

# UC Berkeley

## UC Berkeley Previously Published Works

### Title

A metadata reporting framework (FRAMES) for synthesis of ecohydrological observations

### Permalink

<https://escholarship.org/uc/item/7dr46409>

### Authors

Christianson, Danielle S  
Varadharajan, Charuleka  
Christoffersen, Bradley  
[et al.](#)

### Publication Date

2017-11-01

### DOI

10.1016/j.ecoinf.2017.06.002

Peer reviewed

## A metadata reporting framework (FRAMES) for synthesis of ecohydrological observations

S.Christianson<sup>ab</sup> CharulekaVaradharajan<sup>a</sup> BradleyChristoffersen<sup>c</sup> MatteoDetto<sup>de</sup>  
BorisFaybishenko<sup>a</sup> Bruno O.Gimenez<sup>f</sup> ValHendrix<sup>b</sup> Kolby J.Jardine<sup>a</sup>  
RobinsonNegrón-Juarez<sup>a</sup> Gilberto Z.Pastorello<sup>b</sup> Thomas L.Powell<sup>a</sup>  
MeghaSandesh<sup>b</sup> Jeffrey M.Warren<sup>g</sup> Brett T.Wolfe<sup>d</sup> Jeffrey Q.Chambers<sup>a</sup> Lara  
M.Kueppers<sup>ah</sup> Nathan G.McDowell<sup>cg</sup> Deborah A.Agarwal<sup>b</sup>

### Abstract

Metadata describe the ancillary information needed for data preservation and independent interpretation, comparison across heterogeneous datasets, and quality assessment and quality control (QA/QC). Environmental observations are vastly diverse in type and structure, can be taken across a wide range of spatiotemporal scales in a variety of measurement settings and approaches, and saved in multiple formats. Thus, well-organized, consistent metadata are required to produce usable data products from diverse environmental observations collected across field sites. However, existing metadata reporting protocols do not support the complex data synthesis and model-data integration needs of interdisciplinary earth system research. We developed a metadata reporting framework (FRAMES) to enable management and synthesis of observational data that are essential in advancing a predictive understanding of earth systems. FRAMES utilizes best practices for data and metadata organization enabling consistent data reporting and compatibility with a variety of standardized data protocols. We used an iterative scientist-centered design process to develop FRAMES, resulting in a data reporting format that incorporates existing field practices to maximize data-entry efficiency. Thus, FRAMES has a modular organization that streamlines metadata reporting and can be expanded to incorporate additional data types. With FRAMES's multi-scale measurement position hierarchy, data can be reported at observed spatial resolutions and then easily aggregated and linked across measurement types to support model-data integration. FRAMES is in early use by both data originators (persons generating data) and consumers (persons using data and metadata). In this paper, we describe FRAMES, identify lessons learned, and discuss areas of future development.

### Abbreviations

FRAMES, Framework for Reporting dAta and Metadata for Earth Systems

FATES, Functionally Assembled Terrestrial Ecosystem Simulator

QA/QC, quality assurance/quality control

ENSO, El Niño Southern Oscillation

STRI, Smithsonian Tropical Research Institute

CTFS, Center for Tropical Forest Study

BADM, Biological, Ancillary, Disturbance, and Metadata

ISCN, International Soil Carbon Network

Keywords: Metadata, Data management system, Model-data integration, Data synthesis, Data preservation, Informatics

## 1. Introduction

Current earth systems research challenges, like understanding and predicting carbon cycling in tropical forests under a changing climate, require synthesis of complex and diverse earth system observations. Researchers use synthesized data products to understand the controls and rates of environmental processes, as well as constrain, parameterize, and benchmark process-rich models (e.g., Medlyn et al., 2005). Data synthesis refers to the process of connecting diverse observations collected across field sites and a wide range of spatial and temporal scales to answer a science question or to generate model inputs. Prior to synthesis, each observation must be quality checked, processed (e.g., units transformed, gap-filled, erroneous data flagged or removed), and organized in standardized, comparable formats (e.g., variable names, units). An example of a synthesized data product is the FLUXNET2015 dataset, which includes data collected at sites from a network of single-locale, eddy covariance towers that monitor an ecosystem over many years (FLUXNET, 2016). In addition to ecosystem and global scale datasets, earth system science requires syntheses of individual-based measures like point observations of leaf carbohydrate content, continuous tree sap flow, and demography censuses (e.g., Walker et al., 2014). Physical measures, such as meteorological observations, measurements of soil water content, and 3D structural representations (e.g., LiDAR), are also needed (e.g., Hunter et al., 2015, Powell et al., 2013).

Metadata are essential to describe the different approaches taken to obtain, process, and report diverse ecohydrological and biogeochemical observations and the resulting data products (Michener et al., 1997, Michener, 2006, Papale et al., 2012, Kervin et al., 2013). Metadata allow for interpretation and integration of heterogeneous data obtained from different measurement approaches across disparate study sites, which occur even in well-organized science projects. Additionally, metadata are often critical for quality assurance and quality control (QA/QC). For example, particular equipment can have biases under certain conditions, or events such as power outages or equipment maintenance can affect data quality. Metadata that describe the location and time period of the observations or data products are used for aggregation both in time and space. Furthermore, metadata also describe the people who conducted the work, which is important for provenance (record of data credits) and proper attribution to data originators. Given its broad range of utility, metadata can describe many aspects of observations or data products, including descriptions of the measurement setting (e.g., measurement location and approach), the data

reported (e.g., measurement variable and units), and the datasets (e.g., data processing level and details).

Due to data management requirements from federal funding agencies, a variety of data collection repositories now exist, each with their own metadata requirements (e.g., KBase (Department of Energy Systems Biology Knowledgebase), n.d, KNB (The Knowledge Network for Biocomplexity), n.d, NOAA NCEI (National Centers for Environmental Information), n.d, USGS, n.d). Over the last several years, the digital preservation community has developed a general consensus around best practices for metadata that define how to reliably ingest data into these data repositories, track provenance, build and maintain metadata, and enable future consumers to independently access and use the data. For example, the Open Archival Information System (OAIS) reference model describes the concept of information packages as a collection of content and metadata. The metadata is further delineated as 1) content metadata, 2) descriptive metadata that enable search and retrieval of the content, 3) preservation description metadata necessary for long-term archiving such as provenance, checksums and unique identifiers, and 4) other ancillary metadata needed to define and hold the package together (OAIS/ISO (International Organization for Standards) 14721:2012). Some data repositories provide tools for data originators to prepare and submit a *Submission Information Package* (hence referred to as “data package”) containing content data and all the metadata, and for data consumers to download a *Dissemination Information Package* containing citation information in addition to the content data and metadata.

Several standards and formats currently exist to describe data collection, processing, and reporting for environmental data and promote interoperability between data repositories. Examples include the Open Geospatial Consortium “Observation and Measurements” standard for observations and sampling features (OGC, 2013, ISO (International Organization for Standards) 19156:2011, 2011), International Standards Organization/Federal Geographic Data Committee standards for geospatial (Federal Geographic Data Committee, 1998, ISO 19115-1:2014) and temporal metadata (ISO 8601), netCDF formats for climate and forecast metadata (Unidata, 2016), and the Ecological Metadata Language (EML; Michener et al., 1997, EML (Ecological Metadata Language) Project, 2009). Data information models built upon these standards describe content data and metadata standard formats and relationships, and are easily converted to searchable relational databases (Horsburgh et al., 2016). Data information models suitable for environmental data include Morpho (NCEAS, 2015) that is designed to interface smoothly with EML, and the Observational Data Model 2 (ODM2; Horsburgh et al., 2016). Data information models support a wide range of data types and enable data search, discovery, and synthesis. However, these models still require that additional standard data collection and naming protocols be defined and that metadata for both

observations as well as modeled products be collected in a standardized way before it can be ingested into the searchable database. Moreover, these models require the data originator to be proficient in data science terminology or concepts, and to expend significant additional effort into translating their data and notes into the required formats.

In contrast, other domain-specific templates and accompanying databases have been developed to enable easier reporting of data and metadata by data originators for ecophysiology, hydrology, and meteorology datasets. These efforts include forest plot inventories that collect forest census data like taxa identification, locations, causes of mortality, and size (e.g., Smithsonian Tropical Research Institute - Center for Tropical Forest Study (STRI-CTFS; Condit et al., 2014), CTFS-ForestGEO (CTFS Forest Global Earth Observatories; Anderson-Teixeira et al., 2014) and the Amazon Forest Inventory Network (RAINFOR (Amazon Forest Inventory Network), 2016, Malhi et al., 2002, Peacock et al., 2007)). The AmeriFlux/Biological, Ancillary, Disturbance and Metadata (BADM) protocol has been developed and implemented across several flux-based networks (e.g., AmeriFlux, FLUXNET, ICOS) (Law et al., 2008, AmeriFlux, 2016). AmeriFlux/BADM reporting templates focus primarily on ecosystem-level observations often aggregated in space and time to describe the area within a flux tower footprint. A variety of frameworks support regional and global data repositories, such as BiofuelEcophysiological Traits and Yields (BETYdb) Database (LeBauer et al., 2010), Sapfluxnet (Poyatos et al., 2016), and International Soil Carbon Network (ISCN) (ISCN, 2016). These frameworks are designed to capture metadata specific to their respective measurement types. However, the reporting templates do not necessarily conform to published standards, and are sometimes unstructured, making data synthesis, search within the data, and integration into a database difficult.

Thus, the existing data informational models are too complex for ecohydrological data originators to use directly, and none of the existing standardized data/metadata templates have the necessary structure to support reporting of the diverse observations required for earth system modeling. To bridge this gap between data information models and domain-specific data/metadata reporting templates, we developed a new metadata reporting framework, FRAMES (A Framework for Reporting dAta and Metadata for Earth Science). FRAMES is a set of templates that standardizes reporting of diverse ecohydrological data for synthesis across a range of spatiotemporal scales, and ultimately enables ingestion into a searchable data information model.

We conducted this work as part of an interdisciplinary team-based project whose overarching goal is “to develop a predictive understanding of how tropical forest carbon balance and climate system feedbacks will respond to changing environmental drivers over the 21st Century” (NGEE Tropics, 2016). By employing an iterative scientist-centered design approach, we identified and implemented features into FRAMES that support not only

environmental process understanding but also earth system model development. These features include 1) standardization and organization of metadata according to best data science practices, 2) a modular design that can expand to accommodate diverse measurements, 3) data entry formats that facilitate efficient metadata reporting, 4) a multiscale hierarchy that links observations across spatiotemporal scales, and 5) collection of metadata needed for model-data integration. Although extensible to various earth system data types, the first version of FRAMES described here is focused on primarily automated measurements collected by permanently located sensors, including sap flow (tree water use), leaf surface temperature, soil water content, dendrometry (stem diameter growth increment), and solar radiation. In addition to describing FRAMES, we discuss key challenges, solutions, lessons learned, and areas for future development that are broadly applicable to team-based projects and science networks.

## 2. Methods

Our team-based project supports a dedicated data team that is tightly integrated with an interdisciplinary group of earth scientists. The data team encompasses responsibilities of data manager and data distributor, and refers to persons assisting data originators in metadata and data reporting, preserving data, and making data available to consumers (Peng et al., 2016). The data team led the development of FRAMES by working closely with data originators (the empiricists collecting the observations), as well as data consumers (the empiricists and also modelers using the data and metadata).

We developed FRAMES to support the project's first coordinated data collection effort centered around tree responses to drought conditions in Central and South America during the El Niño Southern Oscillation (ENSO) event of 2015–2016. Prior to developing FRAMES, we identified relevant aspects of existing protocols and standards to use as design foundations including ISO standards (ISO 8601:2004, ISO 19115-1:2014), FGDC standards (FGDC, 1998), AmeriFlux/BADM templates (AmeriFlux, 2016), ISCN reporting templates (ISCN, 2016), STRI-CTFS protocols (Anderson-Teixeira et al., 2014, Condit et al., 2014), RAINFOR-GEM protocols (Marthens et al., 2014, RAINFOR (Amazon Forest Inventory Network), 2016), and Sapfluxnet (Poyatos et al., 2016).

The approach we used to develop FRAMES involved a combination of agile development principles and scientist-centered design (Ramakrishnan et al., 2014). Agile development uses short incremental development cycles with reassessment of priorities and solicitation of feedback after each cycle. The scientist centered-design process works closely with a group of researchers (data originators and data consumers) that provide direction and feedback throughout product development to define the desired end products. The process begins with extensive interviews to understand each participant's standard processes and workflows. It works to 1) understand data sources,

QA/QC needed, and development priorities; 2) develop data algorithms, and 3) build products that enable the science goals.

Based on requests from members of the project's science team, we focused our efforts on collecting metadata necessary to provide interpretation, cross-site comparison, and QA/QC for a prioritized list of ENSO observations. These observations were primarily automated measurements collected by permanently located sensors, including sap flow (tree water use), leaf surface temperature, soil water content, dendrometry (stem diameter growth increment), and solar radiation. Working closely with data originators and data consumers, we addressed one or two measurement types at a time, building out FRAMES as we added additional measurement types. Initial template designs were based on existing data collection protocols and informational interviews conducted with data originators to understand the measurement procedure, identify existing metadata collection methods, and discuss additional metadata collection. Through discussions with data originators and consumers as well as our expertise in data management, required metadata were distinguished from optional metadata based on which information was needed to interpret data, perform cross-site comparison, and conduct QA/QC assessment.

FRAMES was designed to fit as seamlessly as possible into the existing data collection processes of the data originators. We iteratively tested FRAMES with data originators, incorporating additional measurement types and feedback based on field metadata entry trials. Once we had tested FRAMES with four of the ENSO measurement types as well as location and equipment information, we solicited feedback from modelers (data consumers). We also conducted informational interviews with other data originators and consumers of anticipated measurement types (primarily sample-based observations including leaf water potential, gas exchange, and non-structural carbohydrates) to check for compatibility with FRAMES. To minimize the effort of data originators, we transferred information already submitted in previous versions of FRAMES to the newer versions throughout the iterative development.

Finally, FRAMES was designed to facilitate submission to data repositories, including the NGEE Tropics Archive, the project's data repository. The NGEE Tropics Archive has a web portal that allows data originators to upload and download data packages. The Archive is supported by a programmatic REST API built on top of Django Python web framework with an easy-to-use web user interface built with Foundation (Zurb, 2016) front-end framework. The Foundation front-end framework is flexible, highly customizable and provides support for responsive, light-weight HTML for mobile application support. Django is a fully featured open-source Python web application framework that supports rapid development. Django makes the low-level framework decisions so that the development is primarily focused on the application domain rather than composing the framework features. NGEE Tropics

Archive manages the data package by storing the data package metadata in a Postgres database and the data files on the local file system.

In general, completeness and accuracy of metadata submitted via FRAMES templates are considered to be the responsibility of the data originator, although the data team manually inspects data package submissions via the NGEET Tropics Archive portal. The peer-review process enabled by data-sharing provides input to data originators to make corrections to their data.

### 3. Results: a Framework for Reporting dAta and Metadata for Earth Science (FRAMES)

#### 3.1. Key requirements and characteristics of FRAMES

Through initial interviews, we identified key requirements of a metadata framework that would enable multisite comparisons of tree response to drought and testing of spatially explicit models. First, the framework had to support a variety of measurement types and data processing levels that were anticipated to be made and used throughout the project. Many of these measurement types shared similar metadata while some metadata was measurement specific. Secondly, the framework had to enable efficient data entry in recognition of the fact that metadata reporting is time consuming and can add significant overhead to a data originator's field collection and data reporting duties. Additionally, scientists needed the ability to use the data reported at various scales. For example, they wanted, on smaller scales to track multiple, co-located measurement types on a specific tree for assessment of plant trait co-variation, and on larger scales to track relationships across study sites. Finally, the framework had to support integration of data into carbon cycle models, which was identified as a top project priority.

Thus, FRAMES was designed to address these requirements, resulting in the following key characteristics: 1) Standardization and organization of metadata according to best data science practices (Section 3.2), 2) a modular organization in which data originators can report information about data file contents, measurement settings for a variety of observations, and high-level data descriptions and citation information (Section 3.3), 3) reporting formats designed to match existing data collection practices for efficient and streamlined metadata entry (Section 3.4), 4) the concept of a multiscale measurement position hierarchy to enable data aggregation and usage across scales (Section 3.5), and 5) incorporation of additional data and metadata fields that would normally not be collected as part of a field measurement, but were required for model-data integration (Section 3.6).

#### 3.2. Standardization and organization of metadata according to best data science practices

FRAMES uses concepts and terminology from preexisting standards, templates and databases, to support compatibility with external data formats and protocols. First, for sites with a pre-existing, widely-used



identifier such as an AmeriFlux/FLUXNET Site ID (AmeriFlux, 2016), we used the existing ID, to enable standardization with a global network of sites and cross-database search. Other site and plot metadata, including location information and descriptions, were collected directly from site leads or data originators (see Appendix B). The FGDC standard (FGDC, 1998) was supported for reporting spatial location metadata in different reference systems including geographic coordinates (for latitude/longitude representations), planar coordinates (for coordinate or distance/bearing representations), and vertical coordinates (for heights). All dates and timestamps had to be reported in ISO formats (ISO 8601), and a UTC offset specified. The AmeriFlux/BADM reporting templates (AmeriFlux, 2016) were used as a starting point for determining fields for equipment information, installation, and maintenance, as well as for the multiscale measurement position hierarchy (Section 3.5).

We also supported compatibility of certain domain-specific standard terminology when applicable. For example, we have largely adopted the taxonomic identification protocol and based our tree characteristics on the censusing protocols of STRI-CTFS (Anderson-Teixeira et al., 2014, Condit et al., 2014). Additionally, we leveraged RAINFOR-GEM's tree assessment protocols for the measurement of tree height and canopy illumination indices (Marthews et al., 2014, RAINFOR (Amazon Forest Inventory Network), 2016). For sap flow measurements, we consulted the AmeriFlux/BADM and Sapfluxnet protocols (AmeriFlux, 2016, CREAM (Centre for Research on Ecology and Forestry Applications), 2016, Poyatos et al., 2016). For soil water content and other soil-related observations, we consulted the AmeriFlux/BADM and ISCN data reporting templates (Law et al., 2008, AmeriFlux, 2016, ISCN (International Soil Carbon Network), 2016).

Besides the use of preexisting standards, FRAMES also incorporates other best data science practices including 1) standardization of variable names and file structure to enable automation of metadata extraction via scripts, 2) use of controlled vocabularies in drop down menus to facilitate comparability and search across sites, 3) use of descriptive data filenames and definition of data file contents, for example using header lines describing variables, and 4) tabular, row-based data entry templates with consistent column types (e.g. Borer et al., 2009, Hook et al., 2010, Tenopir et al., 2011).

FRAMES Workflow for Submission to Repository (NGEE Tropics Archive)

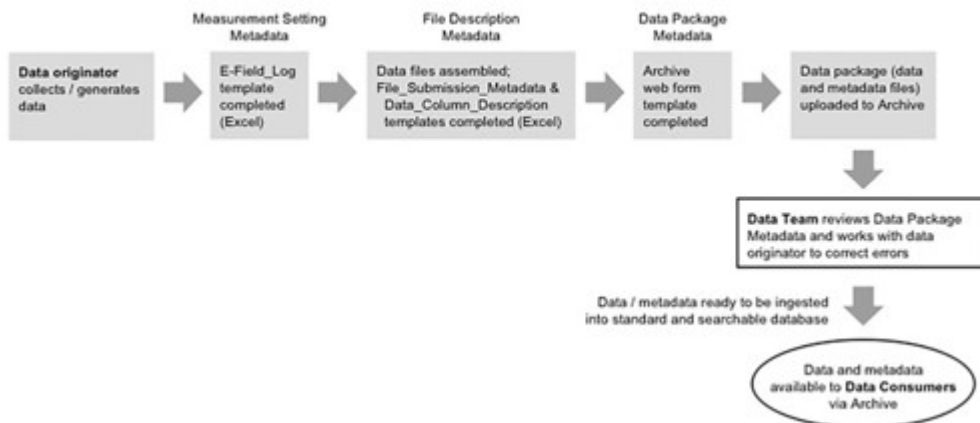


Fig. 1. FRAMES metadata and data package workflow. The Data Originator (grey boxes) collects/generates data and completes FRAME metadata templates (Section 3.3) that are included with data in a data package for submission to a repository (e.g., NGEE Tropics Archive). The Data Team (outlined box) reviews the data package before it is available to Data Consumers (outlined oval) via the Archive.

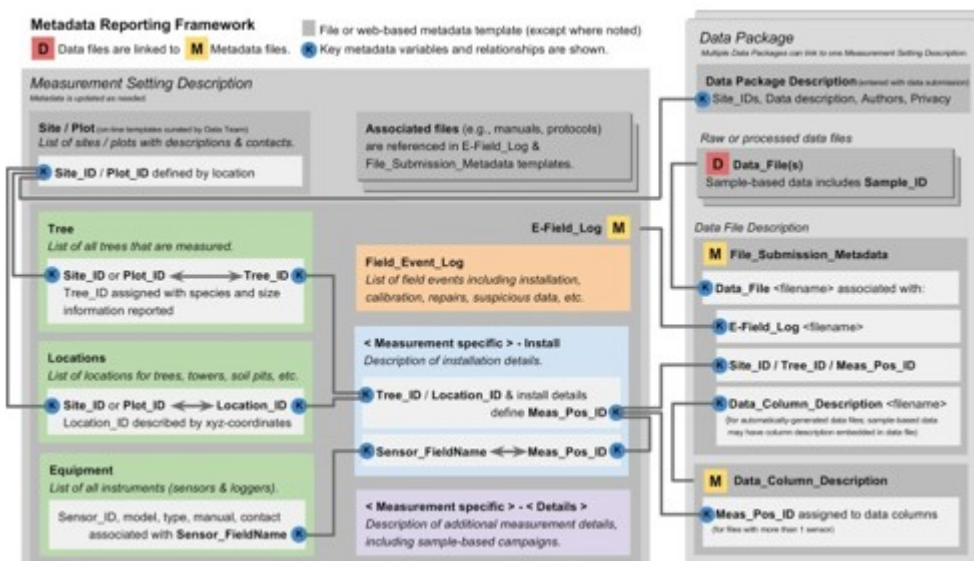


Fig. 2. FRAMES (Framework for Reporting dAta and Metadata for Earth Systems). Data originators enter measurement setting description and data file description metadata via templates files [M] that are linked to raw or processed data files [D]. Site/plot information (part of Measurement Setting Description) as well as Data Package Description are provided via on-line templates. Links between key variables [K] illustrate how the metadata templates work in tandem.

### 3.3. Modular metadata organization

FRAMES is organized into three main groups of related metadata: 1) descriptive information about a data package, 2) content information about the data file organization, and 3) content information about the data collection process and measurement settings. Physically, FRAMES comprises a set of Microsoft (MS) Excel spreadsheet files to describe file contents and measurement setting metadata, and package-level descriptive metadata reported in a web form (spreadsheet templates included in Appendix A, web

screenshots included in Appendix E). The metadata are bundled with data files into a data package and submitted to a data repository (e.g. the NGEE Tropics Data Archive) via a web form. The data reporting workflow is illustrated in Fig. 1, and an overview of FRAMES with relationships between the templates is illustrated in Fig. 2. With this combination of metadata files submitted to a data repository, FRAMES enables digital preservation of the entire data history, including digital reporting of critical information from field notes and raw data files generated by data loggers, to enable reproducibility of scientific analyses.

### 3.3.1. Data package description

FRAMES utilizes the concept of data packages, in which data originators bundle their content (data files) and corresponding content metadata information together for submission to a repository. A data package is often determined by a common theme or activity. Within our project, data packages are typically assembled to support an experiment or set of sensor observations, a data synthesis product, a publication, or a field campaign. A data package may contain many types of data associated with the theme or activity.

The data package description is a set of basic metadata fields that describe its contents and includes information necessary to obtain a unique Digital Object Identifier (DOI), as well as other information needed to identify the package for search and retrieval in the future. These metadata include data package names and descriptions, Site ID and Plot ID, authors, institutions, citations, acknowledgements, and funding sources, as well as QA/QC status (Appendix E). The metadata collected also describes access permissions for data usage. Required fields for the data package description were determined as the minimum set of information needed to obtain a DOI from Datacite (Datacite, 2016) via the U.S. Department of Energy's Office of Science and Technical Information (OSTI).

For the NGEE Tropics project, data are archived using the project's data repository NGEE Tropics Archive, which allows users to upload and access data packages. Currently, data originators can create, save, edit, and submit draft data packages via a web portal (Appendix E). Data originators provide descriptive metadata about the data package in a web form and can upload a single data file of any type (zipped file types allow for upload of multiple files). The web form enables data originators to reuse certain information, such as field site and plot information and person (name, email, institution) information to minimize inconsistent or erroneous data entry. For example, data originators only have to select the site name/ID for all related site information to be auto-populated, including spatial coordinates (numerically and via google maps), PI (principal investigator) information, and general site descriptions.

Once submitted, data package descriptions and data files are manually reviewed for completeness and accuracy as part of the project's archival

approval processes. After approval, data packages with appropriate citation information are made available via the web portal to data consumers who are assigned access privileges.

### 3.3.2. Data file descriptions (file submission metadata and data column description)

For each data file submitted, data originators report the following metadata in the MS Excel template “File Submission Metadata:” 1) Tree ID or other Location ID if applicable, 2) time period of the data and timestamp details (e.g., time zone and whether the timestamp is at the start, middle, or end of the sampling period), 3) data processing level with related processing approaches (e.g. raw, translated/processed, data originator QA/QC, project-level QA/QC), 4) references to the measurement setting description (e.g., E-field Log file)—this information is essential because it links the data to additional metadata reported in the separate templates described in Section 3.3 (see Fig. 2)—, and 5) references to data file descriptions (Data Column Description).

Additionally, for every data file, a corresponding “Data Column Description” template provides the information necessary to understand the data file. This is a semi-standardized template that includes information on header rows (e.g., those automatically generated by instrumentation), column names, units, data averaging (e.g., instantaneous or a mean/standard deviation over the sampling period), measurement type, and a location identifier (e.g., Tree ID, Measurement Position ID, or Sample ID) if multiple measurement positions are recorded in the same file. The location identifier is critical because it links the observations to installation details and other events affecting data quality that are described in the measurement setting templates. Data originators can configure the Data Column Description as a series of tabs in a single MS Excel file, a standalone file, or as a separate tab within the data files (if data file is MS Excel).

Table 1. Measurement setting description template groupings included in the E-field Log file.

<b>E-field Log template</b>	<b>Template description</b>
<b>Tree</b>	Description of observed trees, including species identification and an initial assessment of size and light environment. We include this information because our framework is designed for research in tropical forests. Long-term demographic (census) data is reported elsewhere.
<b>Locations</b>	Location (relative or absolute) information, geomorphology description, and contact information

<b>E-field Log template</b>	<b>Template description</b>
<b>Equipment</b>	<p>for features where observations are made. Example features include trees, towers, cranes, pits, and random observations points.</p> <p>Description of equipment used to make observations, including make, model, contact personnel, and reference to manuals.</p>
<b>Field Event Log</b>	<p>Description of field events that affect data collection and quality. Event examples include equipment installation, maintenance, calibration, and removal, as well as broad categories like “Suspicious Data” that capture events such as power outages or animal interference.</p>
<b>Measurement-specific Install</b>	<p>Detailed description of installation events specific to each measurement type that requires (semi-) permanently installed equipment. For example, a sap flow sensor installation event is recorded on the Field Event Log and the details of that installation, such as sensor height and probe depth illustrated in Fig. 3, are recorded on the SAP-Install template.</p>
<b>Measurement-specific Details</b>	<p>Detailed description of measurement specific information. These templates are designed to capture various types of measurement specific information not recorded on the Field Event Log. For example, leaf gas exchange and leaf water potential observations are conducted in campaigns. Details of the campaign are captured on the Leaf-Campaign template.</p>

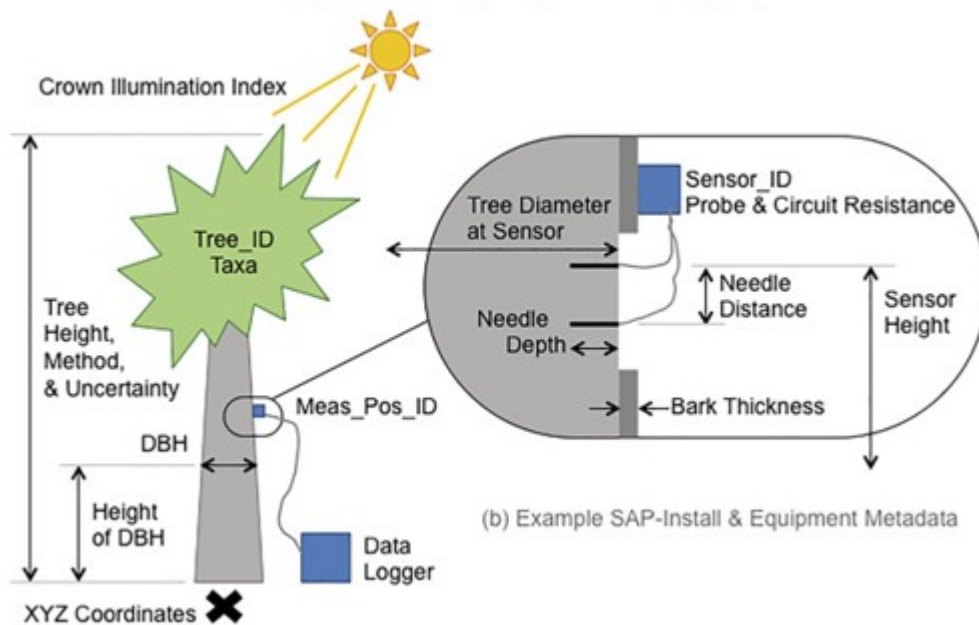


Fig. 3. Examples of (a) Tree and (b) SAP-Install and Equipment metadata variables that are reported as part of the measurement setting description. SAP = Sap flow; Meas\_Pos\_ID = Measurement Position ID; DBH = diameter at breast height.

### 3.3.3. Measurement setting description (E-field Log)

The measurement setting description contains information related to observations: 1) location; 2) equipment details, installation, and maintenance history; 3) approach and technicians; 4) events affecting data quality. We developed a standardized digital format for this information to which data originators could transfer their field notes. Because this information is complex and often hierarchical, we organized the information into a series of templates implemented as tabs in a MS Excel file “E-field Log” (Table 1). Key variables that link the templates together are shown in Fig. 2 (see Appendix C for full relational framework). All variables within each template are described in Appendix B. Examples of measurement setting description variables are illustrated for sap flow in Fig. 3.

### 3.4. Design features that maximize metadata reporting efficiency and data/metadata reuse

To maximize efficiency of reporting metadata and data reuse, we implemented several design features based on data originator interviews and observations of originators entering metadata on beta template versions.

FRAMES enables efficient data entry by being closely aligned with existing field practices as follows. The modular organization of FRAMES (Section 3.3) facilitates co-located entry of related metadata relevant to multiple measurement types or field sites/locations. One example occurs in the web form that data originators use to submit data packages to the project's

repository. Data originators are allowed to submit multiple data files associated with any number of sites and variables. Thus, originators can submit several related data files, for example those associated with a field campaign, in one data package, minimizing time spent on entering metadata and uploading files. As another example, in the measurement setting description spreadsheet (“E-field Log” file), details about measured trees as well as equipment specifications are reported once in the Tree and Equipment templates respectively. Co-location of the measurement setting templates in a single file allows for quick reference between location and equipment metadata when describing installation and other field events. Data originators can also report events that affect multiple measurements in a single entry in the E-field Log file. For example, a power outage affecting soil moisture and sap flow measurements can be reported as suspicious data in one line on the “Field Event Log” template with location and/or sensor identifiers indicated. Through translation of such suspicious data information—automated if desired—, data quality flags can be assigned to the affected data values.

We also intentionally separated the measurement setting description (E-field Log) from metadata describing the data package and data files to allow any data originator to link multiple data files to a single set of metadata templates in the E-field Log. Thus, data originators can submit the E-field Log as a separate data package into the data repository. This structure allows for the data and the measurement setting metadata to be maintained independently of each other, as the latter are typically updated on an infrequent basis. Furthermore it enables reuse of certain metadata across research studies and field sites. For example, two research groups collecting different observations at one or multiple sites can both reference the same E-field Log record in the data repository to share tree, location or equipment information. Finally, multiple types of data, for example raw, processed, or cross-site data synthesis products, can all be linked to the appropriate metadata templates.

Finally, we embedded instructional text and formatting cues to facilitate metadata entry. Within FRAMES, short instructions, metadata variable descriptions, and example entries are provided. Templates within the E-Field\_Log MS Excel file are color coded to indicate similar types of metadata: infrequently changing lists relevant to multiple measurement types, infrequently changing measurement-specific installation templates, and the Field Event Log that is updated at various frequencies. These colors matched highly visual instructional documentation (see Appendix A).

### 3.5. Multiscale measurement position hierarchy

We developed a multiscale measurement position hierarchy to account for the diverse spatial scales that observations represent and to reduce redundancies in reporting of various location identifiers (Fig. 4). In this hierarchy, a “Site” is the largest unit of study, and is assigned a unique Site

ID. We impose no limit on the physical size of a site, which can range from individual locales to regional areas and the entire globe; however, we anticipate most sites to be individual locales on the order of kilometers squared. Smaller “Plot” areas can occur within a site, and each plot has a unique Plot ID. Within a site or plot, a feature located in x-y space, including trees, towers, measurement pits, etc., is assigned a Location ID. Observations occurring repeatedly at a sub-location, e.g., at a specific height or bearing, are assigned a unique Measurement Position ID. Alternatively, observations obtained from a sample of the feature are assigned a unique Sample ID, which may have specific sub-location spatial information.

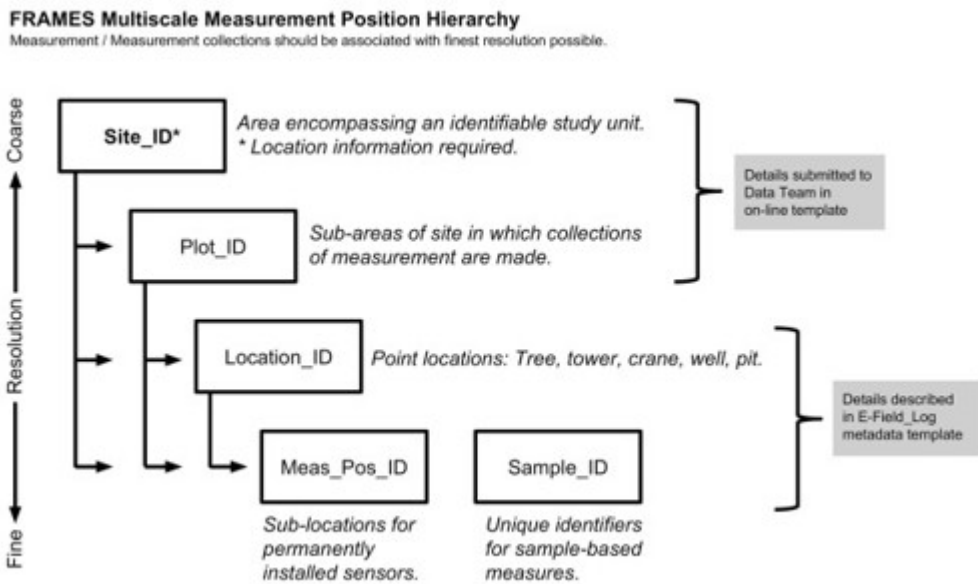


Fig. 4. FRAMES multiscale measurement position hierarchy. Observations including time series are associated with a unique measurement position identifier that may be at any hierarchy level. Any finer level identifier must be linked with at least a Site ID. Within our project focused on forest system, Tree ID is a type of Location ID.

Observations are linked to a unique spatial identifier in the hierarchy and inherit location information from the coarser levels to which that ID is linked. Aggregation to coarser resolutions is thus facilitated by combining all spatial identifiers that are linked to a particular coarser level location. For example, to aggregate individual sensors in a given Plot ID, all measurement position IDs associated with the Plot ID are combined. If multiple levels of locations are defined, an observation or observation time series is associated with the finest resolution spatial identifier defined; however, only Site ID is required. In this measurement position approach, sensors, either permanently installed or mobile, are linked to the appropriate spatial position identifier. Once Site or Plot metadata is collected, it is bundled with Location and Tree metadata (Section 3.3.3) for data originator and consumer reference.

### 3.6. Integration of field observations for model development



Integration of data with models requires translation of empirical observations into the units and time periods required for model inputs or for direct comparison with model output. For example, meteorological time series data, such as air temperature, solar radiation, precipitation, and vapor pressure deficit, are used as boundary conditions to drive earth system models at each time step. In model parameterization, functional characteristics, ideally based on field observations, are assigned to plant functional types (PFT), soil types, and other model components. These functional characteristics, or traits, such as photosynthetic capacity, minimum leaf water potential, and soil organic matter content, may vary with climate conditions, other site characteristics or plant functional traits, component age, or spatial position (e.g., canopy level or depth). For model benchmarking, model predictions through time — for example, size distributions and relative abundance of PFTs, sap flow, and soil water content — are compared to field observations. Field observations are also used to provide insight into modeled ecosystem, ecophysiological, and hydrological processes.

To support model-data integration, we designed FRAMES to capture model-relevant metadata, which are sometimes not collected as part of the data originator's field efforts. In particular, we focused on information to support parametrization and benchmarking of the Functionally Assembled Terrestrial Ecosystem Simulator (FATES) model, which is based on Community Land Model with Ecosystem Demography (CLM(ED); Fisher et al., 2015) and ED (Moorcroft et al., 2001). FATES is a vegetation model that is being developed and used by the project's modelers. In FATES, plant demography (birth, growth, and mortality processes of related plants within a defined area) is modeled with size- and plant functional type-specific responses to environmental conditions. By requiring that the tree height and species information be reported, FRAMES provides input data, like photosynthetic capacity, for FATES to model plant responses, like sap flow and leaf gas exchange. These modeled plant responses are then benchmarked against observed responses made on similar trees under similar environmental conditions. FRAMES ensures that modeled plant responses can be compared to observed responses by linking the measurements to required tree characteristic metadata via the Tree ID. For example, crown illumination index and tree height, which are typically not collected or reported with leaf-level or plant-level response measurements, are required metadata for each measured tree. FRAMES has formalized communication between field scientists and modelers by ensuring that critical information is collected in a standardized, usable way for FATES and similar earth system models, such as ED2 (Medvigy et al., 2009).

#### 4. Discussion: FRAMES applications in interdisciplinary team-based earth science

##### 4.1. Linking complex and diverse observations across spatiotemporal scales for data synthesis

Linking observations across spatiotemporal scales is necessary for earth system process understanding as well as model parameterization and benchmarking (Dietze et al., 2013). FRAMES enables such linkages via the multi-scale measurement position hierarchy, its modular structure, and metadata standardization.

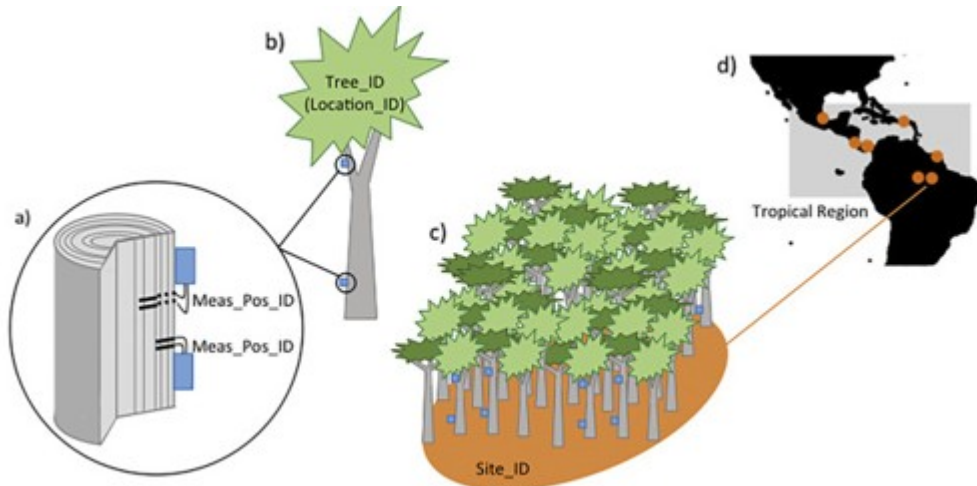


Fig. 5. Spatial scaling of sap flow measurement using multiscale measurement position hierarchy. Using measurement position identifiers that are linked to a common tree identifier, individual sap velocity measurements (a) made at multiple depths and positions on the tree can be processed with sapwood area or dendrometry measurements to characterize sap flow for the entire tree (b). (c) Aggregation across individuals within a single species or plant functional type (light or dark green trees separately) or across an entire site (light and dark trees combined) is enabled by tree identifiers that are linked to species/plant functional types and site identifiers. (d) Regional sap flow characterization can be synthesized by aggregating across site identifiers. As an example, Meas\_Pos\_IDs 00002A and 00002B (a) are linked to Tree\_ID 00002 (b) which in turn is linked to Site\_ID BR-Ma2 (c). If the tropical region (d) includes site BR-Ma2, then observations from Meas\_Pos\_IDs 00002A and 00002B or for Tree\_ID 00002 would be easily accessed for regional aggregation. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)

As a spatial example, sap flow is measured at the sub-tree level (Fig. 5a). Sap flow observations at multiple positions on the tree are used to determine the radial profile of sap flow within the sapwood and at different heights along a tree (e.g., trunk or branch). Integrating these measures yields an understanding of water use for a whole tree (Fig. 5b). The plant hydraulic functionality of FATES predicts tree water use for each combination of tree size and plant functional type. These model predictions can be benchmarked with whole tree water use of similar trees, as estimated from sap flow radial profile observations. Further aggregation at the site and regional scale enables benchmarking of site and regional model configurations, respectively (Fig. 5c-d). Synthesizing sap flow dynamics within a tree, for the whole tree, for groups of functionally related trees, and across the pantropical region enables improved understanding of ecohydrological processes in hyper-diverse tropical forests (Goldstein et al., 1998, Meinzer et al., 2001, Meinzer et al., 2004, Meinzer et al., 2005, Bell et al., 2015). The multiscale measurement position hierarchy facilitates such spatially extensive analyses because observations are defined by their position on the

landscape and are linked by unique measurement position, tree (location), and site identifiers. Additionally, the modular structure and standardization of FRAMES has enabled a pantropical sap flow synthesis effort involving several field sites (and hence many data packages). A data consumer independently automated 1) metadata ingestion from the templates, 2) integration of the metadata with the data files, and 3) additional data processing like removing duplicate timestamps (see Appendix F for R code).

Similarly, integration of observations across temporal scales is fundamental to understanding ecosystem processes (e.g., Detto et al., 2012). Furthermore, models that predict processes well across temporal scales remain elusive, i.e., models that perform well at fine scales (hourly or daily) often perform poorly at coarser scales (Dietze et al., 2011). Using FRAMES's description of data collection time resolutions and methods (e.g. discrete data or data averaged over a time intervals with the timestamp indicating the start, end or middle of the averaging time period), data consumers can temporally aggregate observations as required. For example, FATES predicts plant water fluxdynamics from sub-hourly to seasonal and inter-annual timescales, as driven by interactions between various plant hydraulic traits and environmental variation (as in Christoffersen et al., 2016). Using FRAMES these hydrodynamics may be benchmarked with sap flow data collected across project field sites at different sampling frequencies (10-, 15-, or 30-min) by aggregating to the desired model output time frequency (e.g. see Appendix F for R code that uses FRAMES metadata to automate this).

Alternatively, FATES hydrodynamic predictions can be benchmarked by tracking sub-hourly extremes like daily maximum sap flow, the timing of which is not known a priori, over periods of gradual declines in water availability, which occurred during the ENSO measurement campaign. Thus, by providing data at the finest resolution collected with the corresponding metadata to describe it, modelers have flexibility to customize model benchmarking to best assess a specific process. Additionally, analyses to understand covariation between sap flow and leaf surface temperature are highly sensitive to mismatches or drift in the timestamp. In conjunction with the description of time resolutions and methods, FRAMES includes a consistent reporting method for tracking timestamp drift by tracking the logger and CPU timestamps at data download events (Field Event Log in E-field Log Excel file in Appendix A).

Time-series data collected by different sensors at the same measurement position can be easily linked using the measurement position identifiers at any hierarchical level. For example, continuously measured leaf surface temperature is easily compared with sample-based measured leaf water potential measurements observed on the same tree via the tree identifier. Additionally, location information reported in FRAMES allows for linkages of spatially-explicit measures. For example, sap flow, leaf temperature, leaf water potential, and dendrometry measured on a specific tree can be simultaneously correlated with representative soil moisture conditions.

## 4.2. Expandability of FRAMES to accommodate diverse data

Data needed for earth system science, are not only diverse but also change as models and measurement techniques advance. Thus, the metadata reporting framework for such data must accommodate a variety of existing and new measurement types and approaches. FRAMES is modular to enable expansion to additional measurement types, beyond the few ecohydrological observations for which we have currently defined it.

A key aspect of the modular organization is separation of metadata reporting into three types of descriptions: data package, data file, and measurement setting. The data package description includes a minimal set of generic information, such as site identifier(s), data owner, and privacy settings. Similarly, the data file description is applicable to a wide variety of data types because it also contains generic metadata, like time step and data processing information. Data originators are not restricted to predefined measurement types and formats because the semi-open ended data file column description can describe the content of almost any type of data file. The modular organization of the measurement setting description also readily accommodates new measurement types because the core set of reporting templates (Tree, Equipment, Location, and Field Event Log) describe information relevant to most measurement types in earth system science. New measurement types utilize some or all of these core description templates, and if necessary, a measurement-specific template can be developed to report additional measurement-specific information (see Appendix D for an example of how to add a new measurement to FRAMES).

The modular expandability of FRAMES is similar and compatible with ODM2 (Horsburgh et al., 2016), in that metadata is bundled in related groups. The difference is that ODM2 is a database structure for standardized metadata and data protocols. FRAMES operationalizes such a data structure as a reporting mechanism. In other words, data reported via FRAMES can be translated to a standardized format for assimilation into a database. This pre-database, standard-compatible flexibility differentiates FRAMES from other existing frameworks such as AmeriFlux/BADM, ISCN, and Sapfluxnet, which collect metadata and data in a standardized protocol designed for direct database assimilation.

We took this flexible approach for two reasons. First, it accommodates the needs of data originators by removing barriers to metadata and data sharing, such as the effort required to convert data to specific units and formats. Secondly, via the Data Column Description template which accommodates most types of data files, the flexible approach allows for archiving of raw data directly from loggers. Archiving unaltered data in its original format provides the full history of a data product for repeatability and data quality assessment (measurement errors as well as data processing errors). Archiving the entire data history is not only good science practice

(Dietze et al., 2013, Michener, 2015), but is also important for synthesizing data across sites and approaches because common and transparent processing approaches facilitate comparability. An additional advantage of this flexible approach is that data originators and consumers can assimilate data into variety of databases. A key component of this flexibility is achieved by separating the data column description from the data file description so that the data column description can be customized to the specific data file.

#### 4.3. Lessons learned and future development

FRAMES has supported data package reporting for six core NGEE Tropics field sites in Brazil, Panama, and Puerto Rico across six measurement types. Portions of the templates have also been used broadly in additional data reporting. Information about sensors, approaches, and installation details have informed development of a common sap flow processing approach for a synthesis of sap flow data across nine study sites. Additionally, the uniformity of the reported data enabled a data consumer to, on his own, automate processing of sap flow measurements for model benchmarking (see Appendix F for R code).

The use of FRAMES for the initial NGEE Tropics data collection effort has enabled us to gather feedback regarding what is working and what is not. The most valuable feedback was the effort that six data originators were willing to exert in using FRAMES to archive their data in the project's repository within a few months after the templates were finalized. We attribute this success largely to the scientist-centered design approach, which allowed us to identify data collection processes and design FRAMES to match the scientific goals and practices of both data originators and consumers. Anecdotally, data originators have reported FRAMES useful in organizing their field data. Subsequent data analyses, for example assessing co-dependent physiological responses measured from different sensors on the same tree, has been facilitated by the fact that all relevant information regarding the measurements is organized centrally within the metadata templates and that the tree ID clearly identifies measurements made on the same tree. Furthermore, FRAMES helped data originators to collect important ancillary information (e.g., tree height, diameter, crown illumination index) in conjunction with scheduled field activities rather than requesting the information at a later time, which would require additional field site visits if the measurement could still be made.

Developing an adaptable and efficient reporting framework was necessary for data synthesis across diverse observations, but its complexity has disadvantages. Understanding the modular templates and linkages seemed overwhelming at first to several of our data originators. Thus, further investigation of the instructional features is needed to ascertain and improve their efficacy. We found that the majority of time costs were upfront due to learning the structure of the framework and entering the measurement setting descriptions. However, since most measurement setting information

remains fairly static and is entered in a single template, maintaining the measurement setting description required minimal effort because only infrequent updates were required. For example, once the metadata for equipment and trees were entered, they remained the same over large periods of time, as observations were accumulated and/or new measurements were added.

A potential limitation to the framework is due to the efficient reporting mechanism designed to make reporting easier for data originators. FRAMES does not specify data variable names, units, or formats, which are required for database assimilation. Using FRAMES, reported data can be translated into a standardized protocol for database assimilation, as exemplified by similar case of automation of sap flow processing by a data consumer. The outstanding questions are 1) whether this reporting approach will ultimately result in improved availability of data with accompanying high quality metadata, and 2) what the tradeoffs are in terms of person-hours and who bears that cost—the data originator or dedicated data team personnel. We prioritized reporting formats in FRAMES to maximize reporting efficiency because although improving, the generally low quantity of shared data and poor quality of metadata is problematic in the earth sciences (Tenopir et al., 2011, Kervin et al., 2013, Michener, 2015).

Finally, we implemented several templates in MS Excel because of its ubiquity, operating system neutrality (i.e., it runs on Macs and PCs), copy/paste functionality, and off-line access for remote areas with poor Internet. However, MS Excel is not ideal for selection from a controlled vocabulary menu, collaborative data entry, customization of measurement types, real-time automated data quality verification, and machine readability. The use of MS Excel also makes it cumbersome to release new versions of the templates and ensure backwards compatibility with previous files that were submitted. Additionally, separation of metadata in template files currently requires that the data consumer manage separate sources of metadata information and download different data packages for synthesis efforts. The standardization of metadata alleviates some aspects of this limitation by enabling the data consumer to programmatically link the data and metadata (Section 4.1). As others have reported, new software tools are needed (Michener, 2015), in our case, tools that merge the functionality of MS Excel and eliminate these limitations. Possibilities include web-based or mobile tools that are available offline, can be written to appropriate output formats (e.g., comma-delimited ascii, NetCDF/HDF5, EML, or JSON files), and are customizable to originator preferences and measurement types (e.g., Jones et al., 2007, McIntosh et al., 2007). In the future, we intend that the metadata and data be ingested into a relational database (using a framework like ODM2) to facilitate programmatic data integration, searchability and easy data manipulation, such as sub-setting and aggregation.

## 5. Conclusions

We developed FRAMES, a set of online web forms and Excel-based metadata templates that position data and metadata for easier entry into an operational data repository. FRAMES is designed to facilitate and improve capture of desired metadata for ecohydrological observations, including information about how measurements were conducted, data file contents, and high-level descriptive metadata for citation and attribution. Thus, FRAMES enables synthesis of diverse ecohydrological and biogeochemical observations for study of earth system processes and for integration with predictive earth system models.

The overarching challenges for synthesizing diverse earth system observations were 1) developing a metadata framework that allowed experts to share data with team members from other disciplines, and 2) collecting sufficient metadata to organize and process data comparably across sites and measurement methods. FRAMES incorporates several key features that addresses these challenges and supports interdisciplinary team-based earth system science, including 1) compatibility with standard data protocols, and conformance with data science best practices that enable data interpretation, comparison of observations across sites and approaches, and QA/QC, 2) a modular design that accommodates diverse data types and can expand as required by measurement and model advancement, 3) compatibility of existing field practices to maximize data and metadata reporting efficiency, 4) a multi-scale measurement position hierarchy and comprehensive time step descriptions that facilitate spatiotemporal aggregation and linkage of measurement types for synthesis, and 5) targeted metadata collection that enables model-data integration.

To date, FRAMES templates have been used, in whole or in part, for several submissions to the NGEE Tropics Data repository. An iterative scientist-centered design was central to the successful use of FRAMES within our project, where the goal is to improve a predictive understanding of carbon cycling in tropical forests under climate change. As an interdisciplinary data team of ecologist, hydrologists, and data scientists working closely with data originators and consumers throughout the development process, we were able to identify features critical to the project's science needs and develop pragmatic solutions. This integrated data science approach will underpin further improvement to FRAMES, and we recommend it as a model for harnessing complex and diverse data inherent in team-science and observational networks.

Additionally, FRAMES promotes good data management practices that benefits both data originators and consumers by 1) digitally preserving data with adequate metadata documentation, 2) enabling sharing with the broader community with appropriate citation and attributions, 3) facilitating interoperability with other databases, and 4) broadening data use and reuse for purposes that stretch beyond the initial intentions of the data collection effort (particularly for use in earth system models). Next steps involve making improvements to FRAMES based on data originator and

consumer feedback, and extraction of information in data packages into a queryable database that enables programmatic search, discovery, and processing of data.

### Acknowledgments

We thank the larger Next Generation Ecosystem Experiments-Tropics (NGEE-Tropics) team for helpful feedback throughout the development process. Additionally, we thank Cory Snavely for background on digital datapreservation concepts and terminologies. This research was supported as part of NGEE-Tropics, funded by the U.S. Department of Energy, Office of Science, Office of Biological and Environmental Research under contract no. DE-AC02-05CH11231. We acknowledge institutional support from the Central Office of the Large Scale Biosphere Atmosphere Experiment in Amazonia (LBA) and the National Institute of Amazonian Research (INPA).

### References

- AmeriFlux <http://ameriflux.lbl.gov/data/badm-data-templates/> (2016) (accessed November 18, 2016)
- K.J. Anderson-Teixeira, S.J. Davies, A.C. Bennett, E.B. Gonzalez-Akre, H.C. Muller-Landau, S. Joseph Wright, K. Abu Salim, A.M. Almeyda Zambrano, A. Alonso, J.L. Baltzer, Y. Basset, N.A. Bourg, E.N. Broadbent, W.Y. Brockelman, S. Bunyavejchewin, D.F.R.P. Burslem, N. Butt, M. Cao, D. Cardenas, G.B. Chuyong, K. Clay, S. Cordell, H.S. Dattaraja, X. Deng, M. Detto, X. Du, A. Duque, D.L. Erikson, C.E.N. Ewango, G.A. Fischer, C. Fletcher, R.B. Foster, C.P. Giardina, G.S. Gilbert, N. Gunatilleke, S. Gunatilleke, Z. Hao, W.W. Hargrove, T.B. Hart, B.C.H. Hau, F. He, F.M. Hoffman, R.W. Howe, S.P. Hubbell, F.M. Inman-Narahari, P.A. Jansen, M. Jiang, D.J. Johnson, M. Kanzaki, A.R. Kassim, D. Kenfack, S. Kibet, M.F. Kinnaird, L. Korte, K. Kral, J. Kumar, A.J. Larson, Y. Li, X. Li, S. Liu, S.K.Y. Lum, J.A. Lutz, K. Ma, D.M. Maddalena, J.-R. Makana, Y. Malhi, T. Marthens, R. Mat Serudin, S.M. McMahon, W.J. McShea, H.R. Memiaghe, X. Mi, T. Mizuno, M. Morecroft, J.A. Myers, V. Novotny, A.A. de Oliveira, P.S. Ong, D.A. Orwig, R. Ostertag, J. Ouden, G.G. Parker, R.P. Phillips, L. Sack, M.N. Sainge, W. Sang, K. Sri-ngernyuang, R. Sukumar, I.-F. Sun, W. Sungsalee, H.S. Suresh, S. Tan, S.C. Thomas, D.W. Thomas, J. Thompson, B.L. Turner, M. Uriarte, R. Valencia, M.I. Vallejo, A. Vicentini, T. Vrška, X. Wang, X. Wang, G. Weiblen, A. Wolf, H. Xu, S. Yap, J. Zimmerman **CTFS-ForestGEO: a worldwide network monitoring forests in an era of global change** *Glob. Chang. Biol.*, 21 (2014), pp. 528-549, 10.1111/gcb.12712
- D.M. Bell, E.J. Ward, A.C. Oishi, R. Oren, P.G. Flikkema, J.S. Clark **A state-space modeling approach to estimating canopy conductance and associated uncertainties from sap flux density data** *Tree Physiol.*, 35 (2015), pp. 792-802, 10.1093/treephys/tpv041



- E.T. Borer, E.W. Seabloom, M.B. Jones, M. Schildhauer **Some simple guidelines for effective data management** Bull. Ecol. Soc. Am., 90 (2009), pp. 205-214, 10.1890/0012-9623-90.2.205
- B.O. Christoffersen, M. Gloor, S. Fauset, N.M. Fyllas, D.R. Galbraith, T.R. Baker, B. Kruijt, L. Rowland, R.A. Fisher, O.J. Binks, S. Sevanto, C. Xu, S. Jansen, B. Choat, M. Mencuccini, N.G. McDowell, P. Meir **Linking hydraulic traits to tropical forest function in a size-structured and trait-driven model (TFS v.1-Hydro)** Geosci. Model Dev., 9 (2016), pp. 4227-4255, 10.5194/gmd-9-4227-2016
- R. Condit, S. Lao, A. Singh, S. Esufali, S. Dolins **Data and database standards for permanent forest plots in a global network** For. Ecol. Manag., 316 (2014), pp. 21-31, 10.1016/j.foreco.2013.09.011
- CREAF (Centre for Research on Ecology and Forestry Applications), 2016  
CREAF (Centre for Research on Ecology and Forestry Applications)  
<https://github.com/sapfluxnet/sapfluxnet-public/wiki> (2016) (accessed June 22, 2016)
- Datacite <https://www.datacite.org/> (2016) (accessed July 2016)
- M. Detto, A. Molini, G. Katul, P. Stoy, S. Palmroth, D. Baldocchi **Causality and persistence in ecological systems: a nonparametric spectral granger causality approach** Am. Nat., 179 (2012), pp. 524-535, 10.1086/664628
- M.C. Dietze, R. Vargas, A.D. Richardson, P.C. Stoy, A.G. Barr, R.S. Anderson, M.A. Arain, I.T. Baker, T.A. Black, J.M. Chen, P. Ciais, L.B. Flanagan, C.M. Gough, R.F. Grant, D. Hollinger, R.C. Izaurralde, C.J. Kucharik, P. Lafleur, S. Liu, E. Lokupitiya, Y. Luo, J.W. Munger, C. Peng, B. Poulter, D.T. Price, D.M. Ricciuto, W.J. Riley, A.K. Sahoo, K. Schaefer, A.E. Suyker, H. Tian, C. Tonitto, H. Verbeeck, S.B. Verma, W. Wang, E. Weng **Characterizing the performance of ecosystem models across time scales: a spectral analysis of the North American Carbon Program site-level synthesis** J. Geophys. Res. Biogeosci., 116 (2011), p. G04029, 10.1029/2011JG001661
- M.C. Dietze, D.S. Lebauer, R. Kooprt **On improving the communication between models and data** Plant Cell Environ., 36 (2013), pp. 1575-1585, 10.1111/pce.12043
- EML (Ecological Metadata Language) Project **EML version 2.1.1**  
<https://knb.ecoinformatics.org/#external//emlparser/docs/index.html> (2009) (accessed 19 June 2016)
- Federal Geographic Data Committee. FGDC-STD-001-1998. Content standard for digital geospatial metadata (revised June 1998). Federal Geographic Data Committee. Washington, D.C.
- R.A. Fisher, S. Muszala, M. Verteinstein, P. Lawrence, C. Xu, N.G. McDowell, R.G. Knox, C. Koven, J. Holm, B.M. Rogers, A. Spessa, D. Lawrence, G. Bonan **Taking off the training wheels: the properties of a dynamic**

**vegetation model without climate envelopes, CLM4.5(ED)** Geosci. Model Dev., 8 (2015), pp. 3593-3619, 10.5194/gmd-8-3593-2015

FLUXNET <http://fluxnet.fluxdata.org/data/fluxnet2015-dataset/> (2016) (accessed 20 December 2016)

G. Goldstein, J.L. Andrade, F.C. Meinzer, N.M. Holbrook, J. Cavelier, P. Jackson, A. Celis **Stem water storage and diurnal patterns of water use in tropical forest canopy trees** Plant Cell Environ., 21 (1998), pp. 397-406, 10.1046/j.1365-3040.1998.00273.x

L.A. Hook, S. Santhana-Vannen, T.W. Beaty, R.B. Cook **Best Practices for Preparing Environmental Data Sets to Share and Archive** Oak Ridge National Laboratory Distributed Active Archive Center (2010)

J.S. Horsburgh, A.K. Aufdenkampe, E. Mayorga, K.A. Lehnert, L. Hsu, L. Song, A.S. Jones, S.G. Damiano, D.G. Tarboton, D. Valentine, I. Zaslavsky, T. Whiteman **Observations Data Model 2: a community information model for spatially discrete Earth observations** Environ. Model. Softw., 79 (2016), pp. 55-74, 10.1016/j.envsoft.2016.01.010

M.O. Hunter, M. Keller, D. Morton, B. Cook, M. Lefsky, M. Ducey, S. Saleska, R. C. de Oliveira, J. Schiette **Structural dynamics of tropical moist forest gaps** PLoS One, 10 (2015), p. e0132144, 10.1371/journal.pone.0132144

ISCN (International Soil Carbon Network)

<http://iscn.fluxdata.org/data/dataset-information/> (2016) (accessed April 18, 2016)

ISO 8601:2004, Data elements and interchange formats — Information interchange — Representation of dates and times, <https://www.iso.org/obp/ui/#iso:std:iso:8601:ed-3:v1:en>.

ISO 14721:2012, Space data and information transfer systems -- Open archival information system (OAIS) -- Reference model.

ISO 19115-1:2014, Geographic information — Metadata — Part 1: Fundamentals, <https://www.iso.org/obp/ui/#iso:std:iso:19115:-1:ed-1:v1:en>.

ISO (International Organization for Standards) 19156:2011 **Geographic Information - Observations and Measurements** (2011), 10.13140/2.1.1142.3042

C. Jones, C. Blanchette, M. Brooke, J. Harris, M. Jones, M. Schildhauer **A metadata-driven framework for generating field data entry interfaces in ecology** Eco. Inform., 2 (2007), pp. 270-278, 10.1016/j.ecoinf.2007.06.005

KBase (Department of Energy Systems Biology Knowledgebase) <http://kbase.us>

(accessed July 2016)

K. Kervin, W. Michener, R. Cook **Common errors in ecological data sharing** JESLIB (2013), pp. 1-15, 10.7191/jeslib.2013.1024

KNB (The Knowledge Network for Biocomplexity)  
<https://knb.ecoinformatics.org> (accessed April 2016)

B.E. Law, T. Arkebauer, J.L. Campbell, J. Chen, O. Sun **Terrestrial Carbon Observations: Protocols for Vegetation Sampling and Data Submission** FAO (2008)

D. LeBauer, M. Dietze, R. Kooper, S. Long, P. Mulrooney, G.S. Rohde, D. Wang **Biofuel Ecophysiological Traits and Yields Database (BETYdb)** Energy Biosciences Institute, University of Illinois at Urbana-Champaign (2010), 10.13012/J8H41PB9

Y. Malhi, O.L. Phillips, J. Lloyd, T. Baker, J. Wright, S. Almeida, L. Arroyo, T. Frederiksen, J. Grace, N. Higuchi, T. Killeen, W.F. Laurance, C. Leñaño, S. Lewis, P. Meir, A. Monteagudo, D. Neill, P. Núñez Vargas, S.N. Panfil, S. Patiño, N. Pitman, C.A. Quesada, A. Ruelas, R. Salomão, S. Saleska, N. Silva, M. Silveira, W.G. Sombroek, R. Valencia, R. Vásquez Martínez, I.C.G. Vieira, B. Vinceti **An international network to monitor the structure, composition and dynamics of Amazonian forests (RAINFOR)** J. Veg. Sci., 13 (2002), pp. 439-450, 10.1111/j.1654-1103.2002.tb02068.x

T.R. Marthens, T. Riutta, I. Oliveras Menor, R. Urrutia, S. Moore, D. Metcalfe, Y. Malhi, O. Phillips, W. Huaraca Huasco, M. Ruiz Jaén, C. Girardin, N. Butt, R. Cain, RAINFOR and GEM networks **Measuring Tropical Forest Carbon Allocation and Cycling: A RAINFOR-GEM Field Manual for Intensive Census Plots (v3.0). Manual** Global Ecosystems Monitoring Network (2014)  
<http://gem.tropicalforests.ox.ac.uk/>

A.C.S. McIntosh, J.B. Cushing, N.M. Nadkarni, L. Zeman **Database design for ecologists: composing core entities with observations** Eco. Inform., 2 (2007), pp. 224-236, 10.1016/j.ecoinf.2007.07.003

B.E. Medlyn, A.P. Robinson, R. Clement, R.E. McMurtrie **On the validation of models of forest CO<sub>2</sub> exchange using eddy covariance data: some perils and pitfalls** Tree Physiol., 25 (2005), pp. 839-857, 10.1093/treephys/25.7.839

D. Medvigy, S.C. Wofsy, J.W. Munger, D.Y. Hollinger, P.R. Moorcroft **Mechanistic scaling of ecosystem function and dynamics in space and time: Ecosystem Demography model version 2** J. Geophys. Res., 114 (2009), p. G01002, 10.1029/2008JG000812

F.C. Meinzer, G. Goldstein, J.L. Andrade **Regulation of water flux through tropical forest canopy trees: do universal rules apply?** Tree Physiol., 21 (2001), pp. 19-26, 10.1093/treephys/21.1.19

F.C. Meinzer, S.A. James, G. Goldstein **Dynamics of transpiration, sap flow and use of stored water in tropical forest canopy trees** Tree Physiol., 24 (2004), pp. 901-909, 10.1093/treephys/24.8.901

F.C. Meinzer, B.J. Bond, J.M. Warren, D.R. Woodruff **Does water transport scale universally with tree size?** Funct. Ecol., 19 (2005), pp. 558-565, 10.1111/j.1365-2435.2005.01017.x

W.K. Michener **Meta-information concepts for ecological data management** Eco. Inform., 1 (2006), pp. 3-7, 10.1016/j.ecoinf.2005.08.004

W.K. Michener **Ecological data sharing** Eco. Inform., 29 (2015), pp. 33-44, 10.1016/j.ecoinf.2015.06.010

W.K. Michener, J.W. Brunt, J.J. Helly, T.B. Kirchner, S.G. Stafford **Nongeospatial metadata for the ecological sciences** Ecol. Appl., 7 (1997), pp. 330-342, 10.1890/1051-0761(1997)007[0330:NMFTEs]2.0.CO;2

P.R. Moorcroft, G.C. Hurtt, S.W. Pacala **A method for scaling vegetation dynamics: the ecosystem demography model (ED)** Ecol. Monogr., 71 (2001), pp. 557-585, 10.2307/3100036

NCEAS **Morpho 1.11.0 user guide** <https://knb.ecoinformatics.org/software/dist/MorphoUserGuide.pdf> (2015) (accessed 13 April 2016)

NGEE Tropics <http://eesa.lbl.gov/ngee-tropics/> (2016) (accessed 20 December 2016)

NOAA NCEI (National Centers for Environmental Information)  
<https://www.nodc.noaa.gov>

(accessed November 2016)

OGC **Open Geospatial Consortium (OGC) observations and measurements v2.0 OGC document 10-004r1**

<http://www.opengis.net/doc/AS/OM/2.0> (2013) (also published as ISO/DIS 19156:2010, Geographic information — Observations and Measurements)

D. Papale, D.A. Agarwal, D. Baldocchi, R.B. Cook, J.B. Fisher, C. van Ingen **Database maintenance, data sharing policy, collaboration**

M. Aubinet, T. Vesala, D. Papale (Eds.), Eddy Covariance: A Practical Guide to Measurement and Data Analysis, Springer Netherlands, Dordrecht (2012), pp. 399-424, 10.1007/978-94-007-2351-1

J. Peacock, T.R. Baker, S.L. Lewis, G. Lopez Gonzalez, O.L. Phillips **The RAINFOR database: monitoring forest biomass and dynamics** J. Veg. Sci., 18 (2007), pp. 535-542, 10.1111/j.1654-1103.2007.tb02568.x

G. Peng, N.A. Ritchey, K.S. Casey, E.J. Kearns, J.L. Privette **Scientific stewardship in the Open Data and Big Data era—roles and responsibilities of stewards and other major product stakeholders** D-Lib Magazine, 13 (2016), pp. 1-25, 10.1080/02757259509532294

T.L. Powell, D.R. Galbraith, B.O. Christoffersen, A. Harper, H.M.A. Imbuzeiro, L. Rowland, S. Almeida, P.M. Brando, A.C.L. da Costa, M.H. Costa, N.M. Levine, Y. Malhi, S.R. Saleska, E. Sotta, M. Williams, P. Meir, P.R. Moorcroft **Confronting model predictions of carbon fluxes with measurements of Amazon forests subjected to experimental drought** *New Phytol.*, 200 (2013), pp. 350-365, 10.1111/nph.12390

R. Poyatos, V. Granda, R. Molowny-Horas, M. Mencuccini, K. Steppe, J. Martínez-Vilalta **SAPFLUXNET: towards a global database of sap flow measurements** *Tree Physiol.*, 36 (2016), pp. 1449-1455, 10.1093/treephys/tpw110

RAINFOR (Amazon Forest Inventory Network) **Liana and canopy index protocol** [http://www.rainfor.org/upload/ManualsEnglish/crown%20liana%20protocols\\_Sep%202014\\_EN.pdf](http://www.rainfor.org/upload/ManualsEnglish/crown%20liana%20protocols_Sep%202014_EN.pdf) (accessed March 2016)

L. Ramakrishnan, S. Poon, V. Hendrix, D. Gunter, G.Z. Pastorello, D. Agarwal **Experiences With User-centered Design for the Tigres Workflow API, Presented at the 2014 IEEE 10th International Conference on e-Science (e-Science)** *IEEE* (2014), pp. 290-297, 10.1109/eScience.2014.56

C. Tenopir, S. Allard, K. Douglass, A.U. Aydinoglu, L. Wu, E. Read, M. Manoff, M. Frame **Data sharing by scientists: practices and perceptions** *PLoS One*, 6 (2011), pp. e21101-e21121, 10.1371/journal.pone.0021101

Unidata <https://www.unidata.ucar.edu/software/netcdf/> (2016) (accessed March 2016)

USGS **Science data catalog** <http://data.usgs.gov/> (accessed November 2016)

A.P. Walker, P.J. Hanson, M.G. De Kauwe, B.E. Medlyn, S. Zaehle, S. Asao, M.C. Dietze, T. Hickler, C. Huntingford, C.M. Iversen, A.K. Jain, M. Lomas, Y. Luo, H.R. McCarthy, W.J. Parton, I.C. Prentice, P.E. Thornton, S. Wang, D. Wårlind, E. Weng, J.M. Warren, F.I. Woodward, R. Oren, R.J. Norby **Comprehensive ecosystem model-data synthesis using multiple data sets at two temperate forest free-air CO<sub>2</sub> enrichment experiments: model performance at ambient CO<sub>2</sub> concentration** *J. Geophys. Res. Biogeosci.*, 119 (2014), pp. 937-964, 10.1002/(ISSN)2169-8961

Zurb <http://foundation.zurb.com/> (2016) (accessed October 2016)