

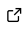
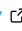

GridFlow: A modular high-performance toolkit for downloading and processing large-scale climate and geospatial datasets

Bhuwan P. Shah ¹ and Ryan P. McGehee ¹

¹ Iowa State University, Ames, Iowa, USA

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review](#) 
- [Repository](#) 
- [Archive](#) 

Editor: 

Submitted: 29 January 2026

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](https://creativecommons.org/licenses/by/4.0/)).

Summary

Climate and geospatial research increasingly relies on large observational and model-derived datasets such as CMIP6, CMIP5, ERA5 reanalysis, PRISM climate records, and global digital elevation models (DEMs). While these resources are widely available through web portals and distributed archives, acquiring and preparing them for analysis remains a major bottleneck. Researchers frequently spend considerable time navigating search interfaces, handling authentication requirements, writing brittle download scripts, and performing repetitive post-processing steps such as cropping, clipping, unit conversion, and temporal aggregation.

GridFlow is an open-source Python based toolkit that streamlines the complete workflow of climate and geospatial data preparation. It provides both a command-line interface (CLI) and a graphical user interface (GUI) to download major climate products and to process NetCDF datasets into analysis-ready subsets. GridFlow emphasizes modular design, parallel execution, and usability to support a broad user community, including researchers, students, and practitioners who need reliable access to large datasets without extensive custom scripting.

GridFlow is designed to be accessible to a wide range of users through both a GUI (Figure 1) and a command-line interface (Figure 2), enabling users to reproducibly download and process climate and geospatial datasets without the need for custom scripts.



Figure 1: GridFlow GUI showing the unified workflow for dataset acquisition and preprocessing.

```
(gridflow_env) D:\Research\Github\GridFlow>gridflow -h

=====
GridFlow
=====
Welcome to GridFlow v1.0! Copyright (c) 2025 Bhuwan Shah
GridFlow is a modular Python toolkit for downloading and processing key climate and geospatial
datasets. It provides both a powerful command-line interface and a user-friendly GUI to fetch
data from PRISM, DEM, CMIP5, CMIP6, and ERA5, and perform post-processing tasks like cropping,
clipping, and temporal aggregation.
=====

usage: gridflow [-h] [-v]

Run 'gridflow <command> -h' for detailed help on each command.

options:
  -h, --help      show this help message and exit
  -v, --version    show program's version number and exit

Downloading Tools:
  prism          Download PRISM climate data
  cmip5           Download CMIP5 climate model data
  cmip6           Download CMIP6 climate model data
  era5           Download ERA5-Land climate data
  dem            Download Digital Elevation Models via OpenTopography

Processing Tools:
  crop           Crop NetCDF files to a spatial bounding box
  clip           Clip NetCDF files using a shapefile
  convert         Convert units in NetCDF files
  aggregate      Temporally aggregate NetCDF files
  catalog        Generate a catalog from NetCDF files

Examples:
  gridflow cmip6 --demo      # Download sample CMIP6 data
  gridflow cmip5 --demo      # Download sample CMIP5 data
  gridflow prism --demo      # Download sample PRISM data
  gridflow crop --demo       # Crop NetCDF files to sample bounds
```

Figure 2: GridFlow CLI interface and built-in command help (gridflow -h).

23 Statement of need

24 Despite the continued growth of publicly available climate archives, data acquisition and
 25 processing remain disproportionately time-consuming compared to downstream analysis. For
 26 example, Earth System Grid Federation (ESGF) portals supporting CMIP5/CMIP6 provide
 27 powerful search tools, but workflows often involve repeated manual filtering, pagination, and
 28 batch downloads. Similarly, reanalysis data systems can impose account setup, API keys, or
 29 queue-based access patterns that complicate reproducible retrieval.

30 In many applied workflows (e.g., hydrology, water resources, agriculture, and land-surface
 31 modeling), users require climate variables for specific regions, watersheds, or administrative
 32 boundaries and often need derived temporal summaries (e.g., monthly means, seasonal sums).
 33 These tasks are commonly addressed through ad-hoc scripts using general-purpose scientific
 34 Python tools such as Xarray (Hoyer & Hamman, 2017) and GeoPandas (Jordahl & others,
 35 2020). However, implementing robust pipelines across multiple datasets and formats can create
 36 barriers to entry for new users and reduce reproducibility across projects.

37 GridFlow addresses this gap by offering a single, consistent interface for acquiring and preparing

multiple widely used climate and geospatial datasets. The toolkit couples high-level download modules with post-processing utilities to produce analysis-ready NetCDF outputs and metadata summaries, enabling rapid adoption in research and educational contexts.

Functionality and design

GridFlow is designed as a modular toolkit with two primary capability groups: (1) data downloading modules and (2) processing modules.

Downloading modules

GridFlow provides dedicated download modules for:

- **CMIP6 climate model data** via ESGF search and retrieval
- **CMIP5 climate model data** via ESGF search and retrieval
- **ERA5 reanalysis data** accessed directly from cloud-hosted sources
- **PRISM historical climate datasets** for the contiguous United States
- **Digital Elevation Models (DEM)** including global (Copernicus) and US-focused DEM products

Users interact with these modules using consistent CLI patterns, while GridFlow manages file organization, logging, metadata generation, and parallel retrieval. Download operations can be configured via command-line options or JSON configuration files to support reproducible workflows.

Processing modules

GridFlow also includes post-processing utilities for common climate-data preparation tasks:

- **Spatial subsetting**
 - Cropping NetCDF datasets to a latitude/longitude bounding box
 - Clipping datasets using irregular polygons defined by shapefiles
- **Temporal aggregation**
 - Aggregation from daily to monthly, seasonal, or annual frequency using multiple methods
- **Unit conversion**
 - Converting common climate variable units (e.g., Kelvin to Celsius)
- **Catalog generation**
 - Scanning NetCDF directories and producing a JSON catalog of dataset metadata for inventorying large multi-model libraries

Figure 3 demonstrates a typical workflow in which users download global or continental-scale NetCDF datasets and then spatially subset the files to a target watershed or region using cropping or clipping.

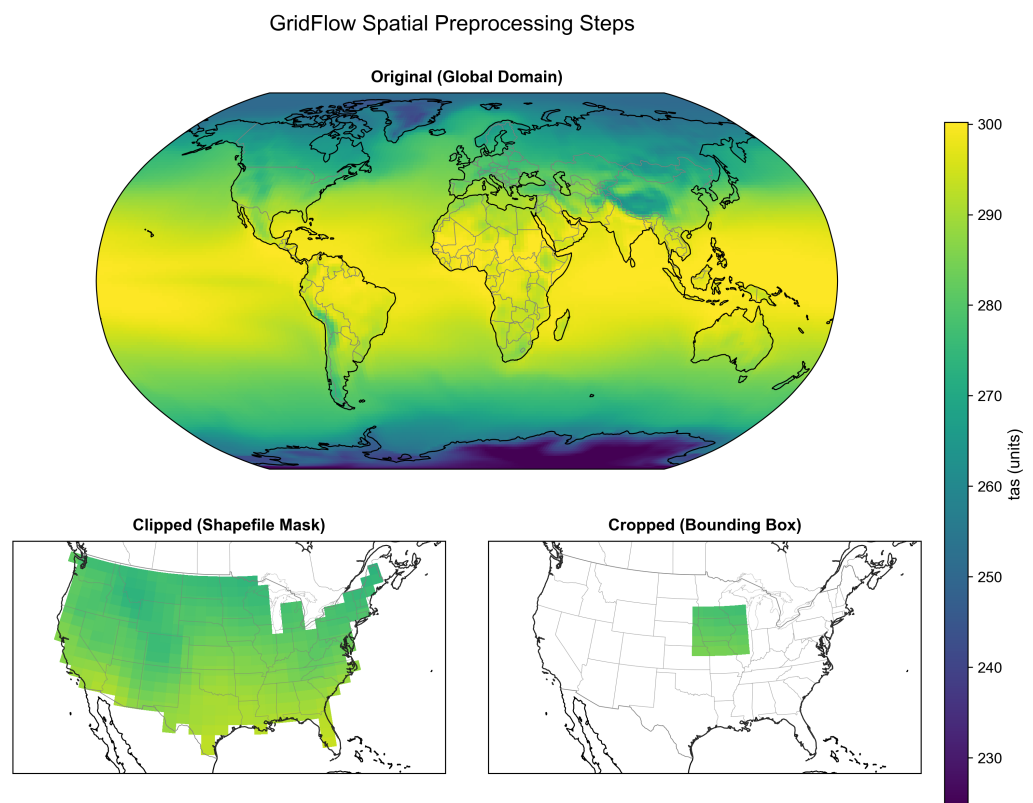


Figure 3: Example of spatial subsetting using GridFlow: clipping to CONUS polygon and cropping to a bounding box.

These operations are frequently required in climate impact studies, watershed-based modeling, and regional climate analyses, and GridFlow provides standardized implementations to reduce boilerplate and ensure consistent outputs.

Related work

GridFlow addresses a recurring challenge in climate and geospatial workflows: while many high-value datasets are publicly available, the end-to-end process of discovering, downloading, organizing, and preparing them for analysis remains fragmented across portals, archive-specific clients, and custom scripts. Users often combine multiple utilities for data acquisition (e.g., web portals or API clients) with separate tools for preprocessing (e.g., cropping, clipping, aggregation), which can reduce reproducibility and increase maintenance effort.

GridFlow complements downstream climate analytics libraries such as xclim (Bourgault et al., 2023) and xCDAT (Vo et al., 2024), which focus on derived indicators and analysis once data are already available locally or in a cloud-native format. In contrast, GridFlow focuses on the upstream bottleneck: scalable acquisition and standardized preprocessing of large gridded datasets. By providing consistent commands for both downloading and processing, GridFlow helps users move from raw archives to analysis-ready outputs without assembling a bespoke pipeline for each data source.

Several ecosystem tools support accessing climate archives through cloud-native workflows and metadata catalogs, including Pangeo-based stacks and dataset indexing approaches. For example, intake-esm (?) enables discovery and loading of CMIP-style collections through structured catalogs, and xarray-based tooling enables flexible analysis once datasets are opened.

93 These approaches are powerful, but they typically assume that users already operate within
94 specific cloud or catalog ecosystems and do not aim to provide a unified, user-facing pipeline
95 for bulk downloading and local preparation across heterogeneous sources.

96 GridFlow also relates to interactive GIS software such as QGIS, which offers extensive raster
97 and vector processing capabilities through graphical workflows. While effective for exploratory
98 analysis, desktop GIS tools are not designed for automated large-scale processing across many
99 NetCDF files, nor do they provide an integrated mechanism for consistently retrieving datasets
100 from multiple climate and geospatial repositories. GridFlow differs by offering both a CLI and
101 GUI that expose equivalent configuration options and consistent logging, enabling workflows
102 to be repeated and shared as version-controlled commands.

103 Overall, GridFlow contributes a modular and extensible “downloader hub” that integrates
104 acquisition and preprocessing for climate and geospatial datasets under a single interface. Its
105 design supports both global-scale datasets and subset-based workflows, enabling reproducible
106 preparation pipelines for research, education, and operational prototyping while reducing the
107 need for archive-specific scripts and manual portal interactions.

108 Roadmap and future development

109 GridFlow is actively being expanded toward a general-purpose “downloader hub” for climate and
110 geospatial datasets, where new sources can be added as independent modules while preserving a
111 consistent user experience. Planned future releases will prioritize additional datasets, improved
112 cross-platform packaging, and expanded preprocessing support. Figure 4 outlines a tentative
113 development roadmap; timelines are subject to change depending on community feedback,
114 available compute resources, and potential funding or sponsorship.

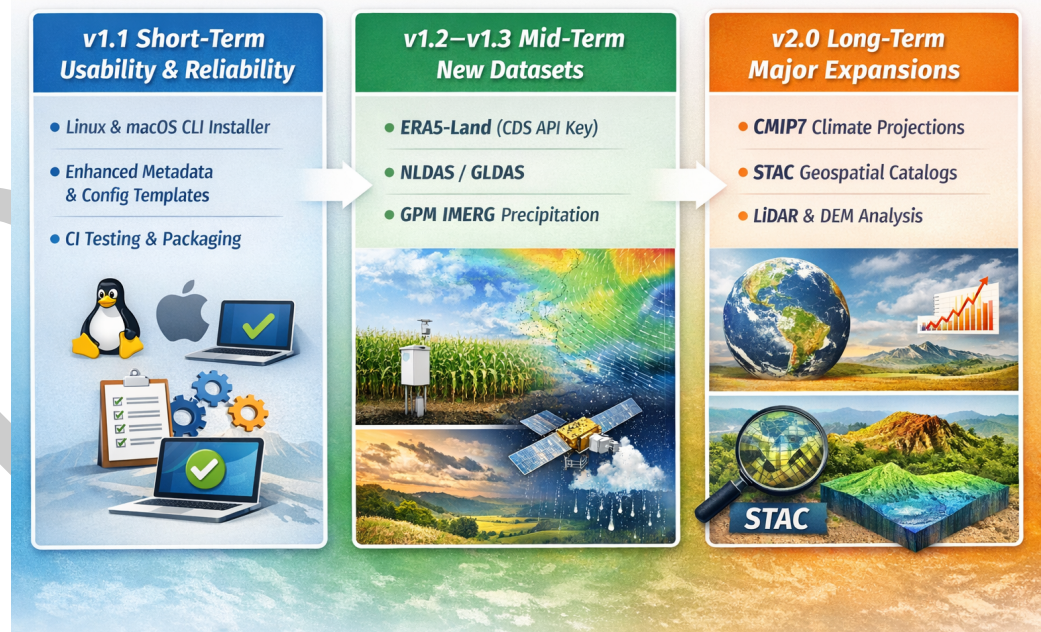


Figure 4: Tentative development roadmap for future GridFlow releases, including planned datasets and processing modules.

115 **Availability**

116 GridFlow (version 1.0) is released under the GNU Affero General Public License v3.0 (AGPLv3)
117 and is available on GitHub at:

118 <https://github.com/shahbhuwan/GridFlow>

119 **Acknowledgements**

120 The authors thank the open-source scientific Python community for foundational libraries that
121 enable GridFlow's functionality, including Xarray, GeoPandas, and the broader NetCDF and
122 geospatial ecosystem. The authors also acknowledge the data providers and maintainers of
123 ESGF, PRISM Climate Group, ECMWF (ERA5), Copernicus DEM, and USGS datasets for
124 making large-scale climate and geospatial resources publicly accessible.

125 This work was supported by the U.S. Department of Agriculture (USDA) Natural Resources
126 Conservation Service (NRCS) under Cooperative Agreement No. NR243A750008C001.

127 Bourgault, P., Huard, D., Smith, T. J., & others. (2023). Xclim: Xarray-based climate data
128 analytics. *Journal of Open Source Software*, 8(85), 5415. [https://doi.org/10.21105/joss.](https://doi.org/10.21105/joss.05415)
129 [05415](https://doi.org/10.21105/joss.05415)

130 Hoyer, S., & Hamman, J. (2017). Xarray: N-d labeled arrays and datasets in python. *Journal*
131 *of Open Research Software*, 5. <https://doi.org/10.5334/jors.148>

132 Jordahl, K., & others. (2020). *GeoPandas: Python tools for geographic data*. [https:](https://geopandas.org/)
133 [//geopandas.org/](https://geopandas.org/)

134 Vo, T., Po-Chedley, S., Boutte, J., Lee, J., & Zhang, C. (2024). xCDAT: A python package
135 for simple and robust analysis of climate data. *Journal of Open Source Software*, 9(98),
136 6426. <https://doi.org/10.21105/joss.06426>