

1 GriddingMachine, a database and software
2 for Earth system modeling at global and
3 regional scales

4 Yujie Wang^{1*}, Philipp Köhler¹, Renato K.
5 Braghiere^{2,3}, Marcos Longo^{2,4}, Russell Doughty^{1,5}, A.
6 Anthony Bloom¹ and Christian Frankenberg^{1,2}

7 ¹Division of Geological and Planetary Sciences, California
8 Institute of Technology, 1200 E California Blvd, Pasadena, 91125,
9 California, USA.

10 ²Jet Propulsion Laboratory, California Institute of Technology,
11 4800 Oak Grove Dr, Pasadena, 91109, California, USA.

12 ³Joint Institute for Regional Earth System Science and
13 Engineering, University of California, Los Angeles, 405 Hilgard
14 Avenue, Los Angeles, 90095, California, USA.

15 ⁴Climate and Ecosystem Sciences Division, Lawrence Berkeley
16 National Laboratory, 1 Cyclotron Rd, Berkeley, 94720, California,
17 USA.

18 ⁵College of Atmospheric & Geographic Sciences, University of
19 Oklahoma, 660 Parrington Oval, Norman, 73019, Oklahoma, USA.

20 *Corresponding author(s). E-mail(s): wuyjie@caltech.edu;
21 Contributing authors: philipp.koehler@eumetsat.int;
22 renato.k.braghiere@jpl.nasa.gov; mlongo@lbl.gov;
23 russell.doughty@ou.edu; alexis.a.bloom@jpl.nasa.gov;
24 cfranken@caltech.edu;

Abstract

Land and Earth system modeling is moving towards more explicit biophysical representations, requiring increasing variety of datasets for initialization and benchmarking. However, researchers often have difficulties in identifying and integrating non-standardized datasets from various sources. We aim towards a standardized database and one-stop distribution method of global datasets. Here, we present the GriddingMachine as (1) a database of global-scale datasets commonly used to parameterize or benchmark the models, from plant traits to vegetation indices and geophysical information and (2) a cross-platform open source software to download and request a subset of datasets with only a few lines of code. The GriddingMachine datasets can be accessed either manually through traditional HTTP, or automatically using modern programming languages including Julia, Matlab, Octave, Python, and R. The GriddingMachine collections can be used for any land and earth modeling framework and ecological research at the regional and global scales, and the number of datasets will continue to grow to meet the increasing needs of research communities.

Keywords: Dataset, Earth System Modeling, Julia, Matlab, Octave, Python, R

Introduction

Land components in Earth system models (ESMs) are moving towards more explicit biophysical representation. For example, in the modeling of the soil-plant-air continuum, vegetation canopy complexity has increased from a single big-leaf to a multi-layered modeling approach including leaf angular distributions [1], thus allowing for more realistic and predictive simulations, bridging Earth surface processes with remote sensing [2]. Further, the use of plant trait- and process-based stomatal models are also drawing more attention given the improved predictive skills compared to empirical representations [3–5] and more direct link to plant physiological status [6]. However, implementing these biophysical representations globally has been challenging due to the lack of a complete suite of global soil and plant traits and the difficulties to assess model simulations.

58 The capability of recently developed terrestrial biosphere models in simu-
59 lating plant physiological processes and canopy optical properties and bridging
60 them to remote sensing makes it possible to constrain Earth surface processes
61 using remote sensing data [7–10]. Notably, increasing volume and types of
62 global scale ecological datasets, products, and databases are largely being made
63 publicly available in recent years, such as multiple decades of satellite-based
64 maps of vegetation indices from Moderate Resolution Imaging Spectroradiome-
65 ter (MODIS) instruments. These global scale datasets may serve as initial or
66 boundary conditions or as benchmark standards for the ESMs [11, 12].

67 However, in parallel to the great promise of the ever-increasing volume
68 of data, harnessing these datasets is becoming a considerable challenge for
69 the research community. Typical day-to-day challenges for researchers emerge
70 because the datasets (1) are posted and stored at different websites; (2) have
71 different formats (e.g., GeoTIFF and NetCDF files); (3) may not be usable
72 directly out of the box (e.g., data is converted to Bytes for smaller size, and thus
73 re-scaling is required); (4) have different orientations (e.g., some from South-
74 ern to Northern Hemisphere, and vice versa); (5) have different latitude and
75 longitude setups (e.g., some exclude high latitude regions); (6) have different
76 projections (e.g., cylindrical and pseudo-cylindrical); and (7) may have differ-
77 ent and non-standard units. These limitations pose substantial barriers even
78 for experienced programmers, given the time required to find and reformat the
79 data. Particularly, it is very inconvenient for both beginners and experienced
80 researchers who may just require a tailored subset of these broadly available
81 scientific datasets. For instance, one may need to download the entire global
82 dataset to obtain information for a few sites or a small region.

83 The emergence of the Google Earth Engine (GEE) [13] has largely advanced
84 data sharing and reuse, and the cloud computing platform provided further

4 *GriddingMachine*

85 relaxes the need for extensive local storage and computation resource. How-
86 ever, GEE is not always the best solution for researchers given that users do
87 not have direct control over the datasets for easier offline analysis and debug-
88 ging. Most importantly, many public datasets circulate within the research
89 community but are not converted to GEE images. In comparison, the Applica-
90 tion for Extracting and Exploring Analysis Ready Samples (AppEEARS) and
91 European Centre for Medium-Range Weather Forecasts (ECMWF) offer users
92 simple and efficient ways to request raw data subsets (limited to their own
93 datasets; need to register, and download and manage the subsets manually).
94 Therefore, it remains a question how to more effectively share, manage, and
95 reuse these community-based datasets stored on local hard drives. One way to
96 circumvent this issue is to build a database for these community datasets and
97 label each with a unique easy-to-read tag. Using tags to automatically down-
98 load and manage the datasets would largely simplify the user end operations
99 and facilitate the data reuse.

100 Our aim is to create and maintain such a standardized collection of glob-
101 ally spanning datasets for biophysical representations in ESMs. Therefore, we
102 need to resolve the problems listed above and distribute standardized datasets
103 using tags. Our solution, *GriddingMachine*, is a one-stop shop for downloading,
104 storing, exploring, and jointly extracting multiple datasets we have collected
105 and reprocessed. *GriddingMachine* was initially developed in the Julia pro-
106 gramming language [14] to parameterize a new-generation ESM developed by
107 the Climate Modeling Alliance (CliMA), particularly the land component—
108 CliMA Land [9, 15, 16]. We have further generalized the model for the purpose
109 of processing and distributing datasets by adding more global scale datasets
110 other than the land parameters (such as remote sensing based products), and

111 providing more user interfaces through HTTP, Matlab, Octave, Python, and
 112 R in addition to Julia (Fig. 1).

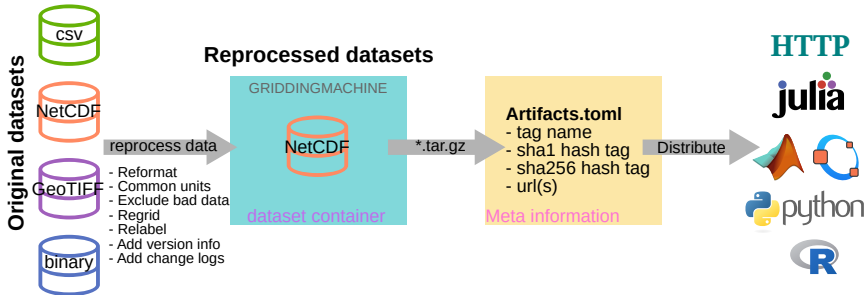


Fig. 1 Pathway used to assemble and distribute the GriddingMachine database. Each dataset is (a) reprocessed to meet our standards for distribution, (b) compressed as a tar.gz file along with an empty label file GRIDDINGMACHINE, and (c) stored on publicly available HTTP servers. Then the meta information for each dataset is stored in Artifacts.toml, which includes the tag name, sha1 hashtag, sha256 hashtag, and downloading URLs. Users are able to access the datasets manually via HTTP protocols or automatically through Julia, Matlab, Octave, Python, and R functions aided by Artifacts.toml.

113 Results

114 Each dataset in GriddingMachine database is labeled with a unique tag follow-
 115 ing a general naming pattern: LABEL_(EXTRALABEL_)IX_JT_(YEAR_)VK,
 116 where EXTRALABEL and YEAR are optional (Table 1). Here, LABEL is the
 117 major data identifier (e.g., GPP for gross primary productivity), EXTRALA-
 118 BEL is the secondary data identifier to distinguish datasets within a collection
 119 (for example, GPP from different models), IX indicates the spatial resolution
 120 is $1/I^\circ$, JT represents the temporal resolution (H for hour, D for day, M for
 121 month, and Y for year), and VK means the dataset version within our collec-
 122 tion (V1 for data from publication 1, and VN for data from publication N).
 123 For example, LAIMODIS_20X_1M_2008_V1 stands for leaf area index (LAI)
 124 from MODIS at $1/20^\circ$ spatial resolution and 1 month temporal resolution for
 125 the year 2008, and V1 indicates that the dataset is the first LAI collection;

6 *GriddingMachine*

126 VCMAX_2X_1Y_V2 stands for maximum carboxylation rate (V_{cmax}) at $1/2^\circ$
 127 spatial resolution and 1-year temporal resolution, and V2 indicates that the
 128 dataset is the second V_{cmax} collection.

129 Although we aim for automatic and convenient data access within diverse
 130 programming languages, we still provide traditional data access from HTTP
 131 servers. The compressed datasets can be manually downloaded from an open
 132 data archive hosted on CaltechDATA [38].

133 **Julia support**

134 We provide Julia language support using GriddingMachine.jl (v0.2, requires
 135 Julia 1.6+), which include three tested sub-modules: Collector, Indexer, and
 136 Requestor (Fig. 2). Collector is responsible for downloading and manag-
 137 ing datasets; Indexer is responsible for reading downloaded datasets; and
 138 Requestor is responsible for requesting a subset of data directly from our
 139 server without downloading the datasets. To install GriddingMachine.jl, one
 140 may simply type the following in Julia REPL (read-eval-print loop):

```
# Julia script
using Pkg;
Pkg.add("GriddingMachine");
```

141 Below, we introduce only the fundamental features of GriddingMachine.jl,
 142 and refer the readers to our online documentation for more details about
 143 GriddingMachine.jl at <https://clima.github.io/GriddingMachine.jl/stable/>.

144 **Download data**

145 We classify the datasets into different collections using the GriddedCollection
 146 structure, which includes the following fields

- 147 • LABEL: a string that appears in file names, such as LAI;

Table 1 Datasets within GriddingMachine collections.

Dataset type	LABEL	EXTRALABEL	IX	JT	YEAR	VK	Reference
Biomass	BIOMASS	ROOT	120X	1Y	-	V1	[17]
	BIOMASS	SHOOT	120X	1Y	-	V2	[18]
Canopy height	CH	-	20X	1Y	-	V1	[19]
	CH	-	2X	1Y	-	V2	[20]
Clumping index	CI	-	2X, 240X	1Y	-	V1	[21]
	CI	-	2X	1Y	-	V2	[22]
Elevation	ELEV	-	4X	1Y	-	V1	[23]
Gross primary productivity	GPP	MPLRS	2X	1M, 8D	2001–2019	V1	[24]
	GPP	VPM	5X, 12X	8D	2000–2019	V2	[25]
Leaf area index	LAI	MODIS	2X, 10X, 20X	1M, 8D	2000–2020	V1	[26]
Land mask	LM	-	4X	1Y	-	V1	ERA5
Leaf nitrogen content	LNC	-	2X	1Y	-	V1	[27]
	LNC	-	2X	1Y	-	V2	[20]
Leaf phosphorus content	LPC	-	2X	1Y	-	V1	[27]
Plant functional type	PFT	-	2X	1Y	-	V1	[28]
Surface area	SA	-	1X, 2X	1Y	-	V1	[28]
Soil color class	SC	-	2X	1Y	-	V1	[28]
Solar-induced chlorophyll fluorescence	SIF	TROPOMI_740	1X, 2X, 4X, 5X, 12X	1M, 8D	2018–2020	V1	[29]
	SIF	TROPOMI_740DC	1X, 2X, 4X, 5X, 12X	1M, 8D	2018–2020	V1	[29]
	SIF	TROPOMI_683	1X, 2X, 4X, 5X, 12X	1M, 8D	2018–2020	V2	[30]
	SIF	TROPOMI_683DC	1X, 2X, 4X, 5X, 12X	1M, 8D	2018–2020	V2	[30]
	SIF	OCO2_757	5X	1M	2014–2020	V3	[31]
SIF	OCO2_757DC	5X	1M	2014–2020	V3	[31]	
SIF	OCO2_771	5X	1M	2014–2020	V3	[31]	
SIF	OCO2_771DC	5X	1M	2014–2020	V3	[31]	

Dataset type	LABEL	EXTRALABEL	IX	JT	YEAR	VK	Reference
Solar-induced luminescence	SIL	-	20X	1Y	-	V1	[32]
Specific leaf area	SLA	-	2X	1Y	-	V1	[27]
	SLA	-	2X	1Y	-	V2	[20]
Soil hydraulics	SOIL	SWCR	12X, 120X	1Y	-	V1	[33]
	SOIL	SWCS	12X, 120X	1Y	-	V1	[33]
	SOIL	VGA	12X, 120X	1Y	-	V1	[33]
	SOIL	VGN	12X, 120X	1Y	-	V1	[33]
	SOIL	KSAT	100X	1Y	-	V2	[34]
Tree density	TD	-	2X, 120X	1Y	-	V1	[35]
Maximum carboxylation rate	VCMAX	-	2X	1Y	-	V1	[36]
	VCMAX	-	2X	1Y	-	V2	[37]
Wood density	WD	-	2X	1Y	-	V1	[20]

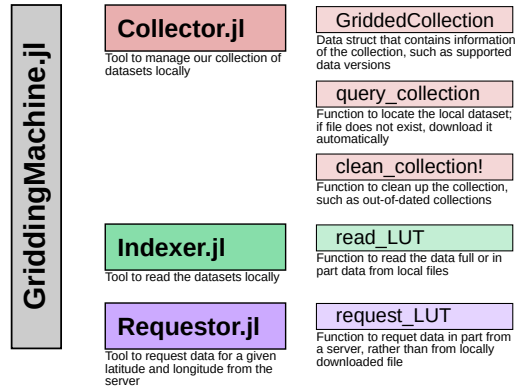


Fig. 2 Framework of GriddingMachine.jl package (v0.2). GriddingMachine.jl contains three sub-modules: Collector, Indexer, and Requestor. Collector downloads and manages the datasets; Indexer reads the downloaded datasets; and Requestor requests a subset of data directly from the server without downloading the datasets.

- 148 • **SUPPORTED_COMBOS**: an array of string for supported combinations, such as [MODIS_2X_8D_2018_V1, MODIS_2X_8D_2019_V1,
149 MODIS_2X_8D_2020_V1];
150
- 151 • **DEFAULT_COMBO**: a string of the default combination, such as
152 MODIS_2X_8D_2020_V1.

153 Currently, we have a total of 19 categories of dataset (column Dataset type
154 in Table 1), and 19 functions to construct these GriddedCollection structures
155 accordingly:

- 156 • biomass_collection
157 • canopy_height_collection
158 • clumping_index_collection
159 • elevation_collection
160 • gpp_collection
161 • lai_collection
162 • land_mask_collection
163 • leaf_nitrogen_collection

10 *GriddingMachine*

- 164 • leaf_phosphorus_collection
- 165 • pft_collection
- 166 • sif_collection
- 167 • sil_collection
- 168 • sla_collection
- 169 • soil_color_collection
- 170 • soil_hydraulics_collection
- 171 • surface_area_collection
- 172 • tree_density_collection
- 173 • vcmx_collection
- 174 • wood_density_collection

175 For example, one can define a LAI collection and check supported combinations
 176 using the following commands:

```
# Julia script
using GriddingMachine.Collector;
@show lai_collection();
```

177 Function `query_collection` returns the path to the dataset (Julia will down-
 178 load the file automatically to folder `~/.julia/artifacts` if the dataset does not
 179 exist):

```
# Julia script
using GriddingMachine.Collector;
data_path_1 = query_collection(lai_collection());
data_path_2 = query_collection(lai_collection(), "MODIS_2X_8D_2020_V1");
data_path_3 = query_collection("LAI_MODIS_2X_8D_2020_V1");
```

180 The first method returns the path to default combination; and the second and
 181 third methods return the path to target dataset.

182 Function `clean_collection!` cleans up the downloaded datasets:

```
# Julia script
using GriddingMachine.Collector;
clean_collections!();
```

```
clean_collections!("old");
clean_collections!("all");
clean_collections!(["LAI_MODIS_2X_8D_2019_V1", "LAI_MODIS_2X_8D_2020_V1"]);
clean_collections!(lai_collection());
```

183 The first and second methods clean up all outdated datasets that were not in
184 Artifacts.toml; the third method cleans up all GriddingMachine datasets; the
185 fourth method cleans up all selected datasets; and the last method cleans up
186 all datasets in a collection.

187 Read downloaded data

188 With the dataset path from query_collection, we are able to load the data using
189 function read_LUT (meaning read look-up-table) in Indexer. The supported
190 methods are

```
# Julia script
using GriddingMachine.Collector;
using GriddingMachine.Indexer;
data_path = query_collection(gpp_collection());
dat_1 = read_LUT(data_path);
dat_2 = read_LUT(data_path, 8);
dat_3 = read_LUT(data_path, 30.1, 116.1; interpolation = false);
dat_4 = read_LUT(data_path, 30.1, 116.1, 8; interpolation = false);
```

191 The first method loads the entire dataset (3D array in the example above);
192 the second method loads the 8th layer of the dataset (only applicable for 3D
193 datasets); the third method returns all the data at a given latitude (30.1°)
194 and longitude (116.1°); and the last method returns only data on the given
195 latitude (30.1°), longitude (116.1°), and cycle index (8). For the third and last
196 methods, option interpolation is false by default and the function returns data
197 that falls into the grid (latitude from 30 to 30.5° and longitude from 116 to
198 116.5° in this example). If option interpolation is true, we use the bi-linear
199 interpolation method, and the function returns linearly interpolated results
200 from its four nearest neighbours.

201 **Request partial data**

202 Note that functions `query_collection` and `read_LUT` are meant for downloading
 203 (if not exist) and reading the local dataset, which may hamper the data reusing.
 204 For example, if one only needs the data for a few sites, downloading gigabytes
 205 of data for these few sites would (1) be time consuming, (2) waste local storage,
 206 and (3) increase unnecessary load for data servers. Therefore, we provide a way
 207 to request data for specific sites in Requestor through function `request_LUT`:

```
# Julia script
using GriddingMachine.Requestor;
art_name = "LAI_MODIS_2K_8D_2020_V1";
lais,stds = request_LUT(art_name, 30.1, 116.1; interpolation = true);
lai8,std8 = request_LUT(art_name, 30.1, 116.1, 8; interpolation = true);
```

208 The first method takes the tag name (full name), latitude (30.1°), and lon-
 209 gitude (116.1°) as input, and returns the time series of LAI and its error;
 210 the second method take an extra cycle index as input, and returns only the
 211 data and error for the target cycle (day 57 to 64 in this example). Similar
 212 to function `read_LUT`, one may choose to interpolate the data by setting the
 213 interpolation option to false or true. What `request_LUT` does are (i) user end
 214 passes the input information to the server, (2) server end reads the data using
 215 `query_collection` and `read_LUT`, (3) server end returns a structured JSON file
 216 back to the user end, and (4) user end translates the structured JSON back
 217 to data (numbers or arrays).

218 **Other language supports**

219 Besides Julia, we also provide simple user interfaces for Matlab, Octave,
 220 Python and R to aid dataset distribution and reuse. We provide three func-
 221 tions for Matlab, Octave, Python and R, and they are (1) `update_GM` that
 222 downloads the latest `Artifacts.toml` from Github, (2) `query_collection` that

223 downloads and returns the path of the dataset (same as that in Julia), and
224 (3) `request_LUT` that requests partial data from the server (same as that
225 in Julia). Different from Julia, `update_GM` and `query_collection` download
226 the data to `~/GMCollections` rather than `~/julia/artifacts`. `Artifacts.toml`
227 is stored at `~/GMCollections/Artifacts.toml`; the compressed `tar.gz` files
228 are stored in `~/GMCollections/archives`; and the datasets are extracted to
229 `~/GMCollections/artifacts`. With the dataset path, users can choose the
230 packages they prefer to read the dataset. In comparison, `GriddingMachine.jl`
231 updates `Artifacts.toml` through package releasing, and one needs to update
232 `GriddingMachine.jl` to use the latest datasets.

233 Matlab

234 We provide Matlab language support via Matlab Toolbox (source code avail-
235 able at <https://github.com/Yujie-W/octave-griddingmachine>). Matlab users
236 may install and use the toolbox using

```
% Matlab script
% Install the toolbox
url = 'https://github.com/Yujie-W/octave-griddingmachine/raw/main/GriddingMachine.mltbx';
urlwrite(url, 'GriddingMachine.mltbx');
matlab.addons.toolbox.installToolbox('GriddingMachine.mltbx');
delete('GriddingMachine.mltbx');
% Use the toolbox
update_GM();
art_name = 'VCMAX_2X_1Y_V1';
file_path = query_collection(art_name);
[vcmax,error] = request_LUT(art_name, 35.1, 115.2);
[vcmax,error] = request_LUT(art_name, 35.1, 115.2, 'interpolation', true);
```

237 Octave

238 Octave language support is also provided via [https://github.com/Yujie-W/](https://github.com/Yujie-W/octave-griddingmachine)
 239 [octave-griddingmachine](https://github.com/Yujie-W/octave-griddingmachine), which is Matlab and Octave compatible. Octave users
 240 may install and use the package using

```
% Octave script
% Install the package
pkg install "https://github.com/gnu-octave/pkg-json/archive/v1.5.0.tar.gz";
pkg install "https://github.com/Yujie-W/octave-griddingmachine/archive/v0.1.1.tar.gz";
% Use the package
pkg load griddingmachine;
update_GM();
art_name = 'VCMAX_2X_1Y_V1';
file_path = query_collection(art_name);
[vcmax,error] = request_LUT(art_name, 35.1, 115.2);
[vcmax,error] = request_LUT(art_name, 35.1, 115.2, 'interpolation', true);
```

241 Python

242 We provide Python language support via a registered Python pack-
 243 age through PyPI (source code available at [https://github.com/Yujie-W/](https://github.com/Yujie-W/python-griddingmachine)
 244 [python-griddingmachine](https://github.com/Yujie-W/python-griddingmachine)), and Python users can install the package using pip:

```
# Shell command
pip install python-griddingmachine
```

245 To query the path to downloaded dataset or subset data directly from the
 246 server, one may use

```
# Python script
from griddingmachine import update_GM, query_collection, request_LUT;
update_GM();
art_name = "VCMAX_2X_1Y_V1";
file_path = query_collection(art_name);
vcmax,error = request_LUT(art_name, 35.1, 115.2);
vcmax,error = request_LUT(art_name, 35.1, 115.2, interpolation = True);
```

247 **R**

248 We provide R language support via a package hosted on Github, and R users
249 can install and use the package using

```
# R script
# Install the package
library(devtools);
install_github("Yujie-W/r-griddingmachine");
# Use the Package
library("griddingmachine");
update_GM();
art_name <- "VCMAX_2X_1Y_V1";
file_path <- query_collection(art_name);
results <- request_LUT(art_name, 35.1, 115.2);
results <- request_LUT(art_name, 35.1, 115.2, interpolation = TRUE);
```

250 As R does not allow for returning multiple variables, data and error of
251 the request are stored as fields of a list, and one can read them out using
252 `results$data` and `results$std`.

253 **Discussion**

254 We recommend GriddingMachine users to use provided interfaces through
255 Julia, Matlab, Octave, Python, and R to access the data. In particular, if the
256 research is designed to run at global or large regional scales, we recommend
257 using `query_collection` to download the datasets, which allows the users to use
258 the latest datasets and perform offline analysis more efficiently. If the research
259 is meant to run at smaller scales such as a few sites, we recommend using
260 `request_LUT` for convenience. For example, Julia code below shows the simple
261 steps of comparing datasets at the site level. We compared two gross primary
262 productivity products vs. one solar-induced chlorophyll fluorescence for a flux
263 tower site in North America (US-NR1; latitude and longitude are 40.0329° and
264 -105.5464°, respectively; data from the year 2019), and the result is shown in

265 Fig. 3. In the example presented, requesting the data directly from the server
 266 avoids downloading approximately 600 MB reprocessed data (original datasets
 267 from three different sources are 2.6 GB):

```
# Julia script
using GriddingMachine.Requestor;
# load FLUXCOM and VPM GPPs and TROPOMI DC SIF
gpp_mpi,_ = request_LUT("GPP_MPI_RS_2X_8D_2019_V1", 40.0329, -105.5464);
gpp_vpm,_ = request_LUT("GPP_VPM_12X_8D_2019_V2", 40.0329, -105.5464);
sif_dc,_ = request_LUT("SIF_TROPOMI_740DC_12X_8D_2019_V1", 40.0329, -105.5464);
# plot the comparison
```

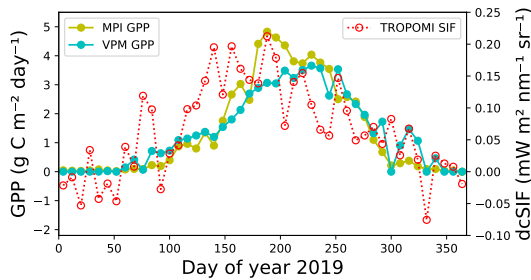


Fig. 3 Example of simple dataset requests using GriddingMachine. The requested data are gross primary productivity (GPP) from [24] (MPI GPP, olive symbols) and [25] (VPM GPP, cyan symbols) and day length corrected solar-induced chlorophyll fluorescence (dcSIF) from [29] (TROPOMI SIF, red symbols).

268 We note here that the data management of GriddingMachine is subject to
 269 future changes and users can follow on the GitHub page at <https://github.com/CliMA/GriddingMachine.jl>. Users can find tutorials and examples about
 270 [https://github.com/CliMA/GriddingMachine.jl#](https://github.com/CliMA/GriddingMachine.jl#data-contribution)
 271 [data-contribution](https://github.com/CliMA/GriddingMachine.jl#data-contribution).
 272 [data-contribution](https://github.com/CliMA/GriddingMachine.jl#data-contribution).

273 In conclusion, we present GriddingMachine, a generalized database and
 274 software to distribute standardized globally spanning datasets in a user-
 275 friendly way. We provide examples that can serve as templates to load different
 276 types of datasets we have collected as well an easy way to request partial
 277 data from our server in five programming languages: Julia, Matlab, Octave,

278 Python, and R. Given the aims of GriddingMachine, we welcome the contri-
279 bution of globally gridded data to our collection through [https://github.com/
280 CliMA/GriddingMachine.jl/issues/62](https://github.com/CliMA/GriddingMachine.jl/issues/62), and believe that GriddingMachine will
281 be a useful tool for Earth System Modeling and ecology communities.

282 **Methods**

283 To maximally simplify and facilitate the data reuse, we first reprocessed the
284 collected data to standardized NetCDF formatted datasets (Fig. 1). We used
285 NetCDF format because the datasets can be 3D arrays, and the file can be
286 read using non-proprietary software. Then, we compressed each dataset as a
287 tar.gz file, and stored the compressed datasets on public HTTP servers. Last,
288 we stored the meta information in an Artifacts.toml file, and provided user
289 friendly functions to automatically retrieve the datasets for multiple program-
290 ming languages such as Julia, Matlab, Octave, Python, and R. Note here that
291 Artifacts.toml file is often used in Julia to host meta information of data con-
292 tainers named artifacts. The containers may contain any other kind of data
293 that would be convenient to place within an immutable, life-cycled data store.
294 These containers, (called artifacts) can be created locally, hosted anywhere,
295 and automatically downloaded and unpacked upon request. In GriddingMa-
296 chine.jl, we used this artifact feature to redistribute the reprocessed datasets,
297 which are not supposed to change with time. Alternatively, users may also
298 download the datasets manually from the HTTP servers.

299 **Data reprocessing**

300 The raw datasets we collected were in various formats and not standard-
301 ized (for example, some datasets were scaled but the scaling factors were not

302 noted within the dataset itself). Thus, we reprocessed each dataset before
303 distributing it to the public, and the standards are

- 304 • The dataset is stored in a NetCDF file
- 305 • The dataset is either a 2D or 3D array
- 306 • The dataset is cylindrically projected (WGS84 projection)
- 307 • The first dimension of the dataset is longitude
- 308 • The second dimension of the dataset is latitude
- 309 • The third dimension (if available) is the cycle index, e.g., time
- 310 • The longitude is oriented from west to east hemisphere (-180° to 180°)
- 311 • The latitude is oriented from south to north hemisphere (-90° to 90°)
- 312 • The dataset covers the entire globe (missing data allowed)
- 313 • Missing data is labeled as NaN (not a number) rather than an unrealistic
314 fill value
- 315 • The dataset is not scaled (linearly, exponentially, or logarithmically)
- 316 • The dataset has common units, such as $\mu\text{mol m}^{-2} \text{ s}^{-1}$ for maximum
317 carboxylation rate
- 318 • The spatial resolution is uniform longitudinally and latitudinally, e.g., both
319 at $1/2^\circ$
- 320 • The spatial resolution is an integer division of 1° , such as $1/2^\circ$, $1/12^\circ$, $1/240^\circ$
- 321 • Each grid cell represents the average value of everything inside the grid cell
322 area (as opposing to a single point in the middle of the cell)
- 323 • The label for the data is “data” (for conveniently loading the data)
- 324 • The label for the error is “std” (for conveniently loading the error)
- 325 • The dataset must contain one data array and one error array besides the
326 dimensions
- 327 • The dataset contains citation information in the attributes

- The dataset contains a log summarizing changes if different from original source

Each reprocessed NetCDF file contains four (2D dataset) or five (3D dataset) fields (see Table 2 for the details of the fields and attributes of the reprocessed datasets). We note that there could be infinite number of fields in a NetCDF file; but for the ease of automatically reading the datasets, we only included data and std fields in one reprocessed datasets. For example, the original TROPOMI solar-induced chlorophyll fluorescence (SIF) datasets contain both uncorrected SIF and day length corrected SIF [29, 30], and we partitioned the file to separate files to allow for more automated data requests.

Data packaging

After data reprocessing, we compressed each dataset along with an empty GRIDDINGMACHINE file as a tar.gz file (contains 2 files). The empty GRIDDINGMACHINE file labels the NetCDF file as a GriddingMachine dataset (used to clean up outdated datasets). For example, any folder with the empty file GRIDDINGMACHINE will be treated as a GriddingMachine artifact, and the hashtag for this artifact is the same as the folder name. If this hashtag does not exist in Artifacts.toml, then this GriddingMachine artifact will be marked as outdated. Users can delete the outdated GriddingMachine using `clean_collections!`. The NetCDF dataset file name and the tag name are the same and follow a general naming pattern: LABEL_(EXTRALABEL_)IX_JT_(YEAR_)VK. See Table 1 for current collections and <https://github.com/CliMA/GriddingMachine.jl/issues/62> for the growing collection list and detailed change logs for each dataset.

Meta information of all the datasets is stored in Artifacts.toml, and each item of the toml file includes the following: tag name, SHA1 hash value,

Table 2 Fields and attributes of the reprocessed NetCDF datasets.

Field	Dimension	Description	Attributes
lon	1D array	Longitude in the center of a grid	unit: ° description: Latitude
lat	1D array	Latitude in the center of a grid	unit: ° description: Longitude
ind	1D array	Cycle index (only available in 3D datasets)	unit: - description: Cycle index
data	2D/3D array	Data in the center of a grid	longname: long name of the data unit: unit of the data about: general information authors: authors of the dataset source publication year: year of the data source publication title: title of the data source publication journal: journal of the data source publication doi: DOI tag of the data source publication changeN: Change log of the Nth change we made same as "data" field
std	2D/3D array	Error of data in the center of a grid	

354 SHA256 hash value, and downloading URLs. Here SHA stands for Secure Hash
355 Algorithm, and it was used to compute a unique hashtag for each dataset,
356 thus aiding data verification when users use the data. Through Artifacts.toml,
357 we provided functions to automatically download and load the datasets in
358 multiple programming languages.

359 Data availability

360 The compressed datasets can be downloaded from CaltechDATA [38].

361 Code availability

362 The code can be found at <https://github.com/CliMA/GriddingMachine.jl>
363 under the Apache 2.0 License. The exact version of the package used to pro-
364 duce the results presented in this paper is also archived on CaltechDATA along
365 with the datasets [38].

366 **Acknowledgments.** We gratefully acknowledge the generous support of
367 Eric and Wendy Schmidt (by recommendation of the Schmidt Futures) and
368 the Heising-Simons Foundation. This research has been supported by the
369 National Aeronautics and Space Administration (NASA) Earth Sciences Divi-
370 sion grant NNX15AH95G and Carbon Cycle Science grant 80NSSC21K1712
371 awarded to Christian Frankenberg. Marcos Longo was supported by the NASA
372 Postdoctoral Program, administered by Universities Space Research Associ-
373 ation under contract with NASA. We acknowledge NASA for the publicly
374 available datasets. The distributed datasets include modified Copernicus Cli-
375 mate Change Service Information [2020]. Neither the European Commission
376 nor the European Centre for Medium-Range Weather Forecasts (ECMWF)
377 are responsible for any use that may be made of the Copernicus information

378 or data in this publication. We thank the data owners for generously shar-
379 ing the datasets. Part of this research was carried out at the Jet Propulsion
380 Laboratory, California Institute of Technology, under a contract with NASA.
381 California Institute of Technology. Government sponsorship acknowledged.

382 **Author contributions.** YW designed the code structure, wrote the code,
383 and led the writing. All authors identified suitable datasets and contributed
384 to the writing.

385 **Competing interests.** The authors declare no competing interests

386 References

- 387 [1] Bonan, G. B., Patton, E. G., Finnigan, J. J., Baldocchi, D. D. & Harman,
388 I. N. Moving beyond the incorrect but useful paradigm: reevaluating big-
389 leaf and multilayer plant canopies to model biosphere-atmosphere fluxes—a
390 review. *Agricultural and Forest Meteorology* **306**, 108435, [https://doi.
391 org/10.1016/j.agrformet.2021.108435](https://doi.org/10.1016/j.agrformet.2021.108435) (2021) .
- 392 [2] Gu, L., Han, J., Wood, J. D., Chang, C. Y.-Y. & Sun, Y. Sun-induced chl
393 fluorescence and its importance for biophysical modeling of photosynthesis
394 based on light reactions. *New Phytologist* **223** (3), 1179–1191 (2019) .
- 395 [3] Mencuccini, M., Manzoni, S. & Christoffersen, B. Modelling water fluxes
396 in plants: From tissues to biosphere. *New Phytologist* **222** (3), 1207–1222
397 (2019) .
- 398 [4] Wang, Y., Sperry, J. S., Anderegg, W. R. L., Venturas, M. D. & Trugman,
399 A. T. A theoretical and empirical assessment of stomatal optimization
400 modeling. *New Phytologist* **227**, 311–325 (2020) .

- 401 [5] Wang, Y. *et al.* Optimization theory explains nighttime stomatal
402 responses. *New Phytologist* **230** (4), 1550–1561 (2021) .
- 403 [6] Sperry, J. S. *et al.* The impact of rising CO₂ and acclimation on the
404 response of US forests to global warming. *Proceedings of the National
405 Academy of Sciences* **116** (51), 25734–25744 (2019) .
- 406 [7] Yang, P., Verhoef, W. & van der Tol, C. The mscope model: A simple
407 adaptation to the scope model to describe reflectance, fluorescence and
408 photosynthesis of vertically heterogeneous canopies. *Remote sensing of
409 environment* **201**, 1–11 (2017) .
- 410 [8] Braghieri, R. K. *et al.* Accounting for canopy structure improves
411 hyperspectral radiative transfer and sun-induced chlorophyll fluorescence
412 representations in a new generation earth system model. *Remote Sensing
413 of Environment* **261**, 112497, <https://doi.org/10.1016/j.rse.2021.112497>
414 (2021) .
- 415 [9] Wang, Y. *et al.* Testing stomatal models at the stand level in deciduous
416 angiosperm and evergreen gymnosperm forests using clima land (v0.1).
417 *Geoscientific Model Development* **14** (11), 6741–6763 (2021) .
- 418 [10] Yang, P., Prikaziuk, E., Verhoef, W. & van der Tol, C. Scope 2.0:
419 a model to simulate vegetated land surface fluxes and satellite signals.
420 *Geoscientific Model Development* **14** (7), 4697–4712 (2021) .
- 421 [11] Stavros, E. N. *et al.* ISS observations offer insights into plant func-
422 tion. *Nature Ecology & Evolution* **1** (7), 0194, [https://doi.org/10.1038/
423 s41559-017-0194](https://doi.org/10.1038/s41559-017-0194) (2017) .

- 424 [12] Schimel, D., Schneider, F. D., Carbon, J. & Participants, E. Flux towers
425 in the sky: global ecology from space. *New Phytologist* **224** (2), 570–584
426 (2019) .
- 427 [13] Gorelick, N. *et al.* Google earth engine: Planetary-scale geospatial analysis
428 for everyone. *Remote sensing of Environment* **202**, 18–27 (2017) .
- 429 [14] Bezanson, J., Edelman, A., Karpinski, S. & Shah, V. B. Julia: A fresh
430 approach to numerical computing. *SIAM Review* **59** (1), 65–98, <https://doi.org/10.1137/141000671>
431 (2017) .
- 432 [15] Wang, Y. & Frankenberg, C. On the impact of canopy model complexity
433 on simulated carbon, water, and solar-induced chlorophyll fluorescence
434 fluxes. *Biogeosciences* **19** (1), 29–45 (2022) .
- 435 [16] Wang, Y. *et al.* Modeling global carbon and water fluxes and hyperspec-
436 tral canopy radiative transfer simultaneously using a next generation land
437 surface model—clima land. *Earth and Space Science Open Archive* **38**,
438 <https://doi.org/10.1002/essoar.10509956.1> (2022) .
- 439 [17] Huang, Y. *et al.* A global map of root biomass across the world’s forests.
440 *Earth System Science Data* **13** (9), 4263–4274 (2021) .
- 441 [18] Santoro, M. *et al.* The global forest above-ground biomass pool for
442 2010 estimated from high-resolution satellite observations. *Earth System
443 Science Data* **13** (8), 3927–3950 (2021) .
- 444 [19] Simard, M., Pinto, N., Fisher, J. B. & Baccini, A. Mapping for-
445 est canopy height globally with spaceborne lidar. *Journal of Geo-
446 physical Research: Biogeosciences* **116**, G4021, [https://doi.org/10.1029/
447 2011JG001708](https://doi.org/10.1029/2011JG001708) (2011) .

- 448 [20] Boonman, C. C. *et al.* Assessing the reliability of predicted plant trait
449 distributions at the global scale. *Global Ecology and Biogeography* **29** (6),
450 1034–1051 (2020) .
- 451 [21] He, L., Chen, J. M., Pisek, J., Schaaf, C. B. & Strahler, A. H. Global
452 clumping index map derived from the modis brdf product. *Remote
453 Sensing of Environment* **119**, 118–130 (2012) .
- 454 [22] Braghieri, R. K., Quaife, T., Black, E., He, L. & Chen, J. Underestimation
455 of global photosynthesis in earth system models due to representation of
456 vegetation structure. *Global Biogeochemical Cycles* **33** (11), 1358–1369
457 (2019) .
- 458 [23] Yamazaki, D. *et al.* A high-accuracy map of global terrain elevations.
459 *Geophysical Research Letters* **44** (11), 5844–5853 (2017) .
- 460 [24] Tramontana, G. *et al.* Predicting carbon dioxide and energy fluxes across
461 global fluxnet sites with regression algorithms. *Biogeosciences* **13** (14),
462 4291–4313 (2016) .
- 463 [25] Zhang, Y. *et al.* A global moderate resolution dataset of gross primary
464 production of vegetation for 2000–2016. *Scientific data* **4**, 170165, <https://doi.org/10.1038/sdata.2017.165>
465 [//doi.org/10.1038/sdata.2017.165](https://doi.org/10.1038/sdata.2017.165) (2017) .
- 466 [26] Yuan, H., Dai, Y., Xiao, Z., Ji, D. & Shangguan, W. Reprocessing the
467 modis leaf area index products for land surface and climate modelling.
468 *Remote Sensing of Environment* **115** (5), 1171–1187 (2011) .
- 469 [27] Butler, E. E. *et al.* Mapping local and global variability in plant trait
470 distributions. *Proceedings of the National Academy of Sciences* **114** (51),
471 E10937–E10946 (2017) .

- 472 [28] Lawrence, P. J. & Chase, T. N. Representing a new MODIS consistent
473 land surface in the community land model (CLM 3.0). *Journal of Geo-*
474 *physical Research: Biogeosciences* **112**, G01023, [https://doi.org/10.1029/](https://doi.org/10.1029/2006JG000168)
475 [2006JG000168](https://doi.org/10.1029/2006JG000168) (2007) .
- 476 [29] Köhler, P. *et al.* Global retrievals of solar-induced chlorophyll fluorescence
477 with TROPOMI: First results and intersensor comparison to OCO-2.
478 *Geophysical Research Letters* **45** (19), 10,456–10,463 (2018) .
- 479 [30] Köhler, P. *et al.* Global retrievals of solar-induced chlorophyll fluores-
480 cence at red wavelengths with TROPOMI. *Geophysical Research Letters*
481 **47** (15), e2020GL087541, <https://doi.org/10.1029/2020GL087541> (2020)
482 .
- 483 [31] Sun, Y. *et al.* OCO-2 advances photosynthesis observation from space via
484 solar-induced chlorophyll fluorescence. *Science* **358** (6360), eaam5747,
485 <https://doi.org/10.1126/science.aam5747> (2017) .
- 486 [32] Köhler, P. *et al.* Mineral luminescence observed from space. *Geophys-*
487 *ical Research Letters* **48** (19), e2021GL095227, [https://doi.org/10.1029/](https://doi.org/10.1029/2021GL095227)
488 [2021GL095227](https://doi.org/10.1029/2021GL095227) (2021) .
- 489 [33] Dai, Y. *et al.* A global high-resolution data set of soil hydraulic and
490 thermal properties for land surface modeling. *Journal of Advances in*
491 *Modeling Earth Systems* **11** (9), 2996–3023 (2019) .
- 492 [34] Gupta, S., Lehmann, P., Bonetti, S., Papritz, A. & Or, D. Global
493 prediction of soil saturated hydraulic conductivity using random forest
494 in a covariate-based geotransfer function (cogtf) framework. *Journal*
495 *of Advances in Modeling Earth Systems* **13** (4), e2020MS002242, <https://doi.org/10.1029/2020MS002242>, <https://doi.org/10.1029/2020MS002242>

- 496 [//doi.org/10.1029/2020MS002242](https://doi.org/10.1029/2020MS002242) (2021) .
- 497 [35] Crowther, T. W. *et al.* Mapping tree density at a global scale. *Nature*
498 **525** (7568), 201–205 (2015) .
- 499 [36] Smith, N. G. *et al.* Global photosynthetic capacity is optimized to the
500 environment. *Ecology Letters* **22** (3), 506–517 (2019) .
- 501 [37] Luo, X. *et al.* Global variation in the fraction of leaf nitrogen allocated
502 to photosynthesis. *Nature Communications* **12**, 4866, [https://doi.org/10.](https://doi.org/10.1038/s41467-021-25163-9)
503 [1038/s41467-021-25163-9](https://doi.org/10.1038/s41467-021-25163-9) (2021) .
- 504 [38] Wang, Y. Artifacts of griddingmachine.jl (v0.2) for land modeling.
505 *CaltechDATA* <https://doi.org/10.22002/D1.2129> (2021) .