Miller, C. E., Griffith, P., Chadburn, S., Chatterjee, A., Ciais, P., Douglas, T. A., Genet, H., Ito, A., Neigh, C. S. R., Poulter, B., ... Zhang, Z. (2018). Missing pieces to modeling the Arctic Boreal puzzle. *Environmental Research Letters*, 13, 020202. <https://doi.org/10.1088/1748-9326/aa9d9a>
- Fisher, J. B., Huntzinger, D. N., Schwalm, C. R., & Stith, S. (2014). Modeling the terrestrial biosphere. *Annual Review of Environment and Resources*, 39, 91–123. <https://doi.org/10.1146/annurev-environ-012913-093456>
- Fox, A., Hoar, T. J., Anderson, J. L., Arellano, A. F., Smith, W. K., Litvak, M. E., MacBean, N., Schimel, D. S., & Moore, D. J. (2018). Evaluation of a data assimilation system for land surface models using CLM4.5. *Journal of Advances in Modeling Earth Systems*, 10, 2471–2494. <https://doi.org/10.1002/2018MS001362>
- Fox, A., Williams, M., Richardson, A. D., Cameron, D., Gove, J. H., Quaife, T., Ricciuto, D., Reichstein, M., Tomelleri, E., Trudinger, C. M., & Van Wijk, M. T. (2009). The REFLEX project: Comparing different algorithms and implementations for the inversion of a terrestrial ecosystem model against eddy covariance data. *Agricultural and Forest Meteorology*, 149, 1597–1615. <https://doi.org/10.1016/j.agrfo.2009.05.002>
- Friedlingstein, P., Cox, P., Betts, R., Bopp, L., von Bloh, W., Brovkin, V., Cadule, P., Doney, S., Eby, M., Fung, I., Bala, G., John, J., Jones, C., Joos, F., Kato, T., Kawamiya, M., Knorr, W., Lindsay, K., Matthews, H. D., ... Zeng, N. (2006). Climate-carbon cycle feedback analysis: Results from the C4MIP model intercomparison. *Journal of Climate*, 19, 3337–3353. <https://doi.org/10.1175/JCLI3800.1>
- Friedlingstein, P., Meinshausen, M., Arora, V. K., Jones, C. D., Anav, A., Liddicoat, S. K., & Knutti, R. (2014). Uncertainties in CMIP5 climate projections due to carbon cycle feedbacks. *Journal of Climate*, 27, 511–526. <https://doi.org/10.1175/JCLI-D-12-00579.1>
- Gil, Y., David, C. H., Demir, I., Essawy, B. T., Fulweiler, R. W., Goodall, J. L., Karlstrom, L., Lee, H., Mills, H. J., Oh, J.-H., Pierce, S. A., Pope, A., Tzeng, M. W., Villamizar, S. R., & Yu, X. (2016). Toward the geoscience paper of the future: Best practices for documenting and sharing research from data to software to provenance. *Earth and Space Science*, 3, 388–415. <https://doi.org/10.1002/2015EA000136>
- Gomes, V. C., Queiroz, G. R., & Ferreira, K. R. (2020). An overview of platforms for big earth observation data management and analysis. *Remote Sensing*, 12(8), 1253. <https://doi.org/10.3390/rs12081253>
- Gries, C., Servilla, M., O'Brien, M., Vanderbilt, K., Smith, C., Costa, D., & Grossman-Clarke, S. (2019). Achieving FAIRData principles at the environmental data initiative, the US-LTER data repository. *Biodiversity Information Science and Standards*, 3, e37047. <https://doi.org/10.3897/biss.3.37047>
- Hanson, P. J., & Walker, A. P. (2020). Advancing global change biology through experimental manipulations: Where have we been and where might we go? *Global Change Biology*, 26, 287–299. <https://doi.org/10.1111/gcb.14894>
- Hart, E. M., Barmby, P., LeBauer, D., Michonneau, F., Mount, S., Mulrooney, P., Poisot, T., Woo, K. H., Zimmerman, N. B., & Hollister, J. W. (2016). Ten simple rules for digital data storage. *PLoS Computational Biology*, 12, e1005097. <https://doi.org/10.1371/journal.pcbi.1005097>
- Hartig, F., Dyke, J., Hickler, T., Higgins, S., O'Hara, R., Scheiter, S., & Huth, A. (2012). Connecting dynamic vegetation models to data – An inverse perspective. *Journal of Biogeography*, 39, 2240–2252. <https://doi.org/10.1111/j.1365-2699.2012.02745.x>
- Hartig, F., Minunno, F., & Paul, S. (2019). *BayesianTools: General-purpose MCMC and SMC samplers and tools for Bayesian statistics*. R package version 0.1.7. <https://cran.r-project.org/web/packages/BayesianTools/>
- Hasselbring, W., Carr, L., Hettrick, S., Packer, H., & Tiropanis, T. (2020). From FAIR research data toward FAIR and open research software. *it – Information Technology*, 62(1), 39–47. <https://doi.org/10.1515/itit-2019-0040>
- Herger, N., Abramowitz, G., Sherwood, S., Knutti, R., Angéil, O., & Sisson, S. A. (2019). Ensemble optimisation, multiple constraints and overconfidence: A case study with future Australian precipitation change. *Climate Dynamics*, 53, 1581–1596. <https://doi.org/10.1007/s00382-019-04690-8>
- Hoffman, F. M., Koven, C. D., Keppel-Aleks, G., Lawrence, D. M., Riley, W. J., Randerson, J. T., Ahlström, A., Abramowitz, G., Baldocchi, D. D., Best, M. J., & Bond-Lamberty, B. (2017). *International Land Model Benchmarking (ILAMB) 2016 workshop report*, DOE/SC-0186. U.S. Department of Energy, Office of Science. <https://doi.org/10.2172/1330803>
- Huang, Y., Stacy, M., Jiang, J., Sundi, N., Ma, S., Saruta, V., Jung, C. G., Shi, Z., Xia, J., Hanson, P. J., Ricciuto, D., & Luo, Y. (2019). Realized ecological forecast through an interactive Ecological Platform for Assimilating Data (EcoPAD, v1.0) into models. *Geoscientific Model Development*, 12, 1119–1137. <https://doi.org/10.5194/gmd-12-1119-2019>
- Huntzinger, D. N., Schwalm, C. R., Wei, Y., Cook, R. B., Michalak, A. M., Schaefer, K., Jacobson, A. R., Arain, M. A., Ciais, P., Fisher, J. B., & Hayes, D. J. (2016). *NACP MstMIP: Global 0.5-deg terrestrial biosphere model outputs (version 1) in standard format*. ORNL DAAC. <https://doi.org/10.3334/ORNLDAAAC/1225>
- Jeltsch, F., Blaum, N., Brose, U., Chipperfield, J. D., Clough, Y., Farwig, N., Geissler, K., Graham, C. H., Grimm, V., Hickler, T., Huth, A., May, F., Meyer, K. M., Pagel, J., Reineking, B., Rillig, M. C., Shea, K., Schurr, F. M., Schröder, B., ... Zurell, D. (2013). How can we bring together empiricists and modellers in functional biodiversity research? *Basic and Applied Ecology*, 14(2), 93–101. <https://doi.org/10.1016/j.baae.2013.01.001>
- Keenan, T. F., Davidson, E. A., Munger, J. W., & Richardson, A. D. (2013). Rate my data: Quantifying the value of ecological data for the development of models of the terrestrial carbon cycle. *Ecological Applications*, 23(1), 273–286. <https://doi.org/10.1890/12-0747.1>
- Kelley, D. I., Prentice, I. C., Harrison, S. P., Wang, H., Simard, M., Fisher, J. B., & Willis, K. O. (2013). A comprehensive benchmarking system for evaluating global vegetation models. *Biogeosciences*, 10, 3313–3340. <https://doi.org/10.5194/bg-10-3313-2013>
- Kraemer, G., Camps-Valls, G., Reichstein, M., & Mahecha, M. D. (2020). Summarizing the state of the terrestrial biosphere in few dimensions. *Biogeosciences*, 17, 2397–2424. <https://doi.org/10.5194/bg-17-2397-2020>
- LaDeau, S. L., Han, B. A., Rosi-Marshall, E. J., & Weathers, K. C. (2017). The next decade of big data in ecosystem science. *Ecosystems*, 20, 274–283. <https://doi.org/10.1007/s10021-016-0075-y>
- Lai, J., Lortie, C. J., Muenchen, R. A., Yang, J., & Ma, K. (2019). Evaluating the popularity of R in ecology. *Ecosphere*, 10(1), e02567. <https://doi.org/10.1002/ecs2.2567>



- LeBauer, D. S., Kooper, R., Mulrooney, P., Rohde, S., Wang, D., Long, S. P., & Dietze, M. C. (2018). BETYdb: A yield, trait, and ecosystem service database applied to second-generation bioenergy feedstock production. *GCB Bioenergy*, 10, 61–71. <https://doi.org/10.1111/gcbb.12420>
- LeBauer, D. S., Wang, D., Richter, K. T., Davidson, C. C., & Dietze, M. C. (2013). Facilitating feedbacks between field measurements and ecosystem models. *Ecological Monographs*, 83, 133–154. <https://doi.org/10.1890/12-0137.1>
- Lovenduski, N. S., & Bonan, G. B. (2017). Reducing uncertainty in projections of terrestrial carbon uptake. *Environmental Research Letters*, 12, 044020. <https://doi.org/10.1088/1748-9326/aa66b8>
- Lowndes, J. S. S., Best, B. D., Scarborough, C., Afflerbach, J. C., Frazier, M. R., O'Hara, C. C., Jiang, N., & Halpern, B. S. (2017). Our path to better science in less time using open data science tools. *Nature Ecology & Evolution*, 1, 0160. <https://doi.org/10.1038/s41559-017-0160>
- Luo, Y. Q., Randerson, J. T., Abramowitz, G., Bacour, C., Blyth, E., Carvalhais, N., Ciais, P., Dalmonech, D., Fisher, J. B., Fisher, R., Friedlingstein, P., Hibbard, K., Hoffman, F., Huntzinger, D., Jones, C. D., Koven, C., Lawrence, D., Li, D. J., Mahecha, M., ... Zhou, X. H. (2012). A framework for benchmarking land models. *Biogeosciences*, 9, 3857–3874. <https://doi.org/10.5194/bg-9-3857-2012>
- McKiernan, E. C., Bourne, P. E., Brown, C. T., Buck, S., Kenall, A., Lin, J., McDougall, D., Nosek, B. A., Ram, K., Soderberg, C. K., Spies, J. R., Thaney, K., Updegrove, A., Woo, K. H., & Yarkoni, T. (2016). How open science helps researchers succeed. *Elife*, 5, e16800. <https://doi.org/10.7554/eLife.16800>
- Medlyn, B. E., Zaehle, S., De Kauwe, M. G., Walker, A. P., Dietze, M. C., Hanson, P. J., Hickler, T., Jain, A. K., Luo, Y., Parton, W., Prentice, I. C., Thornton, P. E., Wang, S., Wang, Y.-P., Weng, E., Iversen, C. M., McCarthy, H. R., Warren, J. M., Oren, R., & Norby, R. J. (2015). Using ecosystem experiments to improve vegetation models. *Nature Climate Change*, 5, 528–534. <https://doi.org/10.1038/nclimate2621>
- Nagaraj, A., Shears, E., & de Vaan, M. (2020). Improving data access democratizes and diversifies science. *Proceedings of the National Academy of Sciences of the United States of America*, 117(38), 23490–23498. <https://doi.org/10.1073/pnas.2001682117>
- Oliver, H., Shin, M., Matthews, D., Sanders, O., Bartholomew, S., Clark, A., Fitzpatrick, B., van Haren, R., Hut, R., & Drost, N. (2019). Workflow automation for cycling systems. *Computing in Science & Engineering*, 21, 7–21. <https://doi.org/10.1109/MCSE.2019.2906593>
- Piccolo, S. R., & Frampton, M. B. (2016). Tools and techniques for computational reproducibility. *GigaScience*, 5. <https://doi.org/10.1186/s13742-016-0135-4>
- Pinnington, E., Quaife, T., Lawless, A., Williams, K., Arkebauer, T., & Scoby, D. (2020). The Land Variational Ensemble Data Assimilation Framework: LAVENDAR v1.0.0. *Geoscientific Model Development*, 13, 55–69. <https://doi.org/10.5194/gmd-13-55-2020>
- Powers, J. G., Klemp, J. B., Skamarock, W. C., Davis, C. A., Dudhia, J., Gill, D. O., Coen, J. L., Gochis, D. J., Ahmadov, R., Peckham, S. E., Grell, G. A., Michalakes, J., Trahan, S., Benjamin, S. G., Alexander, C. R., Dimego, G. J., Wang, W., Schwartz, C. S., Romine, G. S., ... Duda, M. G. (2017). The weather research and forecasting model: Overview, system efforts, and future directions. *Bulletin of the American Meteorological Society*, 98, 1717–1737. <https://doi.org/10.1175/BAMS-D-15-00308.1>
- Powers, S. M., & Hampton, S. E. (2019). Open science, reproducibility, and transparency in ecology. *Ecological Applications*, 29(1), e01822. <https://doi.org/10.1002/eap.1822>
- R Core Team. (2020). *R: A language and environment for statistical computing*. R version 4.0.3. R Foundation for Statistical Computing.
- Raiho, A., Dietze, M., Dawson, A., Rollinson, C. R., Tipton, T., & McLachlan, J. (2020). Determinants of predictability in multi-decadal forest community and carbon dynamics. *bioRxiv*. <https://doi.org/10.1101/2020.05.05.079871>
- Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., & Prabhat. (2019). Deep learning and process understanding for data-driven Earth system science. *Nature*, 566, 195–204. <https://doi.org/10.1038/s41586-019-0912-1>
- Reyer, C. P. O., Silveyra Gonzalez, R., Dolos, K., Hartig, F., Hauf, Y., Noack, M., Lasch-Born, P., Rötzer, T., Pretzsch, H., Meesenburg, H., Fleck, S., Wagner, M., Bolte, A., Sanders, T. G. M., Kolari, P., Mäkelä, A., Vesala, T., Mammarella, I., Pumpanen, J., ... Frierle, K. (2020). The PROFOUND Database for evaluating vegetation models and simulating climate impacts on European forests. *Earth System Science Data*, 12, 1295–1320. <https://doi.org/10.5194/essd-12-1295-2020>
- Rineau, F., Malina, R., Beenaerts, N., Arnauts, N., Bardgett, R. D., Berg, M. P., Boerema, A., Bruckers, L., Clerinx, J., Davin, E. L., De Boeck, H. J., De Dobbelaer, T., Dondini, M., De Laender, F., Ellers, J., Franken, O., Gilbert, L., Gudmundsson, L., Janssens, I. A., ... Vangronsveld, J. (2019). Towards more predictive and interdisciplinary climate change ecosystem experiments. *Nature Climate Change*, 9, 809–816. <https://doi.org/10.1038/s41558-019-0609-3>
- Schimel, D., Schneider, F. D., Carbon, J. P. L., & Participants, E. (2019). Flux towers in the sky: Global ecology from space. *New Phytologist*, 224, 570–584. <https://doi.org/10.1111/nph.15934>
- Schwalm, C. R., Schaefer, K., Fisher, J. B., Huntzinger, D., Elshorban, Y., Fang, Y., Hayes, D., Jafarov, E., Michalak, A. M., Piper, M., Stofferahn, E., Wang, K., & Wei, Y. (2019). Divergence in land surface modeling: Linking spread to structure. *Environmental Research Communications*, 1, 111004. <https://doi.org/10.1088/2515-7620/ab4a8a>
- Seidel, S. J., Palosuo, T., Thorburn, P., & Wallach, D. (2018). Towards improved calibration of crop models – Where are we now and where should we go? *European Journal of Agronomy*, 94, 25–35. <https://doi.org/10.1016/j.eja.2018.01.006>
- Seidl, R. (2017). To model or not to model, that is no longer the question for ecologists. *Ecosystems*, 20, 222. <https://doi.org/10.1007/s10021-016-0068-x>
- Shiklomanov, A. N., Cowdery, E. M., Bahn, M., Byun, C., Jansen, S., Kramer, K., Minden, V., Niinemets, Ü., Onoda, Y., Soudzilovskaia, N. A., & Dietze, M. C. (2020). Does the leaf economic spectrum hold within plant functional types? A Bayesian multivariate trait meta-analysis. *Ecological Applications*, 30(3), 02064. <https://doi.org/10.1002/eap.2064>
- Smith, P., Soussana, J.-F., Angers, D., Schipper, L., Chenu, C., Rasse, D. P., Batjes, N. H., Egmond, F., McNeill, S., Kuhnert, M., Arias-Navarro, C., Olesen, J. E., Chirinda, N., Fornara, D., Wollenberg, E., Álvaro-Fuentes, J., Sanz-Cobena, A., & Klumpp, K. (2020). How to measure, report and verify soil carbon change to realize the potential of soil carbon sequestration for atmospheric greenhouse gas removal. *Global Change Biology*, 26(1), 219–241. <https://doi.org/10.1111/gcb.14815>
- Stucky, B. J., Guralnick, R., Deck, J., Denny, E. G., Bolmgren, K., & Walls, R. (2018). The plant phenology ontology: A new informatics resource for large-scale integration of plant phenology data. *Frontiers in Plant Science*, 9, 517. <https://doi.org/10.3389/fpls.2018.00517>
- Sullivan, I., DeHaven, A., & Mellor, D. (2019). Open and reproducible research on open science framework. *Current Protocols*, 18(1), e32. <https://doi.org/10.1002/cpet.3>
- Tao, F., Zhou, Z., Huang, Y., Li, Q., Lu, X., Ma, S., Huang, X., Liang, Y., Hugelius, G., Jiang, L., Doughty, R., Ren, Z., & Luo, Y. (2020). Deep learning optimizes data-driven representation of soil organic carbon in earth system model over the conterminous United States. *Frontiers in Big Data*, 3, 17. <https://doi.org/10.3389/fdata.2020.00017>
- van Oijen, M. (2017). Bayesian methods for quantifying and reducing uncertainty and error in forest models. *Current Forestry Reports*, 3, 269–280. <https://doi.org/10.1007/s40725-017-0069-9>
- Waide, R. B., Brunt, J. W., & Servilla, M. S. (2017). Demystifying the landscape of ecological data repositories in the United States. *BioScience*, 67(12), 1044–1051. <https://doi.org/10.1093/biosci/bix117>

- White, E. P., Yenni, G. M., Taylor, S. D., Christensen, E. M., Bledsoe, E. K., Simonis, J. L., & Ernest, S. K. M. (2019). Developing an automated iterative near-term forecasting system for an ecological study. *Methods in Ecology and Evolution*, 10, 332–344. <https://doi.org/10.1111/2041-210X.13104>
- Wieder, W. R., Lawrence, D. M., Fisher, R. A., Bonan, G. B., Cheng, S. J., Goodale, C. L., Grandy, A. S., Koven, C. D., Lombardozzi, D. L., Oleson, K. W., & Thomas, R. Q. (2019). Beyond static benchmarking: Using experimental manipulations to evaluate land model assumptions. *Global Biogeochemical Cycles*, 33, 1289–1309. <https://doi.org/10.1029/2018GB006141>

## SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section.

**How to cite this article:** Fer I, Gardella AK, Shiklomanov AN, et al. Beyond ecosystem modeling: A roadmap to community cyberinfrastructure for ecological data-model integration. *Glob Change Biol*. 2021;27:13–26. <https://doi.org/10.1111/gcb.15409>