# OKG-Soft: An Open Knowledge Graph with Machine Readable Scientific Software Metadata

**Daniel Garijo**, Maximiliano Osorio, Deborah Khider,
Varun Ratnakar and Yolanda Gil

University of Southern California,
Information Sciences Institute

**@dgarijov**

# Science is changing: Open Science

**Open data**

**Open access**

**Impact and credit**

**Open publications**

**Open source software**

# The importance of Scientific Software

**Open data**



- Software helps understand data
  - Provenance, reproducibility

- Software helps understanding methods
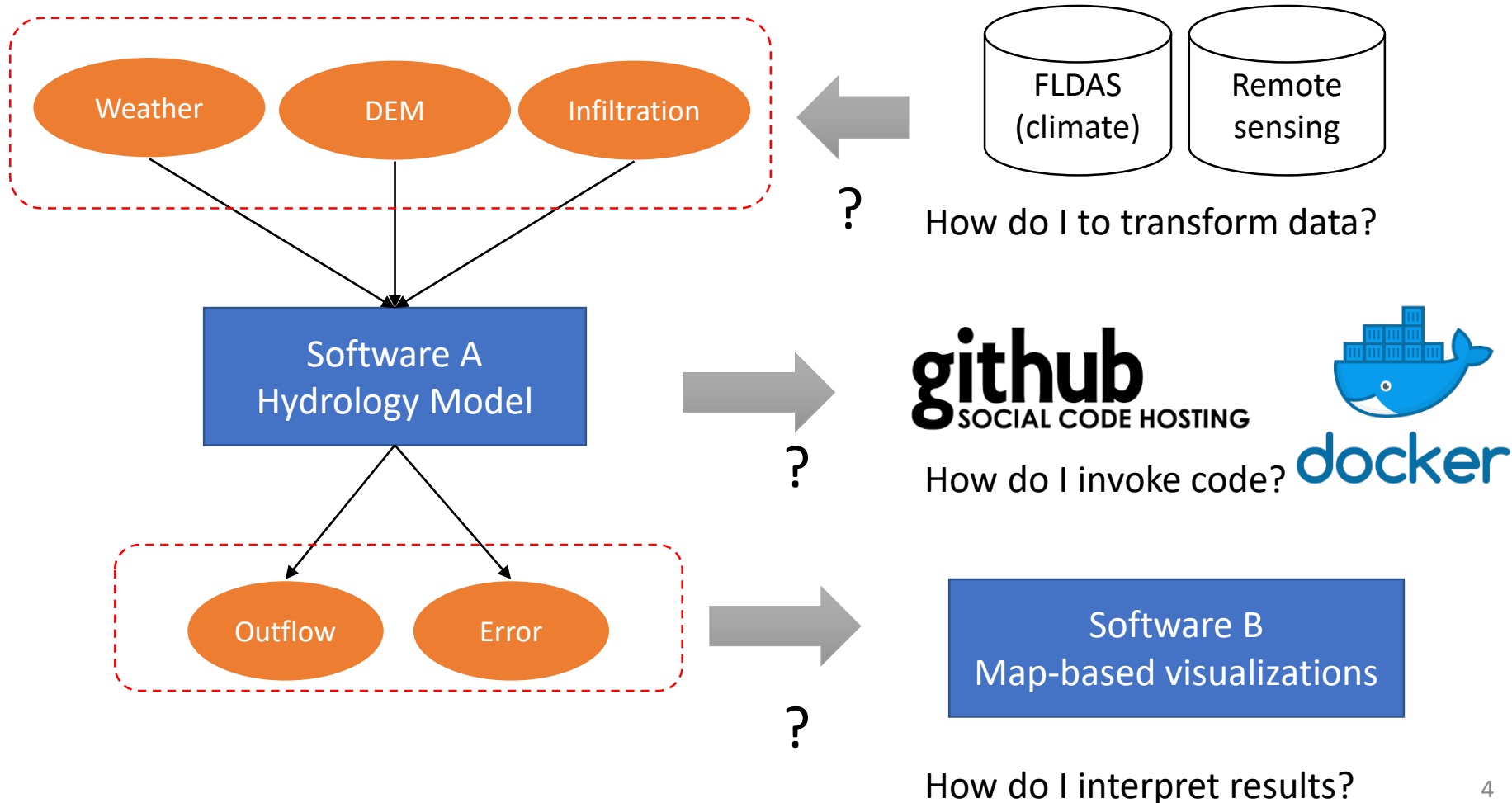  - Assumptions, limitations

**Open source software**



**Open publications**



3

# Why is it difficult to reuse Scientific Software?

Let's imagine we want to reuse existing work:



How do I to transform data?

How do I invoke code?

How do I interpret results?

4

# Outline
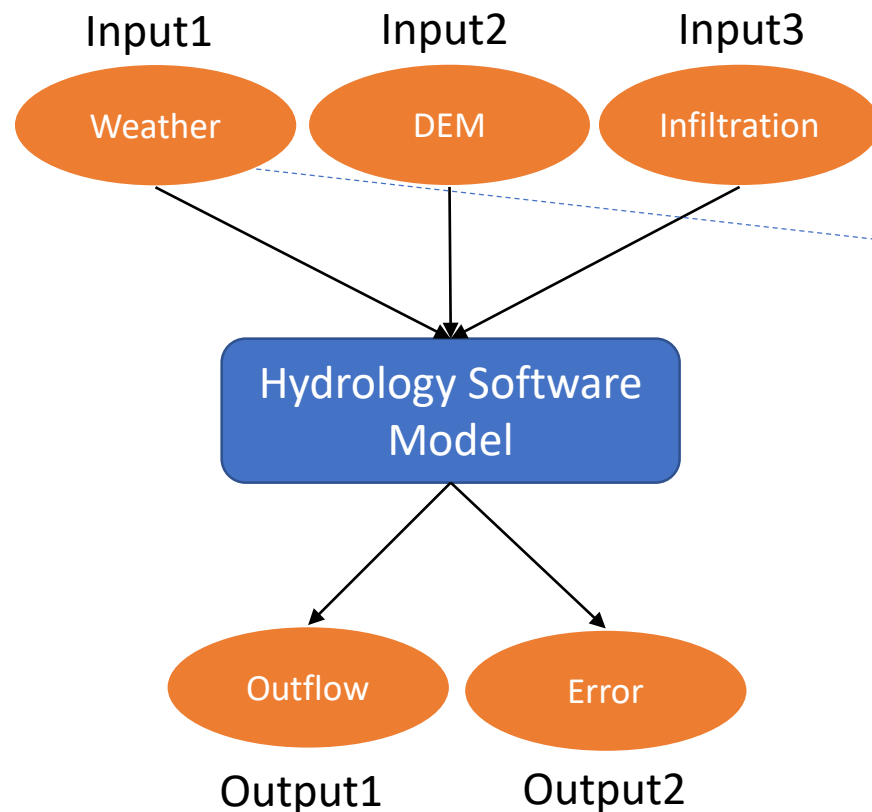
1. Requirements help scientific software reusability

2. Our current approach for representing scientific software metadata

3. A framework to query, explore, exploit and publish software metadata

# Outline

1. **Requirements help scientific software reusability**

2. Our current approach for representing scientific software metadata

3. A framework to query, explore, exploit and publish software metadata

# Requirements for Software Reusability

1. Exposing software inputs, outputs and their corresponding variables



- Land surface temperature (degC)
- Precipitation rate (mm/h)
- Land surface wind speed (m/day)
- Net radiation (MJ/(day m^2)

# Requirements for Software Reusability

1.  Exposing software inputs, outputs and their corresponding variables
2.  Capturing the functions of the software being used

```
Hydrology Software Model

Function A: Richards
Equation for water
movement (unsat soil)

Function B: Saint Venant
equations
(shallow water)
```

# Requirements for Software Reusability

1. Exposing software inputs, outputs and their corresponding variables

2. Capturing the functions of software being used

3. Using principled ontologies with structured names for model variables, processes, and methods
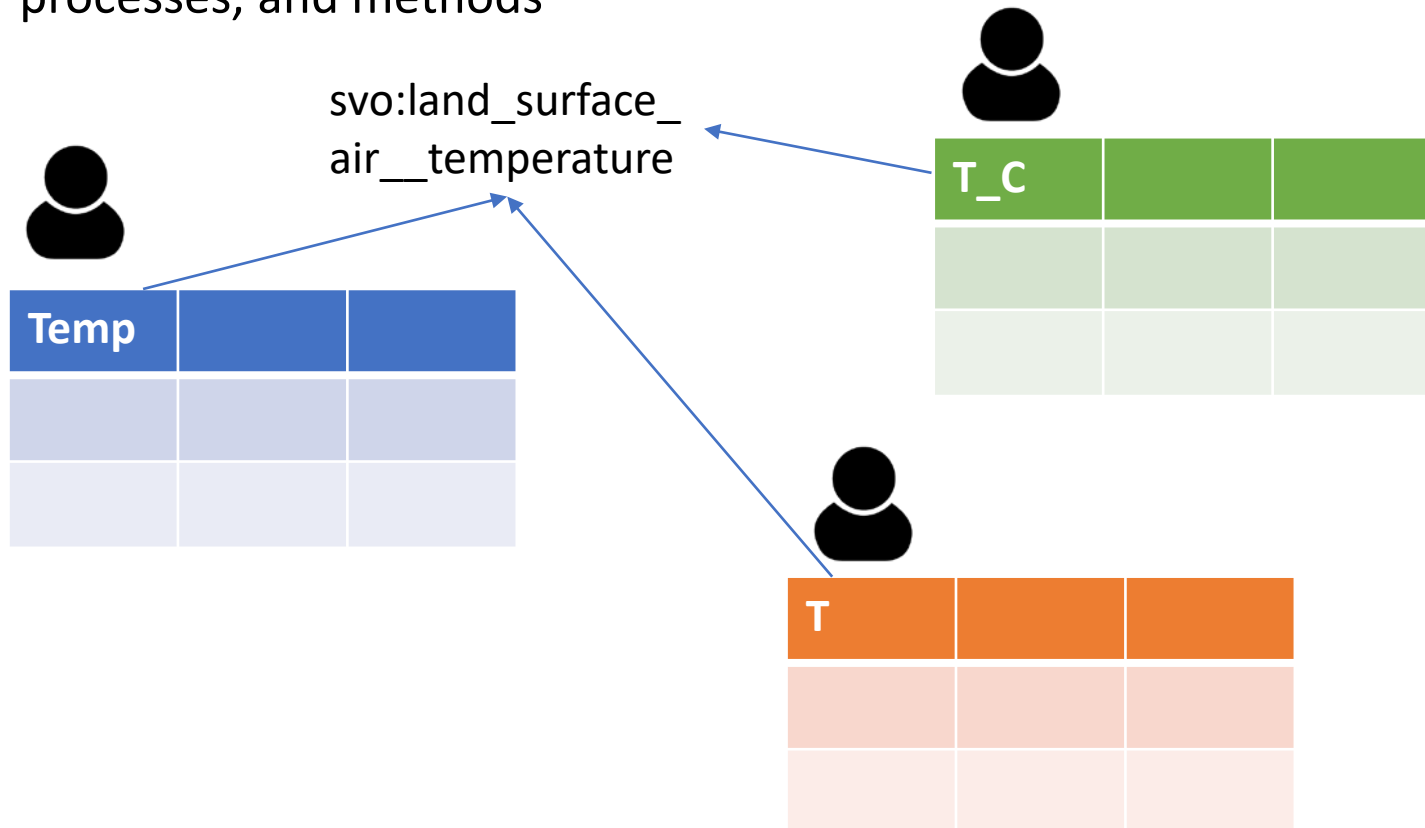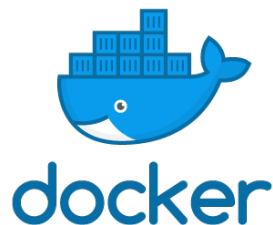


svo:land_surface_
air__temperature

# Requirements for Software Reusability

1. Exposing software inputs, outputs and their corresponding variables

2. Capturing the functions of software being used

3. Using principled ontologies with structured names for model variables, processes, and methods

4. Capture the semantic structure of software invocations

Dependencies?
Sample runs?
Invocation command?
Is data supposed to be in the same folder?
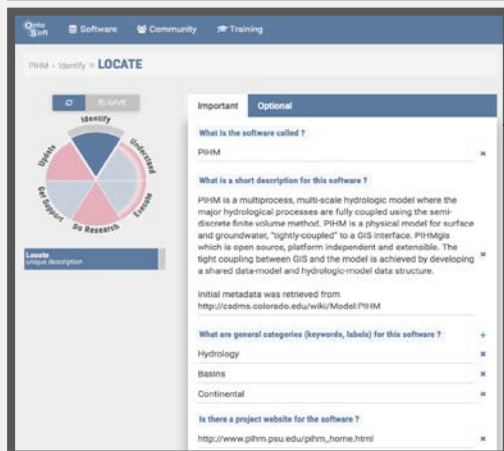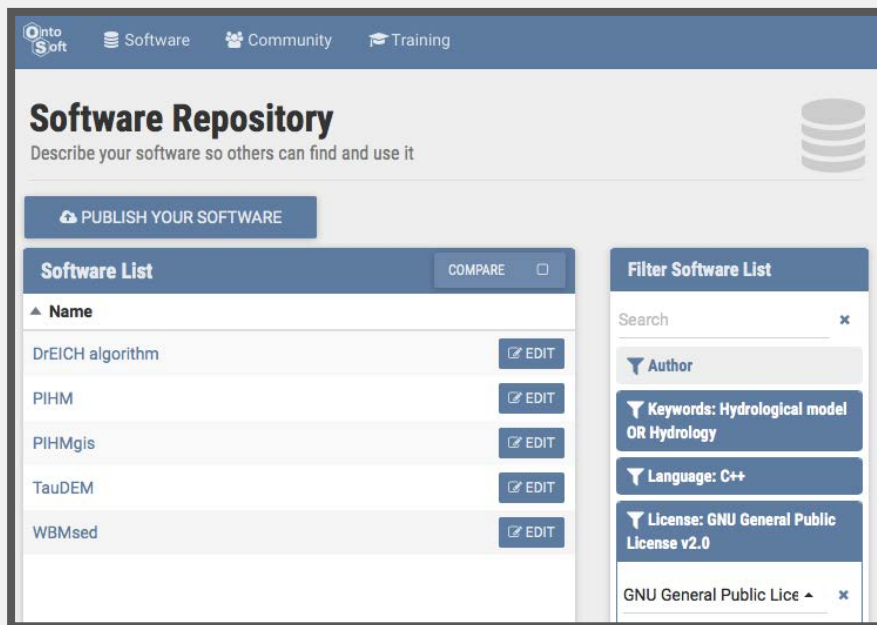Default arguments/Configuration files?
Volumes?
Do I have to log in in the image

# Outline

1. **Requirements** help scientific software reusability

2. **Our current approach for representing scientific software metadata**

3. A framework to query, explore, exploit and publish software metadata

# Prior Work: OntoSoft Software Metadata Registry



**Finding Software**

## OntoSoft

Model and Software Metadata Registry

- Complements code repositories to make them understandable

- Software metadata designed for scientists

- Metadata is curated by decentralized communities of users

- Training scientists on best practices



## http://ontosoft.org

[Gil et al 2015]: OntoSoft: Capturing Scientific Software Metadata Eighth ACM International Conference on Knowledge Capture, Palisades, NY, 2015

12

# Prior Work: OntoSoft Software Metadata Registry

# Evolving OntoSoft: Software Description Ontology



Extensions:

- Schema.org (software metadata)
- W3C Data Cubes (Contents of inputs and outputs)
- NASA QUDT (Units)
- DockerPedia (Software images)
- Scientific Variables Ontology (Standard Variables)

https://w3id.org/okn/o/sd#

14

# Evolving OntoSoft: Extending schema.org and Codemeta

| author | Organization or Person | The author of this content or rating. Please note that author is special in that HTML 5 provides a special mechanism for indicating authorship via the rel tag. That is equivalent to this and may be used interchangeably. |
|---|---|---|
| citation | CreativeWork or URL | A citation or reference to an article, etc. |
| contributor | Organization or Person | A secondary contributor to the CreativeWork. |
| copyrightHolder | Organization or Person | The party holding the legal copyright to the CreativeWork. |
| copyrightYear | Number | The year during which the claimed copyright for the CreativeWork was first asserted. |
| creator | Organization or Person | The creator/author of this CreativeWork. This is the same as the Author property for CreativeWork. |
| dateCreated | Date or DateTime | The date on which the CreativeWork was created or the item was added to a DataFeed. |
| dateModified | Date or DateTime | The date on which the CreativeWork was most recently modified or when the item's entry was modified within a DataFeed. |
| datePublished | Date | Date of first broadcast/publication. |

Codemeta Terms

https://w3id.org/okn/o/sd#

15

# Evolving OntoSoft: Software Description Ontology
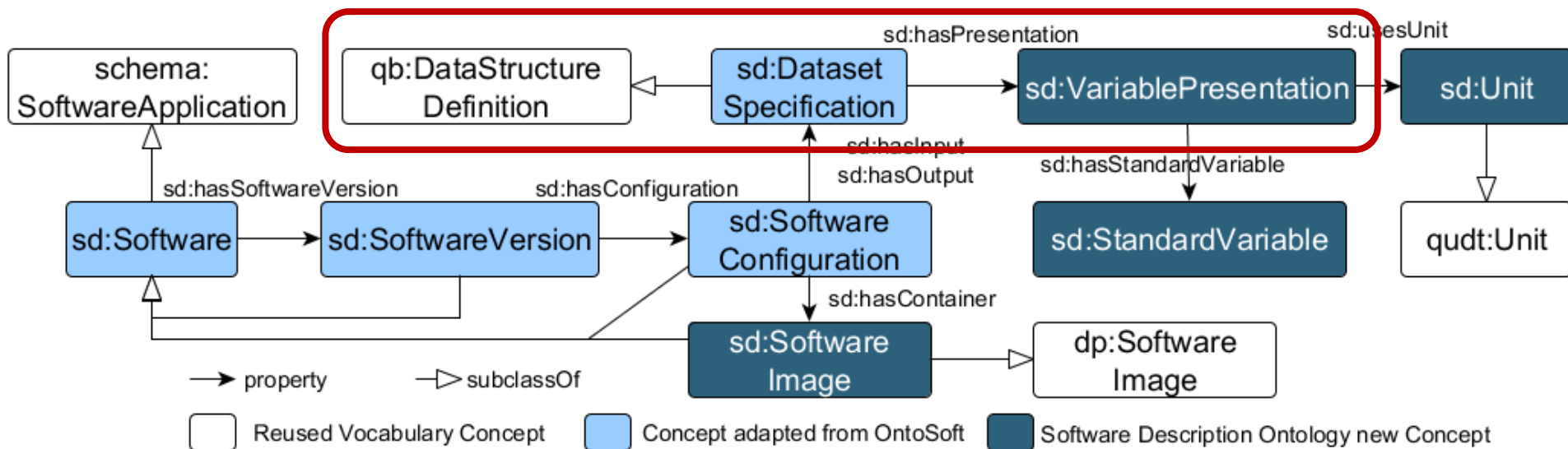


Extensions:

- Schema.org (software metadata)
- W3C Data Cubes (Contents of inputs and outputs)
- NASA QUDT (Units)
- DockerPedia (Software images)
- Scientific Variables Ontology (Standard Variables)

https://w3id.org/okn/o/sd#

16

# Describing Input/Output files, parameters and variables

Input files

Name

cycles_soil

cycles_reinit

cycles_cropname

cycles_weather

cycles_percent_increase_fertilizer

Output files

Name

cycles_som

cycles_seasonConfig

cycles_water

| cycles_soil | Cycles soil file | | | |
|---|---|---|---|---|
| Label | Long Name | Description | Standard Name | Units |
| DZ | soil layer thickness | Soil layer thickness | soil_layer__thickness | m |
| SLOPE | slope of the field | Average slope of field of interest | land_surface__slope | m m-1 |
| BD | bulk density | Soil mass dry and wihtout rock divided by the sampled volume | soil~no-rock~dry__mass-per-volume_density | Mg m-3 |

Cycles weather file

Cycles increase in fertilizer    txt

Cycles annual SOM file

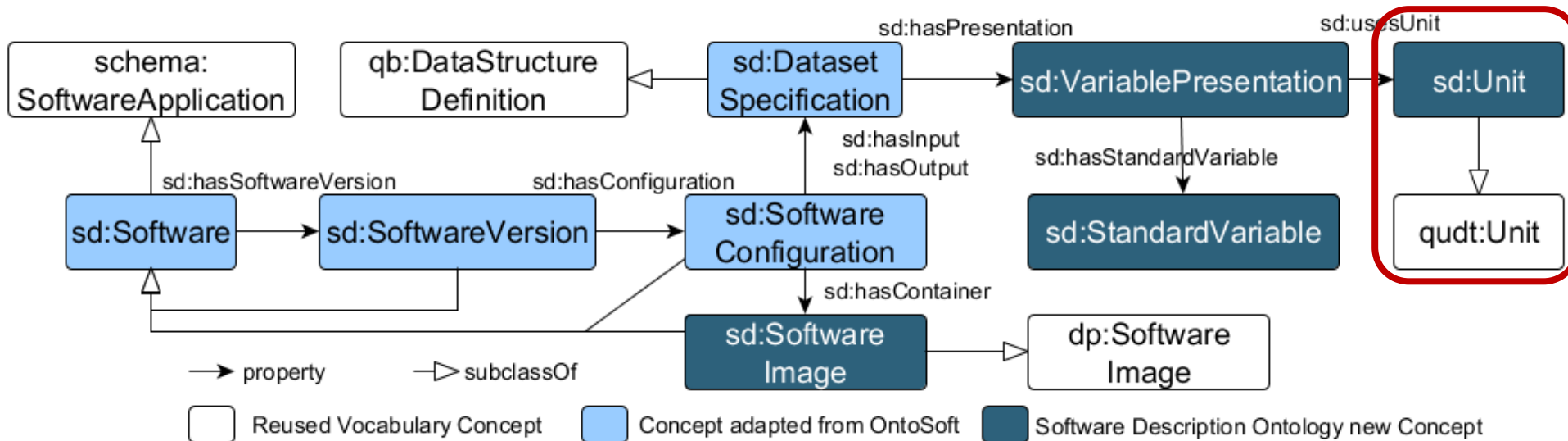Cycles season configuration file

Cycles water file

Scientific Variables Ontology identifiers

M. Stoica and S. D. Peckham, "An Ontology Blueprint for Constructing Qualitative and Quantitative Scientific Variables," in *Proceedings of the ISWC 2018 Posters & Demonstrations, Industry and Blue Sky Ideas Tracks co-located with 17th International Semantic Web Conference (ISWC 2018), Monterey, USA, October 8th - to - 12th, 2018.*, 2018

## Parameters:

| Name | Description | Default value |
|---|---|---|
| cycles_planting-day-1 | Day of the year when the planting started<br>The range is from 1 to 365 | 100 |
| cycles_planting-day-1-duration | Duration of planting (in days) | 10 |

# Evolving OntoSoft: Software Description Ontology



Extensions:

- Schema.org (software metadata)
- W3C Data Cubes (Contents of inputs and outputs)
- NASA QUDT (Units)
- DockerPedia (Software images)
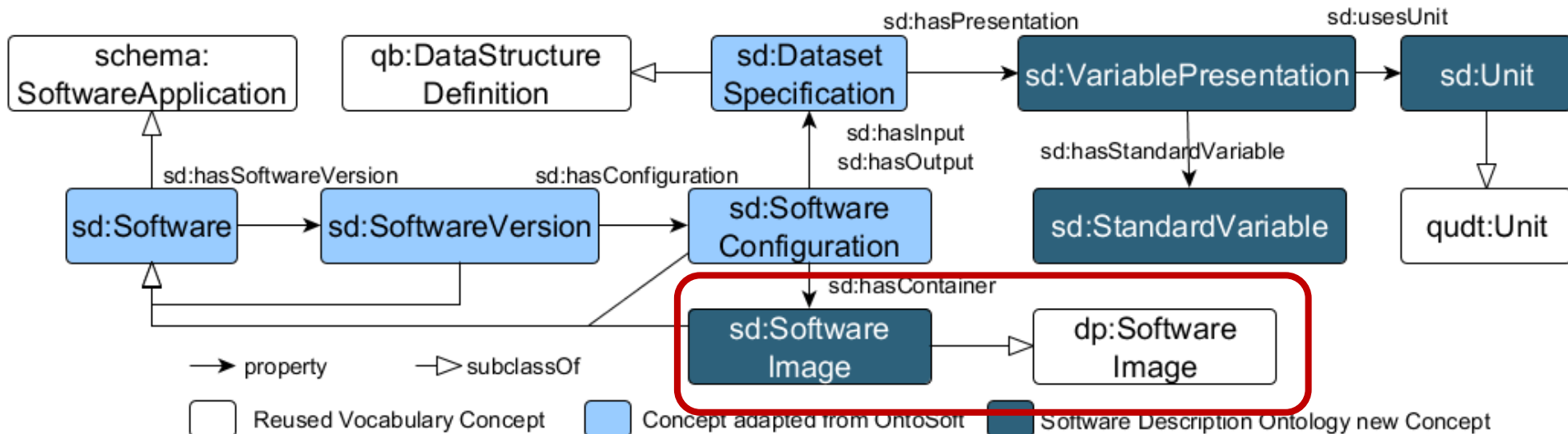- Scientific Variables Ontology (Standard Variables)

https://w3id.org/okn/o/sd#

18

# Machine readable representation of units



"RAIN"
mm/day

CCUT

Linking and augmenting from Wikidata

19

- B. Shbita, A. Rajendran, J. Pujara, and C. Knoblock, Parsing, Representing and Transforming Units of Measure, in Modeling the World's Systems, 2019.

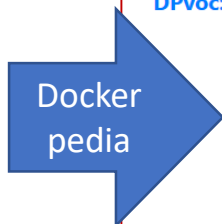# Evolving OntoSoft: Describing Containers



Extensions:

- Schema.org (software metadata)
- W3C Data Cubes (Contents of inputs and outputs)
- NASA QUDT (Units)
- DockerPedia (Software images)
- Scientific Variables Ontology (Standard Variables)

## https://w3id.org/okn/o/sd#

# Evolving OntoSoft: Describing Containers

| | |
|---|---|
| **rdfs:label** | **mintproject/pihm2cycles** |
| **DPvoc:size** | **272551998** |
| **rdf:type** | DPvoc:SoftwareImage |
| **DPvoc:imageIdentifier** | **mintproject-pihm2cycles_1.1** |
| **DPvoc:tag** | **1.1** |
| **DPvoc:composedBy** | DPres:ImageLayer/sha256%3Acc1a78bfd46becbfc3abb8a74d9a70a0e0dc7a5809bbd12e814f9382db003707<br>DPres:ImageLayer/sha256%3A49eab01d36f3e1b840ac3380b77dcb03f23d9f51baeece25086e314562581398<br>DPres:ImageLayer/sha256%3Ac2c2cfea02132c74bdec895f226ba7ee60b13e13f4549fe19d0a353c82bb817d<br>DPres:ImageLayer/sha256%3A419499c9a4cf0cd98be64b355f96799bbb6f9bf34932733bc3ade822a8cbc17f<br>DPres:ImageLayer/sha256%3Af7550509e92e740029588e5ed5a7c0ea6a363724db2b8653cd31029a6230b050<br>DPres:ImageLayer/sha256%3Aab6c636c45aac563fccf031b503bea35d8cd87b31fe8b7e468d288bc5bec4cc0<br>DPres:ImageLayer/sha256%3Addb63a5e3b0c0a77e6c9eef5eaa3e7c1787de61bfe1dc3d5a5af8f0f3a42daa1<br>DPres:ImageLayer/sha256%3Aa49272fbc797b544f7964c6403a3b82b96ecdff83e9fb58a64eb1502f7386589<br>DPres:ImageLayer/sha256%3Aa7e7574dc810f7edbbce9f20b4d25772e4980b46f6ba1c8277a4fd02aebbaf64<br>DPres:ImageLayer/sha256%3A72486d8655e5ddbe258310bcc80201bd6a8834f1c4a0b3693ab5eb26df74f85b |
| **DPvoc:containsSoftware** | zlib_1:1.2.8.dfsg-5<br>shadow_1:4.4-4.1<br>mawk_1.3.3-17<br>apt_1.4.8<br>nettle_3.3-1<br>iproute2_4.9.0-1+deb9u1<br>debconf_1.5.61<br>nghttp2_1.18.1-1<br>libbsd_0.8.3-1<br>util-linux_2.29.2-1+deb9u1 |

Image semantic representation

Image tag

**"mintproject/ pihm2cycles"**

Docker pedia

Page 1 >

https://dockerpedia.inf.utfsm.cl/

• M. Osorio, H. Vargas, and C. Buil Aranda, "DockerPedia: a Knowledge Graph of Docker Images," in Proceedings of the ISWC 2018 Posters & Demonstrations, Industry and Blue Sky Ideas Tracks co-located with 17th International Semantic Web Conference (ISWC 2018), Monterrey, 2018.

# Outline

1. Requirements help scientific software reusability

2. Our current approach for representing scientific software metadata

3. **A framework to query, explore, exploit and publish software metadata**

# OKG-SOFT: Framework

Software Model Catalog contains:

- Models from hydrology, agriculture and economy, their versions and model configurations.
  - More than 200 variables mapped to SVO.
  - All models are executable through scientific workflows
  - Most contents are added manually (expert users) **collaboratively**

- Automated unit transformations
- Automated software image description
- Semi-automated Wikidata linking

APIs:
- SPARQL endpoint
- REST APIs (GET/POST)  https://query.mint.isi.edu/api/mintproject/MINT-ModelCatalogQueries#/
- Python clients

swagger

# Exploitation: Exploring Scientific Software Model Metadata



Find Software Models

Explore Software I/O

Explore variables

http://models.mint.isi.edu

24

# Exploitation: Comparing Scientific Software Models

Comparing:      ✕

### TopoFlow

Select version
Topoflow v3.5

Select configuration
TopoFlow with basic configuration

Select calibration
TopoFlow calibrated model for South Sudan

To:      ✕

### Penn State Integrated Hydrology Model (PIHM)

Select version
PIHM V4

Select configuration
PIHM++ configuration for version v4 with aggregated outputs

Select calibration
PIHM++ v4 configuration (v4) calibrated for South Sudan (Pongo Region) with

**Model comparison:**

|  | TopoFlow | Penn State Integrated Hydrology Model (PIHM) |
|---|---|---|
| **Page:** | Model documentation | Model documentation |
| **Creation date:** | 2002 | 2007 |
| **Funding:** | National Science Foundation | National Science Foundation |
| **Publisher:** | University of Colorado-Boulder | The Pennsylvania State University |
| **Type:** | Theory Guided | Theory Guided |

**Configuration comparison:**

|  | TopoFlow with basic configuration | PIHM++ configuration for version v4 with aggregated outputs |
|---|---|---|
| **Description** | A basic configuration of the TopoFlow model | PIHM++ configuration for version v4 aggregating all outputs in a single zip file |
| **Parameter assignment method** |  | Expert-configured |
| **Target variables** |  | pihm_streamflow_ph |
| **Spatial dimensionality** | 2D | 2D |
| **Spatial grid type** | SpatiallyDistributedGrid | SpatiallyDistributedGrid |
| **Spatial grid resolution** | 100x100m | 50m-200m |

http://models.mint.isi.edu

# Exploitation: Towards Automated Software Composition



**DRIVING VARIABLES**

- rainfall (atmosphere_water__rainfall_mass_flux, atmosphere_water__globe_time_average_of_rainfall_volume_flux, atmosphere_water__geologic_time_average_of_rainfall_volume_flux, atmosphere_water__domain_time_integral_of_rainfall_volume_flux)
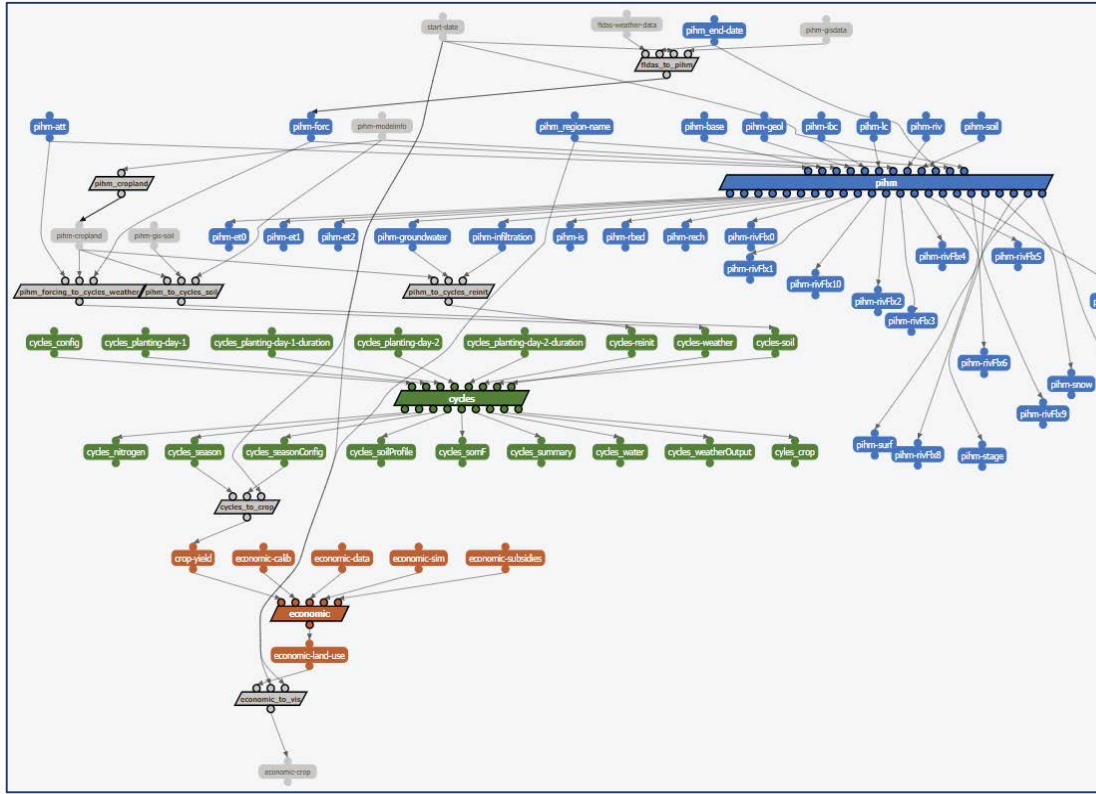
**RESPONSE VARIABLES:**

- crop_production (crop__produced_mass)

**SELECT MODEL COMPOSITION**

1. cycles / economic / pihm

2. cycles / economic / topoflow

# Summary

Scientific software reusability is crucial to understand
- Existing data
- Published methods

1. Requirements for scientific software reusability include
   - Expose inputs, outputs, variables and software invocation details!

2. Our approach for capturing and structuring scientific software

3. A framework to query, explore, exploit and publish software metadata

# Help us making your software more reusable

Contact me: dgarijo@isi.edu