

IWSG 2022

GeoEDF: A Framework for Designing and Executing Reproducible Geospatial Research Workflows in Science Gateways

Rajesh Kalyanam, Jungha Woo, Lan Zhao, Carol X. Song, and Jack Smith*

Rosen Center for Advanced Computing, Purdue University

*Marshall University

June 15-17, 2022, Trento, Italy



GeoEDF Vision



Researchers spend up to 80% of their time "wrangling data"

OUR DATA WORKFLOW - Ver. 1 2 3

1. Make sure date is just after 1st or 15th!
2. Go to: `usgs.gov/prd-tnm.s3.amazonaws.com/index.html?prefix=StagedProducts`
3. Browse: Hydrography...NH4PlusHR_Beta...GDB
4. Download NH4PLUS_H_01##_HUA_GDB.zip where ## is 02 to 14
5. Unzip it - WARNING: Have enough space!!!
6. Run our tool. **WARNING: Takes a loooooong time!!! DO NOT TURN OFF PC!!!**
7. Upload output files to cluster - Note: wait until all successful!
8. Kick off our standard jobs.
9. Occasionally-check-em. Wait for email(s)
10. Download new images.
11. Ask Fred to upload to website. *Mary?*
12. Tell everybody there's new stuff.

Process your files before running our tool

- get the latest code from *Nicole* for filtering the input data
- Note: that's Windows code - use the lab desktop
- you need to get the maps from the group folder `/depot/lyllegroup/project1/maps/exp/...` for aggregation



Remote data repos, smart devices, streaming data



GeoEDF Geospatial Data Framework

Reusable Data Connectors

Pluggable Data Processors

Integrated Active Learning

Resource Interoperability Interfaces

Geospatial Data Pipeline Composer (GUI & API)

Cyberinfrastructure (Campus, XSEDE, HUBzero, Geospatial Tools, storage, Solr, ...)

Make Science FAIR

OUR DATA WORKFLOW - Final

1. Go to the science gateway
2. Define "my_workflow.yml" (or use tool GUI if needed)
3. Ask GeoEDF to execute!
4. Data and workflow automatically published to science gateway



Remote data directly usable in code, seamless workflow

Complexity abstracted away

Reusable data connectors, processors, and workflows

Automatic provenance capture & data annotation => FAIR

GeoEDF Project

An Extensible Geospatial Data Framework Towards
FAIR Science

To help data-driven sciences to be more
Findable, Accessible, Interoperable, Reusable

funded by NSF CSSI program award #: 1835822, Oct 2018 - Sep 2023



GeoEDF Components

Reusable Data Connectors

Implement various data access protocols, enable data acquisition from popular repositories



Reusable Data Processors

Implement domain agnostic & domain specific geospatial processing operations



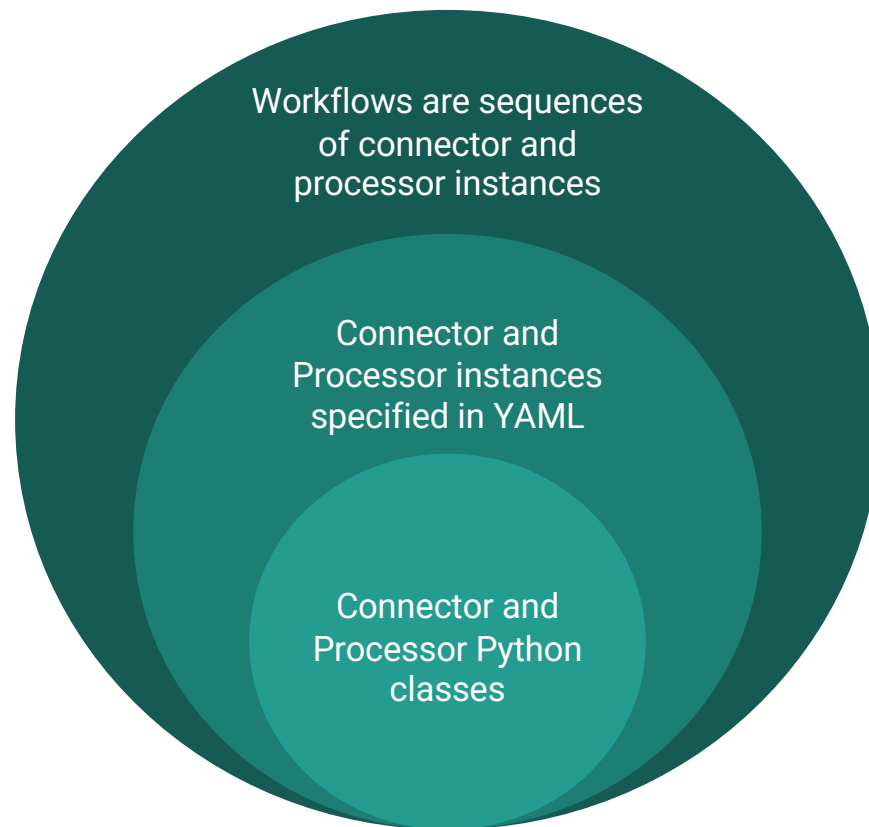
Plug-and-play Workflow Composer

Enable the composition of individual connectors & processors into complex workflows

GeoEDF

Enable researchers to conceive of geospatial data driven workflows as a sequence of data acquisition and processing steps that can be carried out using pre-existing or user contributed connectors and processors

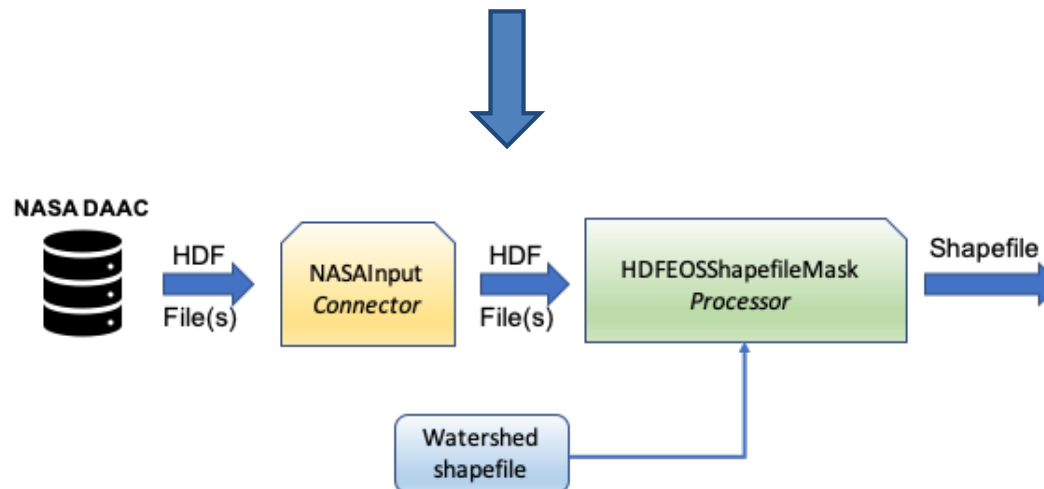
GeoEDF Workflow in a nutshell



Example Hydrologic Workflow



Apply GeoEDF principles



Corresponding GeoEDF Workflow (YAML)

Data connector

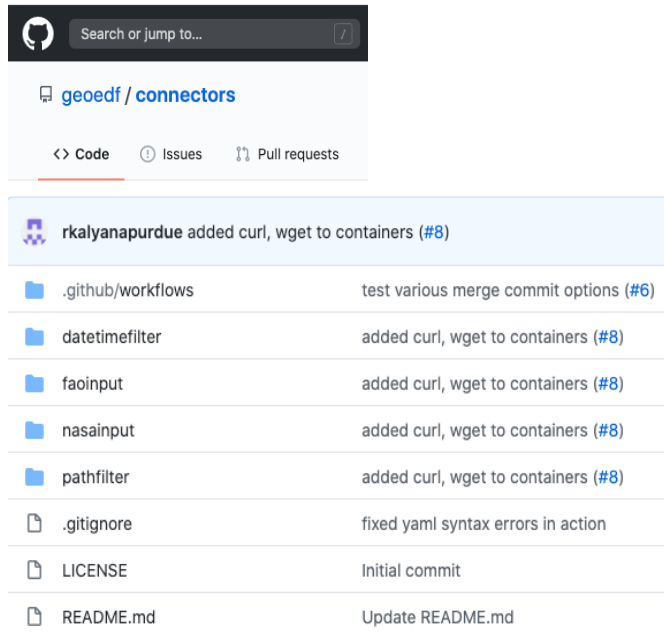
```
$1:  
  Input:  
    NASAInput:  
      url: https://e4ftl01.cr.usgs.gov/MOTAMCD15A3H.006/{file}  
      user: rkalyana  
      password:  
  Filter:  
    file:  
      PathFilter:  
        pattern: '{dtstring}/MCD15A3H.*.h09v07*.hdf'  
      dtstring:  
        DateTimeFilter:  
          pattern: '%Y.%m.%d'  
          start: 07/16/2002
```

- Filters enable spatial and temporal filtering before data acquisition
- This improves workflow generality and efficiency

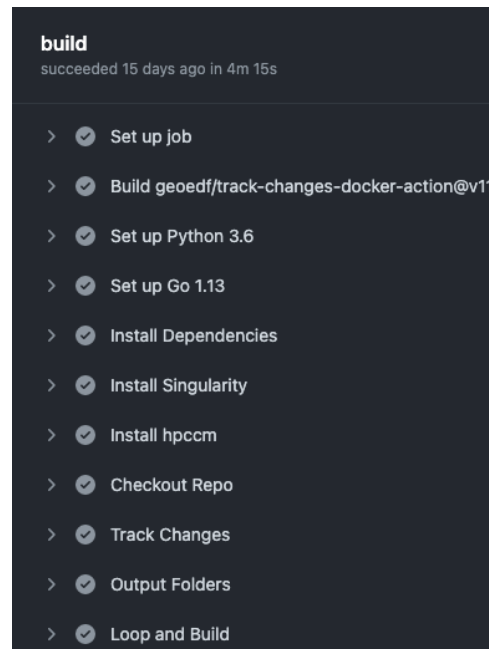
Data processor

```
$2:  
  HDFEOSShapefileMask:  
    hdf file: $1  
    shapefile: /home/mygeohub/rkalyana/subs1_projected_171936.shp  
    datasets: [Lai]
```

Connector, Processor Contribution Process



(1) Contribute connectors/processors via GitHub pull requests

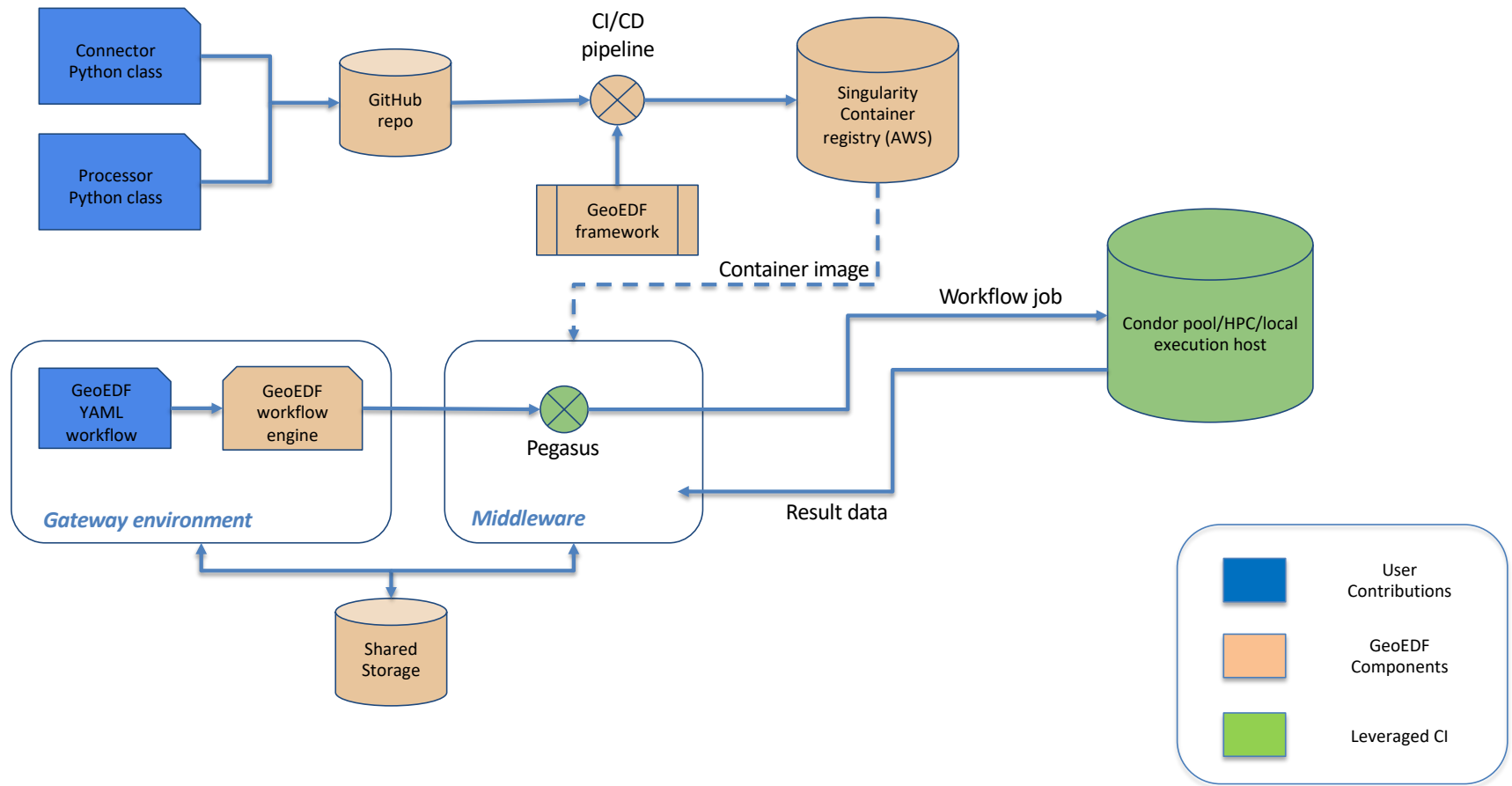


(2) GitHub action detects changes, builds Singularity container, pushes to registry server



(3) Workflow engine queries registry for list of connector, processor containers

Gateway Integration



Deployment Solutions



```
In [1]: import hublib.use
       use pegasus-5.0.1.dev
       use geodefengine-1.0

In [2]: from geodefengine.GeoEDFWorkflow import GeoEDFWorkflow

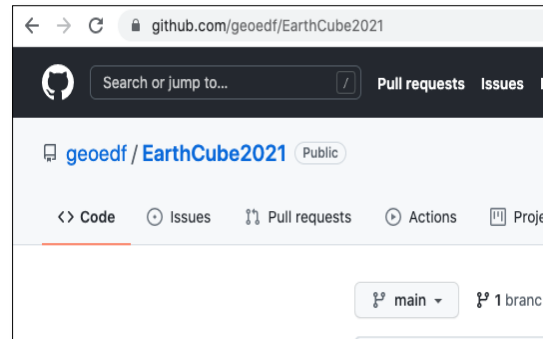
In [4]: workflow = GeoEDFWorkflow('/home/mygeohub/rajkala/workflows/files/wgpinput.yml')

In [5]: workflow.execute()
Workflow created and written to /home/mygeohub/rajkala/workflows/1621541097
Workflow submitted for execution; outputs will be written to /home/mygeohub/rajkala/workflows/1621541097

In [6]: workflow = GeoEDFWorkflow(workflow_dir='/home/mygeohub/rajkala/workflows/1621541097')

In [15]: workflow.monitor()
workflow is complete; check pegasus.analysis file in this directory to check success
```

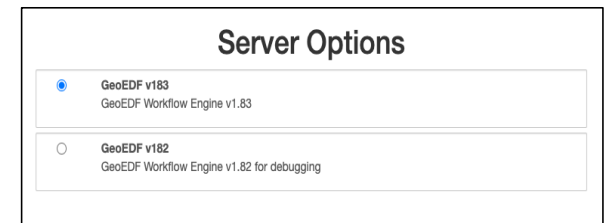
- Publicly available gateway
- Deployed in Jupyter notebook environment as a Python library
- Job submission to Purdue's Halstead cluster



- Self-contained Docker container
- Can use to build and test new connectors, processors
- Run on your own machine



Sign in with CILogon



- Standalone deployment, in the works...
- CILogon authentication
- Workflow execution in local minicondor

GeoEDF Demo

GEOEDF WORKFLOW EXECUTION IN THE MYGEOHUB GATEWAY

GeoEDF Applications – Water Quality

Synthesize hydrologic and water quality data from various federal agencies (USGS, EPA, etc.) for EPSCoR states for ease of visualization and analysis

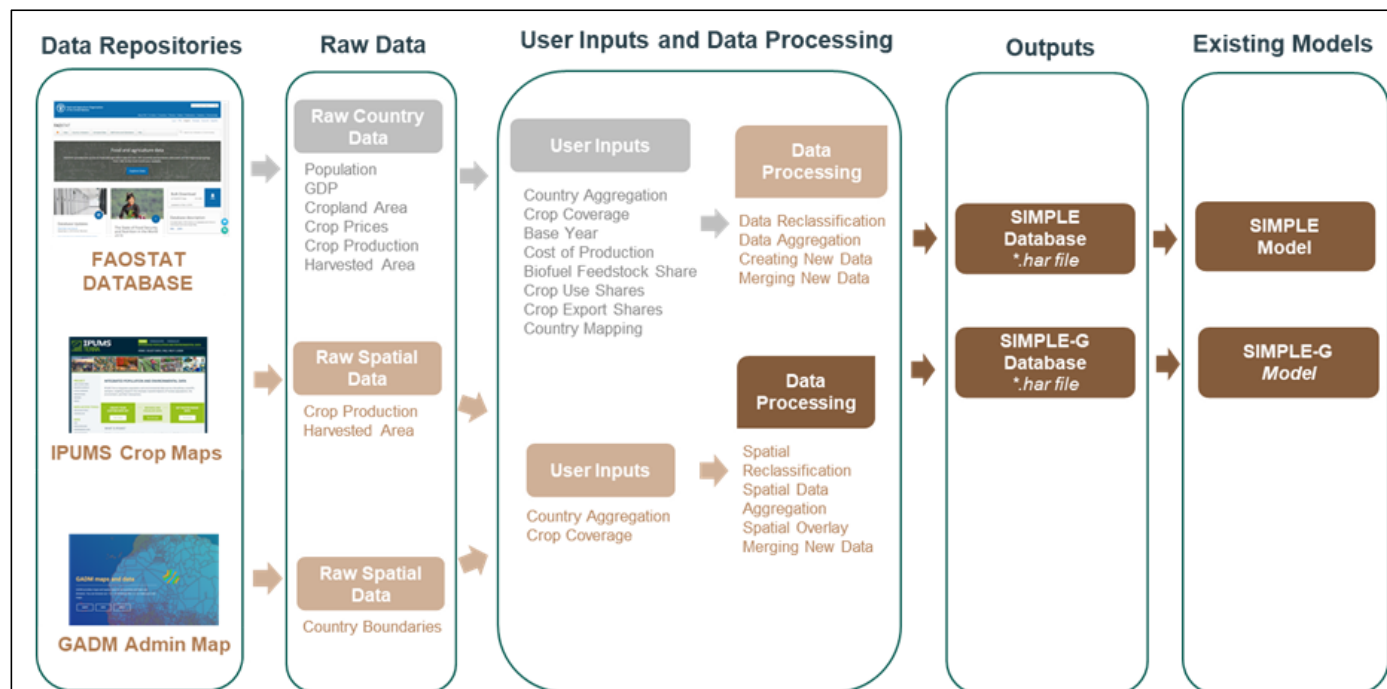
- Workflow produces an interactive map combining water quality data from WQP and stream reach data for a given monitoring station



GeoEDF Applications – Agricultural Economics

Acquire and pre-process the necessary socio-economic, agricultural, and climate data for analyzing global-to-local food security and sustainability

- Workflow acquires diverse U.N. FAOSTAT datasets, aggregates it for the study region, and converts from custom “HAR” format into widely-used csv

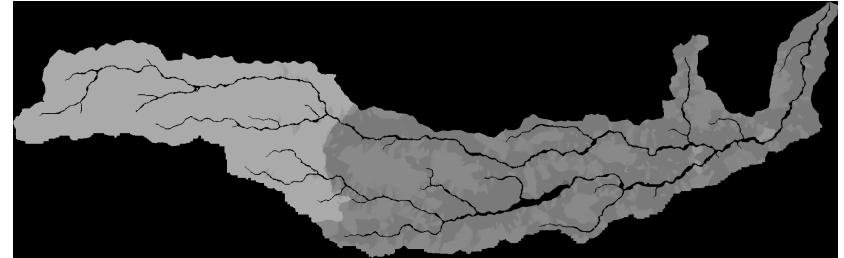


GeoEDF Applications – InVEST

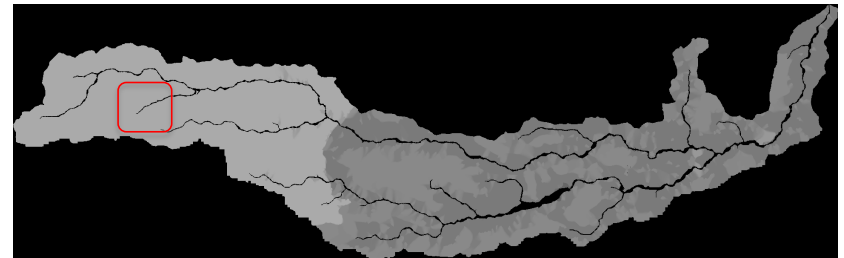
Family of tools (models) for quantifying importance of natural capital

- Workflow wraps the InVEST NDR (nutrient delivery ratio) model to enable efficient parameter sweeps via HPC execution

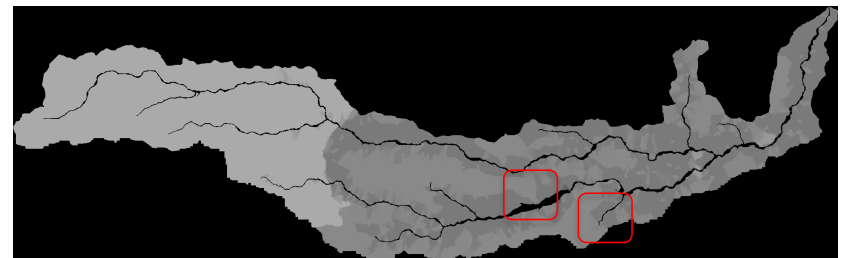
`threshold_flow_accumulation = 1000`



`threshold_flow_accumulation = 1500`



`threshold_flow_accumulation = 2000`



Summary

- Research reproducibility via declarative workflow conceptualization
- Broad applicability across domains that have workflows with a mix of data acquisition and processing steps
- Variety of integration options with CI and gateway platforms
- Ability to leverage various compute resources (local machine, Condor pool, HPC)

Thank You

- GeoEDF GitHub Repository: <https://github.com/geoedf>
- GeoEDF Documentation: <https://geoedf.readthedocs.io>
- Publication: <https://dl.acm.org/doi/10.1145/3311790.3396631>

Email:

Carol Song (cxsong@purdue.edu)