

Semantic web for scientific information

Streamlining how we write, find, link and reuse data and models

Ferdinando Villa
&
Integrated Modelling Partnership

bc³
BASQUE CENTRE
FOR CLIMATE CHANGE
Kiima Aldaketa Ikergai



aquacross

An integrated solution for shared, distributed, collaborative modelling

SEMANTICS FOR DATA AND COMPUTATIONS

- Maintenance of the core conceptualization and language
- Maintenance and delivery of a **shared worldview** (ontologies) for cross-domain communication



OPEN SOURCE SOFTWARE

- User-end (modelers and end users)
- Server technology (institutions)
- Developer team and user support



INTEGRATED MODELING INFRASTRUCTURE

- Assembly of models from networked data and model components
- Partners can manage their servers or use the partnership's

APPLICATIONS

- Ecosystem services assessment (ARIES)
- Real-time monitoring using remotely sensed data
- Food and other environmental securities
- Integrating hydrology, primary production, nutrients with agent models of SES.



COLLABORATIVE MODELING

- Interoperable data and models
- Serving models on the Web
- Direct support of partner projects
- [International Spring University](#) since 2013



Models and data live on a semantic web

An extensible network hosts data, models and model services available to users

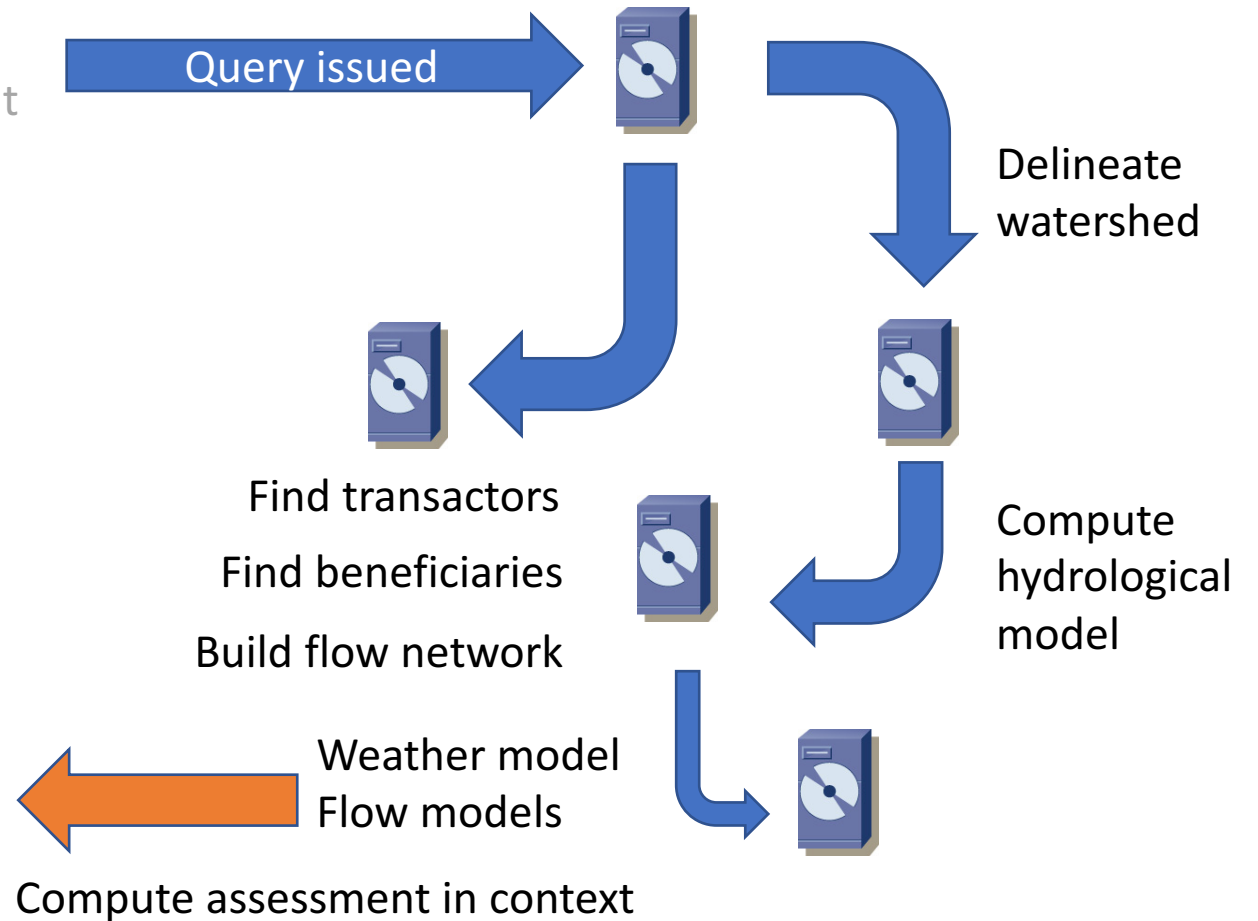
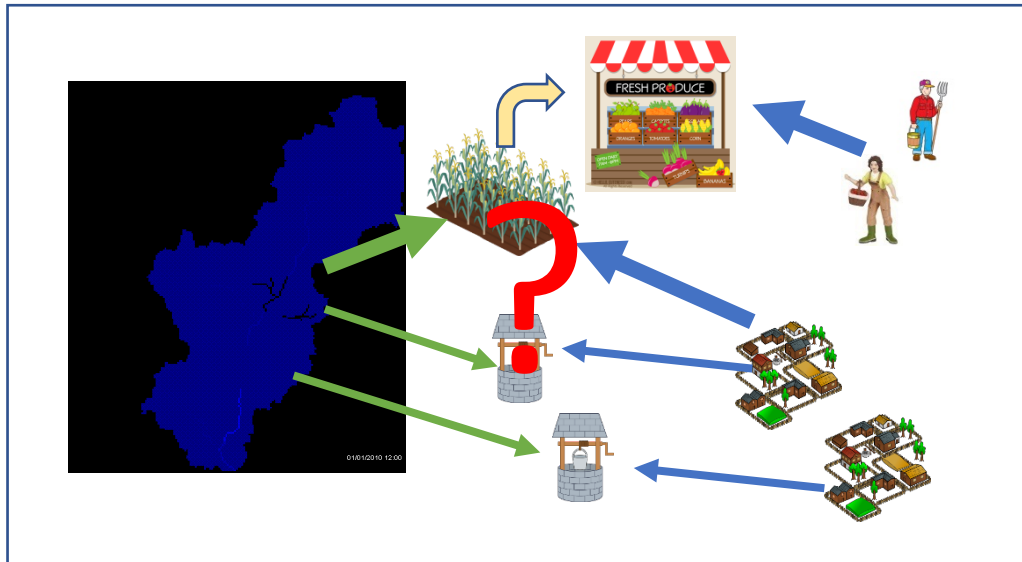


Query:

1. Set context to region X
2. Observe water social dynamics in it



Results!



A user's perspective: two-step assessment

Client software (desktop & soon web-based) allow modeling with minimal configuration and training. Provenance info is compiled into user documentation for each result set.

Ecosystem Services toolkit

This toolkit gives you access to the most common concepts used to assess ecosystem services according to the ARIES methodology. It also contains some example locations for testing and some fully finished case studies so that you can learn to interpret results. You can change the toolkit as you like, adding and removing concepts, observations or roles. Close this description when you have become familiar with it.

Ecosystem Benefits

These processes define all ecosystem services. For most of them, sources and transactors need to be identified or computed. Beneficiaries can also be added to compute actual values. You can drag and drop them on a map to assess each of them in their most likely location, or create a location of your choice and drag/drop one or more on it.

- Water supply
- Carbon services
- Aesthetics
- Hydropower
- Raw materials
- Cultural
- Sediment

Aesthetic Roles

The roles of provider, transactor and beneficiary for aesthetic services define sources of aesthetic enjoyment, locations from where those are appreciated, and likely beneficiaries of those. Use these roles to parameterize your world before you compute an aesthetic service.

- Beauty
- Viewpoint
- Viewers
- Visual blight

Common aesthetic assets

This toolbox contains object types commonly considered in scenic value assessments. Drop the tool to have k.LAB find them, or use Ctrl+drop to create some yourself.

- Mountain Peaks
- Rest areas
- Middle-class groups
- Lakes
- Ocean

Test areas and case studies

Drag-and-drop paradigm for end users

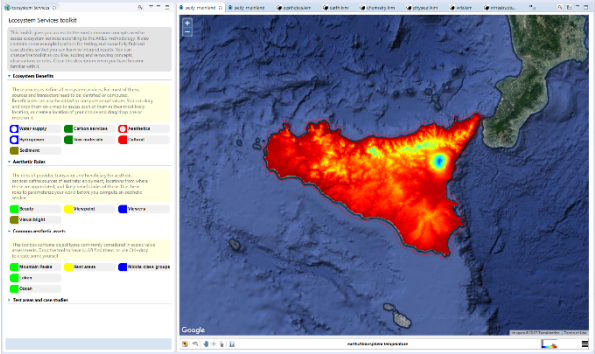
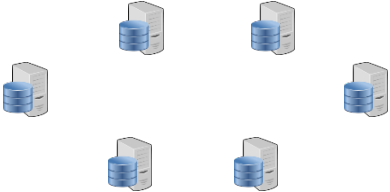
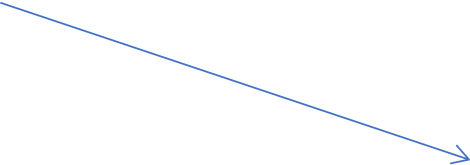
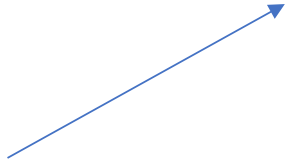
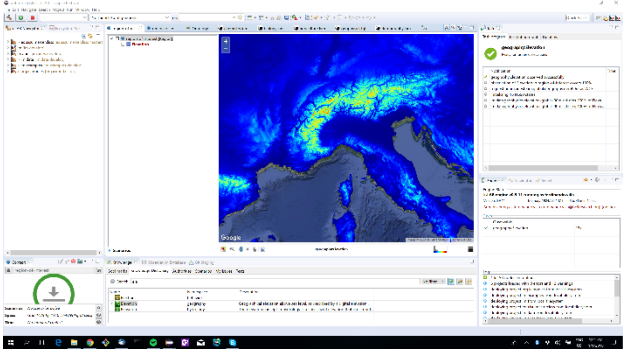
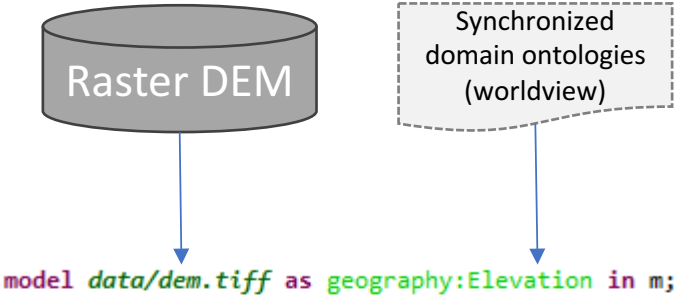
“Palette” of Ecosystem Services tools can store finished studies and scenario results builds best-case model out of components and data and computed when user drops the concept computes it...

Full reports are built to document the computation logged into network secure certificate

St
Se
co
(se
or

A modeler's/data scientist's perspective

streamlining and supporting the annotation workflow for both data and computations



One language for knowledge, data and computations

Uses English-like syntax to express the observation ontology; linguistic approach keeps ontologies small and learnable

```
model hydrology:SurfaceWaterFlow,  
  hydrology:RunoffVelocity,  
  hydrology:RunoffWaterVolume named runoff  
observing  
  earth:PrecipitationVolume,  
  geography:Elevation,  
  earth:AtmosphericTemperature,  
  earth:Stream  
using  
  hydrology.swat.distributed()  
over time (step = "1 day")  
  do [  
    ... actions  
  ],  
  change runoff to [  
    ...  
  ]  
;
```

Hydrological model

```
learn value of ecology:Biodiversity 0 to 100  
observing  
  conservation:Protected ecology:Vegetation earth:Site as im:Archetype,  
  not conservation:Protected ecology:Vegetation earth:Site as im:Archetype,  
  geography:Aspect in degree_angle as im:ExplanatoryQuality,  
  geography:Slope in degree_angle as im:ExplanatoryQuality,  
  im:Annual im:Mean earth:AtmosphericTemperature in Celsius as im:ExplanatoryQuality,  
  im:Annual earth:PrecipitationVolume in mm as im:ExplanatoryQuality,  
  distance from landcover:UrbanFabric earth:Region in m as im:ExplanatoryQuality,  
  distance from earth:Coast in m as im:ExplanatoryQuality,  
  distance from infrastructure:Road in m as im:ExplanatoryQuality,  
  proportion of soil:Silt in soil:TopSoil im:Volume as im:ExplanatoryQuality  
using weka.bayesnet();
```

Machine learning model

Ontological statements read as English and are validated while editing. Inconsistent concepts are flagged as errors and discarded.

Color coding, assisted editing, and informative error messages help user

Often dozen of OWL axioms in 1-2 lines

```
thing Coast  
  "A portion of land adjacent to a major marine or lacustrine water body."  
  is earth:Terrestrial earth:Region  
    adjacent to (earth:Marine or earth:Lacustrine) earth:Region;  
  
thing Coastline  
  "The boundary between land and an adjacent @Coast."  
  is im:Boundary  
    of (earth:Region adjacent to (earth:Marine or earth:Lacustrine) earth:Region)  
    adjacent to (earth:Marine or earth:Lacustrine) earth:Region;
```

from IM worldview (general users receive it from the network and search it)

Worldviews merge domains and vocabularies reliably and intuitively

```
namespace chemistry
  using im, physical
  in domain im:Chemistry;

abstract identity Compound
  "Concrete subclasses of Compound must be identified with
  an InChI code validated by the IUPAC authority."
  is ChemicalSpecies
  requires authority IUPAC
  has disjoint children
    (Water identified as "1S/H2O/h1H2" by IUPAC),
    (CO2 identified as "1S/CO2/c2-1-3" by IUPAC),
    (NH3 identified as "1S/H3N/h1H3" by IUPAC),
    (H3O identified as "1S/H2O/h1H2/p+1" by IUPAC),
    (SO4 identified as "1S/H2O4S/c1-5(2,3)4/h(H2,1,2,3,4)" by IUPAC),
    (NaCl identified as "1S/ClH.Na/h1H;/q;+1/p-1" by IUPAC);

quality MassConcentration
  is ratio of ${inherent extends ChemicalSpecies} im:Mass
  to ${context extends ChemicalSpecies} im:Volume;

@origin("SWEET")
quality Acidity
  "Capability of a molecule to donate a hydron (proton or hydrogen ion H+)."
  is MassConcentration of H3O within Water
;

quantity Ph
  "A measure of acidity or alkalinity of an aqueous solution used universally, with
  values of 7 indicating neutrality. Computed as the negative logarithm of the activity
  of H ions."
  decreases with Acidity
;
```

Explicit domains

Assisted editor with as-you-type syntax validation

- syntax embodies observation semantic
- predicate and observable composition
- usage of attributes

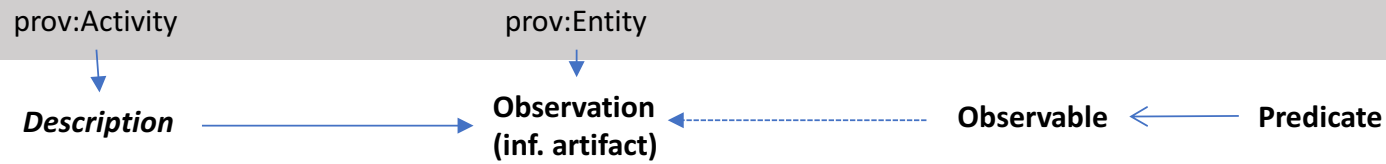
Logical validation at each save

- Uses reasoner of choice
- Consistency of usages with upper ontology
- (potentially nested) inherency...

Bridging to vocabularies through endorsed **authorities**

- GBIF (taxonomic identities)
- IUPAC (chemical identities)
- WRB (soil identities)
- AGROVOC (agri processes and practices)
- ... (plug-in)

Provenance ontology

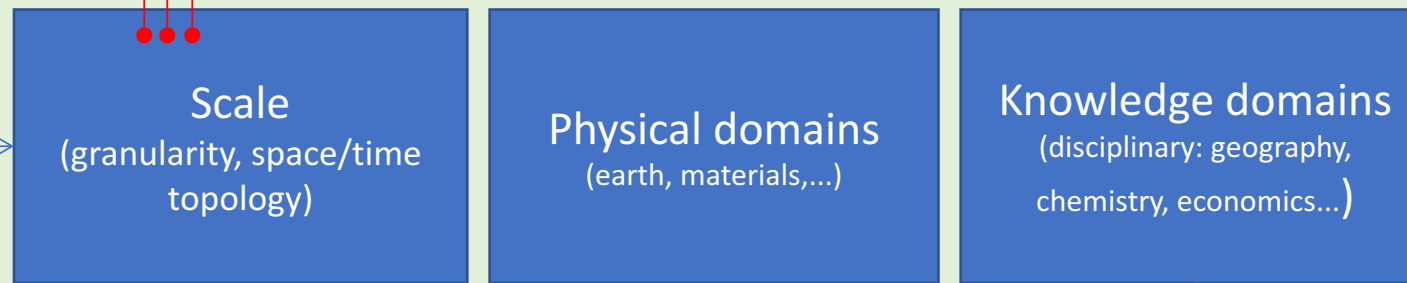
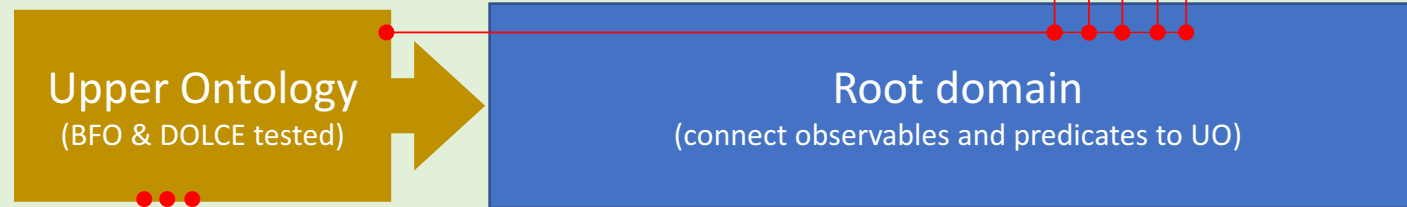


Acknowledgement
Computation
Instantiation
Resolution
Emergence

Structural description (graph, ...)
Dynamic description (DDEs, dataflow...)
Collection (objects: database, vector...)
State (map, timeseries...)

Subject	Attribute
Agent	Realm
Event	Identity
Process	Ordering
Quality	Role
Physical property	
Numerosity	
Value	
Type	
Extent	
Relationship	
Functional	
Structural	

Observation ontology



Worldview

urn:klab:node:user:namespace:id

Data **geometry** becomes **scale** through the worldview

Roles reinterpret physical "reality" according to disciplines

A semantics-first approach - for wide user groups

Address all the “W’s of information – what, where, when, why, and how – without becoming too large or complex to learn and use.



SUBJECTS:

A mountain

A group of humans

A forest

A river

QUALITIES:

Elevation (measurement)

Per capita income (value)

Percent tree canopy cover (%)

Stream order (ranking)

PROCESSES:

Erosion

Migration

Tree growth

Streamflow

EVENTS:

Snowfall

A birth

Death of a tree

A flood event

RELATIONSHIPS:

↖ Skiers using a mountain for recreation ↗

↖ A city using a river for water supply ↗

Semantics for **predicates** allow to compose attributes, realms and identities without inheritance; interface to vocabularies
Roles account for **usages** of general observables in disciplinary contexts without giving up consistency and FAIR goals

Tooling (1): k.IM language and support software

```
role PollinatorSupplier
  is ses:Provider
  applies to earth:Region
  implies PollinatorAbundance as ses:Supply;

role AgriculturalProductionDependent
  is ses:Beneficiary
  implies PollinatedYield as ses:Demand
  applies to observation:Subject;

/**
 * Roles that define the P->T and B->T relationships.
 */

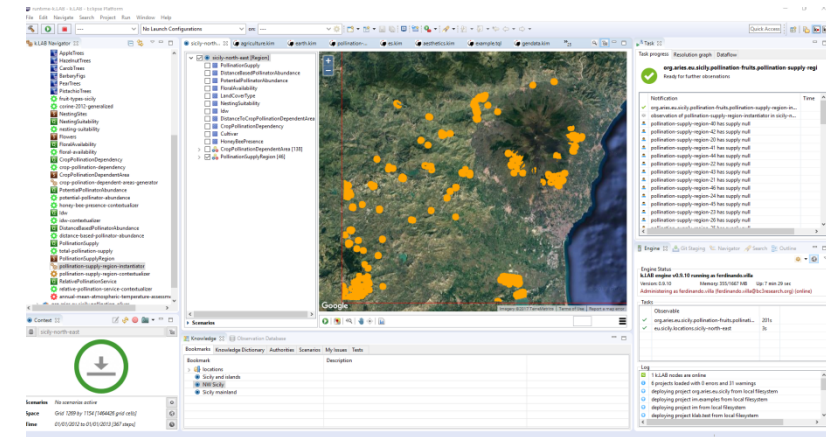
role PollinationSupplyConnection
  is ses:ProvisionFlow
  applies to im:MatterTransferConnection between PollinatorSupplier and PollinationDependent;

role AgriculturalUseConnection
  is ses:UseFlow
  applies to im:MatterTransferConnection between AgriculturalProductionDependent and PollinationDependent;

/**
 * Role for the ES, tying everything together.
 */

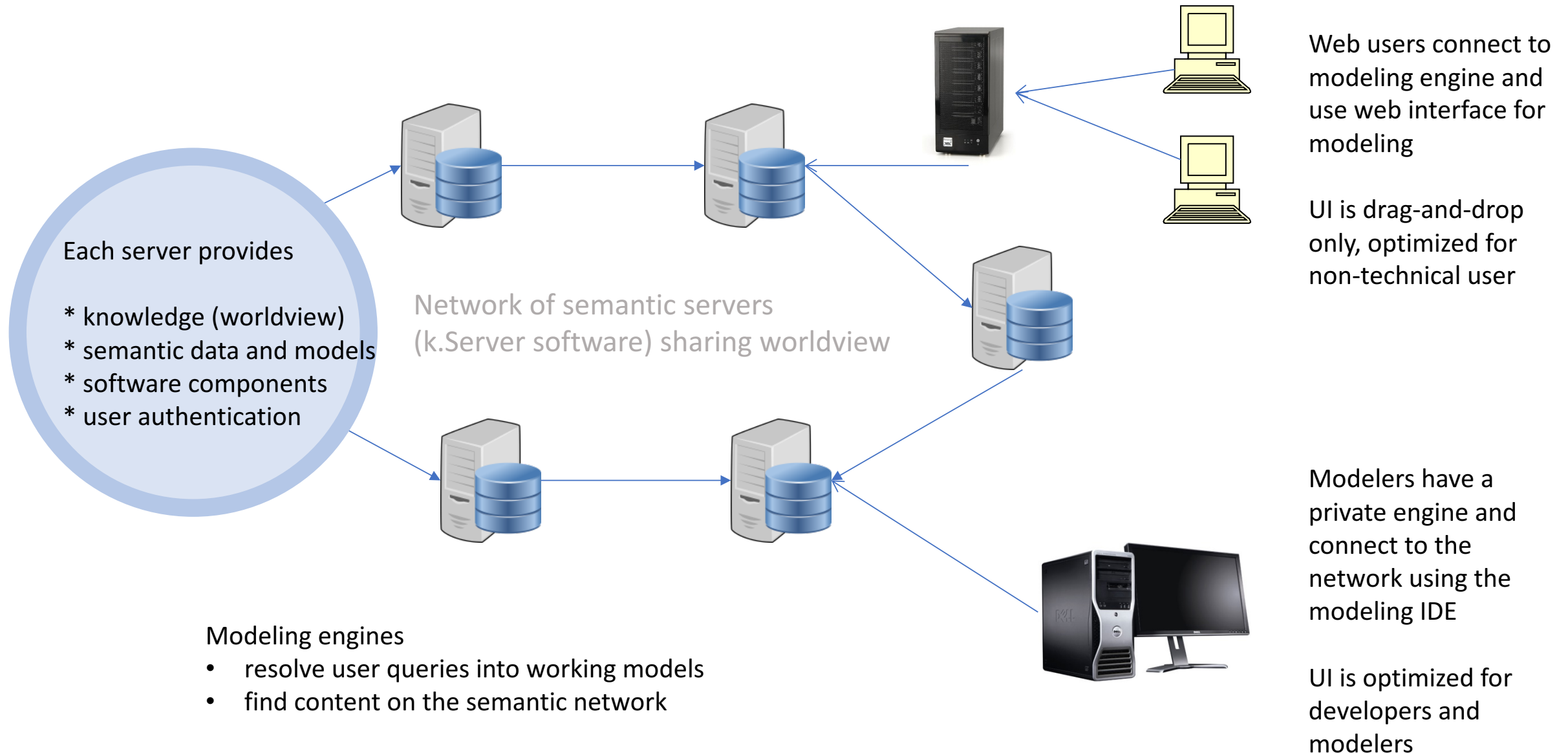
role PollinationEcosystemBenefit
  "The benefit obtained by any user of the yield made possible by pollination. This is
  easier to monetize than most ES when defined this way."
  is ses:ProvisioningEcosystemBenefit
  implies at least 1 PollinationSupplyConnection, at least 1 AgriculturalUseConnection
;
```

The k.IM language is used to express both the worldview and the data/models that use it



- Tools and interfaces enable [end users](#), [modelers](#), and [network administrators](#)
- Simplify the tasks of semantically describing, coding, and publishing data and models.
- Provide and maintain documentation, community resources for [discussion](#), [user support](#) and [bug reporting](#)
- Create [tools for participatory, graphical model building](#) that can be directly translated into templates for working models.

Tooling (2): distributed semantic web infrastructure

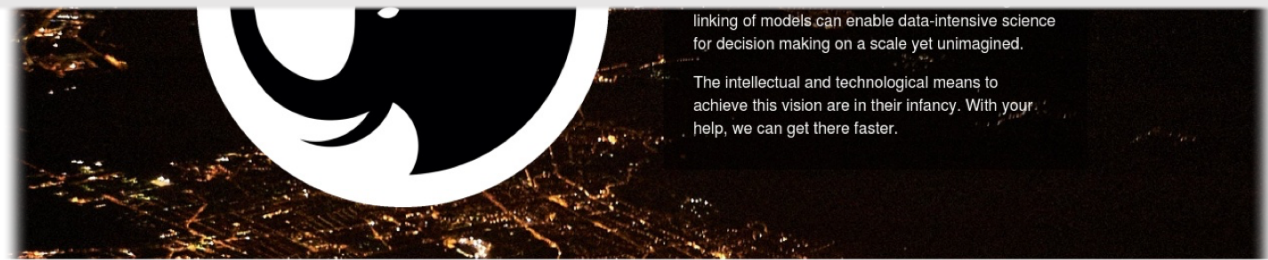


The Integrated Modelling Partnership

http://www.integratedmodelling.org

Partners share: Maintains:

- **Participation**
 - Downloadable software and documentation for users and developers.
 - design their applications in collaboration with the core staff
- Certification for users and institutions
 - Work Packages drive the development of larger initiatives
- The shared *worldview for all domains, prioritizing needs of partners.*
- **Ownership**
 - Online data and models annotated with worldview semantics and identified by unique courses) are open source/open access
 - Supports partners in deploying their own servers and modeling engines.
 - Products bear the copyright of the partnership with its member institutions.
- Online courses and training material on integrated modelling.
- **Control:**
 - partners enter the steering committee that defines activities, governance and directions of online support channels.
 - development



linking of models can enable data-intensive science for decision making on a scale yet unimagined.

The intellectual and technological means to achieve this vision are in their infancy. With your help, we can get there faster.

The **Integrated Modelling Partnership**, begun in 2017, brings together institutions contributing to designing and building a fully integrated information landscape for the science of the future.

The partnership develops and maintains the [IM worldview](#), the [k.IM language](#) and the [k.LAB software stack](#). It provides [training in semantic modelling](#) and supports partners and users in creating unprecedented model-data integration in [projects](#) such as [ARIES](#).

Become a partner to participate in building the vision, knowledge, and tools to support a more efficient, integrated, and democratic scientific process.

[Learn more](#)

[Become a partner](#)

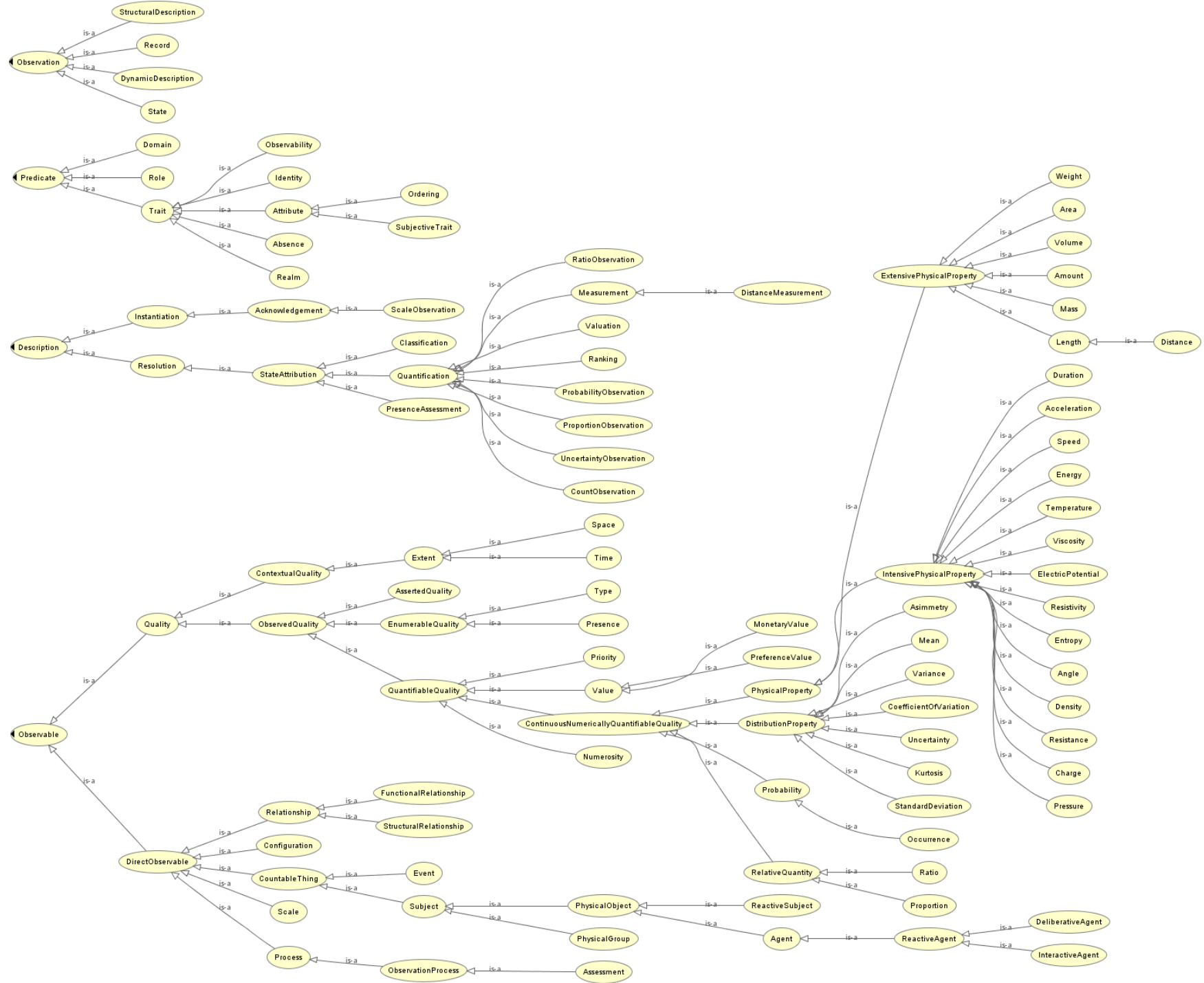
Four building blocks for one new approach



1. Semantics
The language used to describe scientific observations must be flexible and shareable, without ambiguity. It must efficiently address



2. Open, linkable data
Making data and models FAIR is complex and requires understanding of ... and agreement on ... the nature of all scientific



Promoting **semantics-first** solutions for open, linkable data and models

- In today's dialogue on interoperability, *information* equals *data*. Both data and models can be seen as ways to make scientific observations – *definitions* of observations. Doing so enables a consistent discussion on how to semantically connect data to models and how to build complex models by assembling simpler ones.
- Semantics-first data/models are *first-class research objects* that can be found online, read and understood by computers and humans alike. They can reuse existing vocabularies and thesauri while ensuring consistent semantics throughout the information landscape.
- Powered by semantics, artificial intelligence can transparently match data and models to a chosen time, place, problem, and (multiple!) scales.
- **Much of the complexity of building and running models can be handled by machines, with substantial advantages for science and decision making.**

The challenge of data/model integration and reuse

Scientists in the past collected data in notebooks. In the digital age, we want scientific data and models to be **FAIR** - [Findable, Accessible, Interoperable, and Reusable](#), to ensure their maximum value.

A fully connected information landscape using open, safe, accurate, “Wikipedia-like” sharing and linking of models can enable data-intensive science for decision making on a scale yet unimagined:

1. **reuse** the abundance of data and specialized knowledge available and needed to analyse social and natural processes (and their interactions)
2. **avoid** the risk of **fragmentation** hidden in the use of ad-hoc (or no) semantics to describe data
3. enable **simple user workflows** in modelling, supporting **direct** questions like: What is the social dynamics of water in basin X? How does switching to crop Y affect rural food security in region Z?

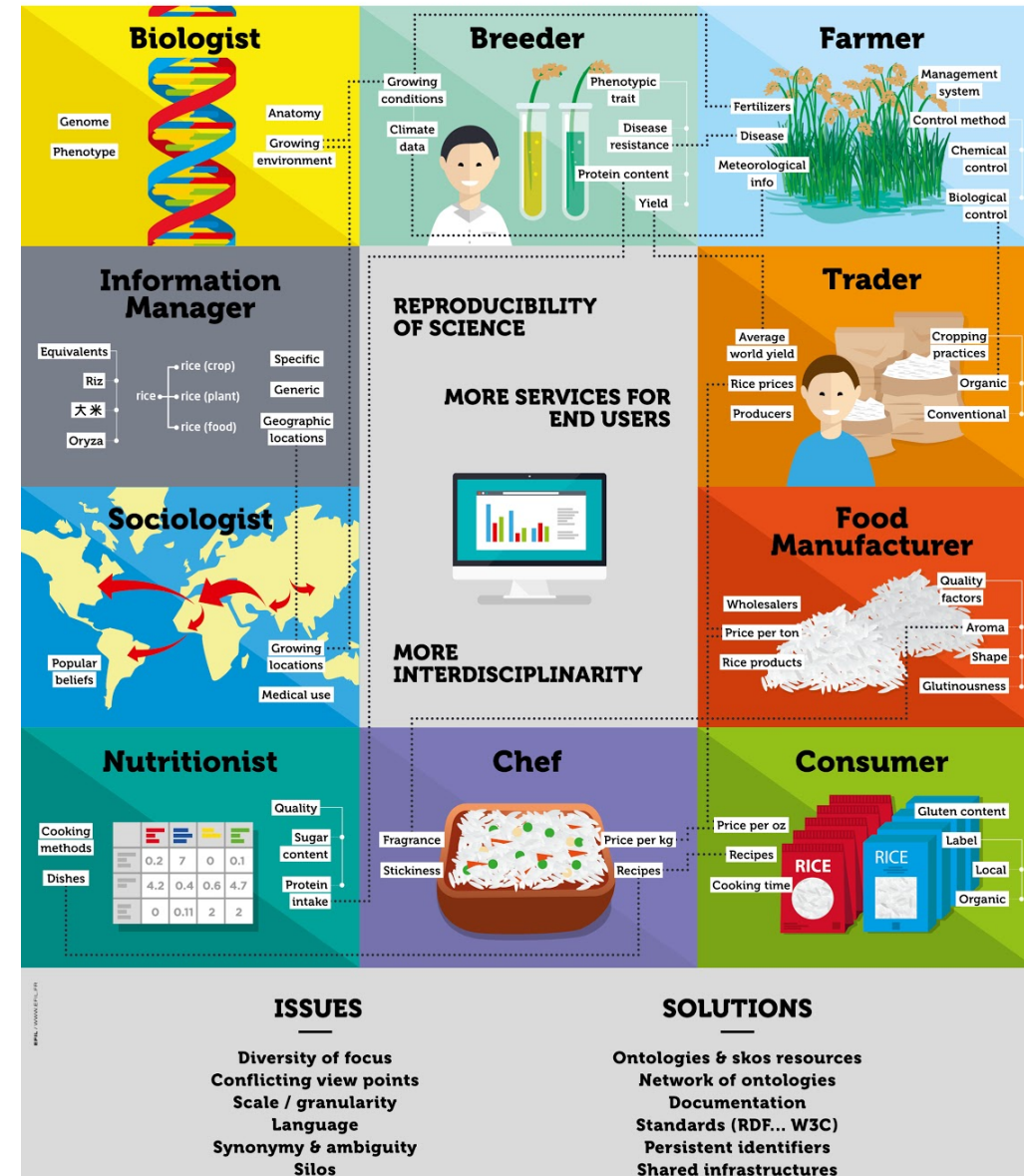
Where are we along this path in 2017?

Using and reusing data: The state of the art

1. Distributed access to datasets over the web (OGC, OpenDAP, ...)
2. Linked Open Data paradigm: open standards, each artifact is coupled with a URI pointing to its “meaning”.
3. Problem: the meaning *differs for each observer* - unless semantics is coherent across domains, uses and goals.
4. If it’s not consistent, it’s not FAIR

Image credits: INRA, AgriSemantics RDA working group

SEMANTICS - THE WAY TO RECONCILE POINTS OF VIEW AND DATA THE EXAMPLE OF "RICE"

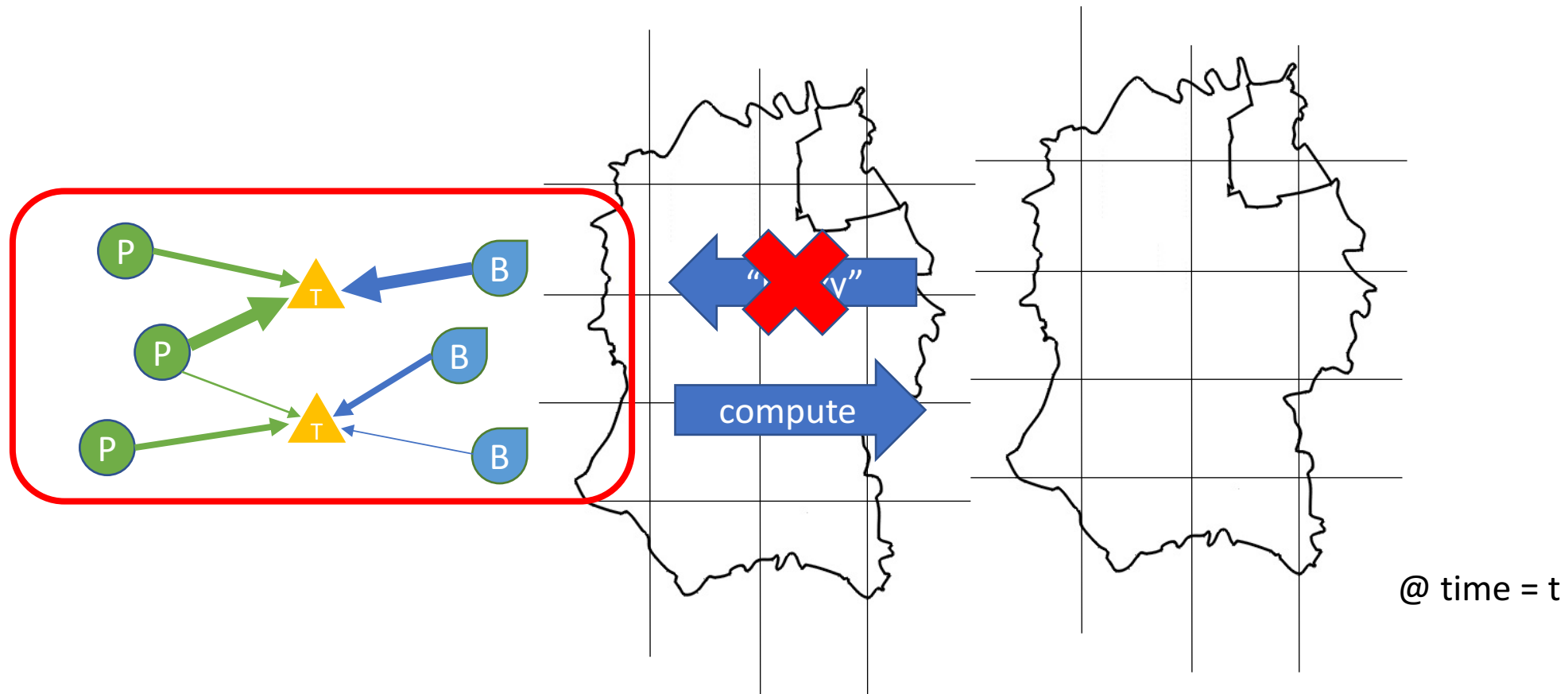


Reusing models

- Modeling paradigms represent different “metaphors” adopted during model design:
 - process-based vs. agent-based
 - stochastic/probabilistic vs. deterministic models
 - spatial vs. non-spatial, raster/vector, continuous vs. discrete time, etc.
- It remains **difficult to mix and match models incarnating different paradigms** across the lifecycle of an application.
- Often, complex problems are handled with one paradigm that fits some components but must be “tricked” to handle the rest.
- As a result models are still brittle **monoliths**, hard to disassemble and reassemble.
- Integrating architectures (OpenMI &C.) only handle the technical aspects of integration, addressing only a subset of the problem.

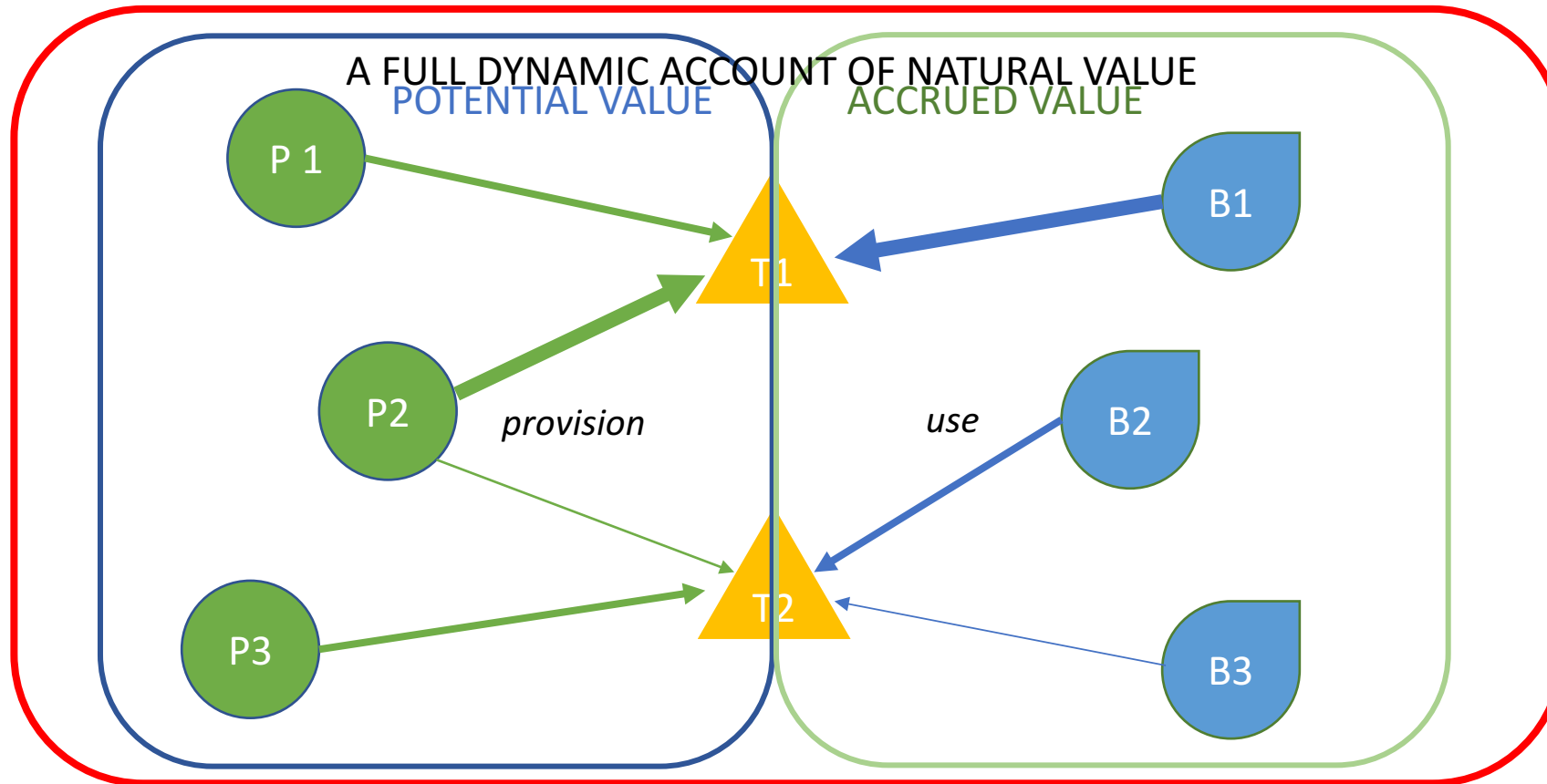
A case in point: accounting for human-natural interactions

- We know the limitations of “proxy” models – and it’s not because of decision makers.
- Still, building models of the *true* system models is hard - impossible in rapid assessments



Adaptive, assisted system characterization

Driven by semantics and by *roles*, supporting a specific view of physical phenomena without introducing ambiguities



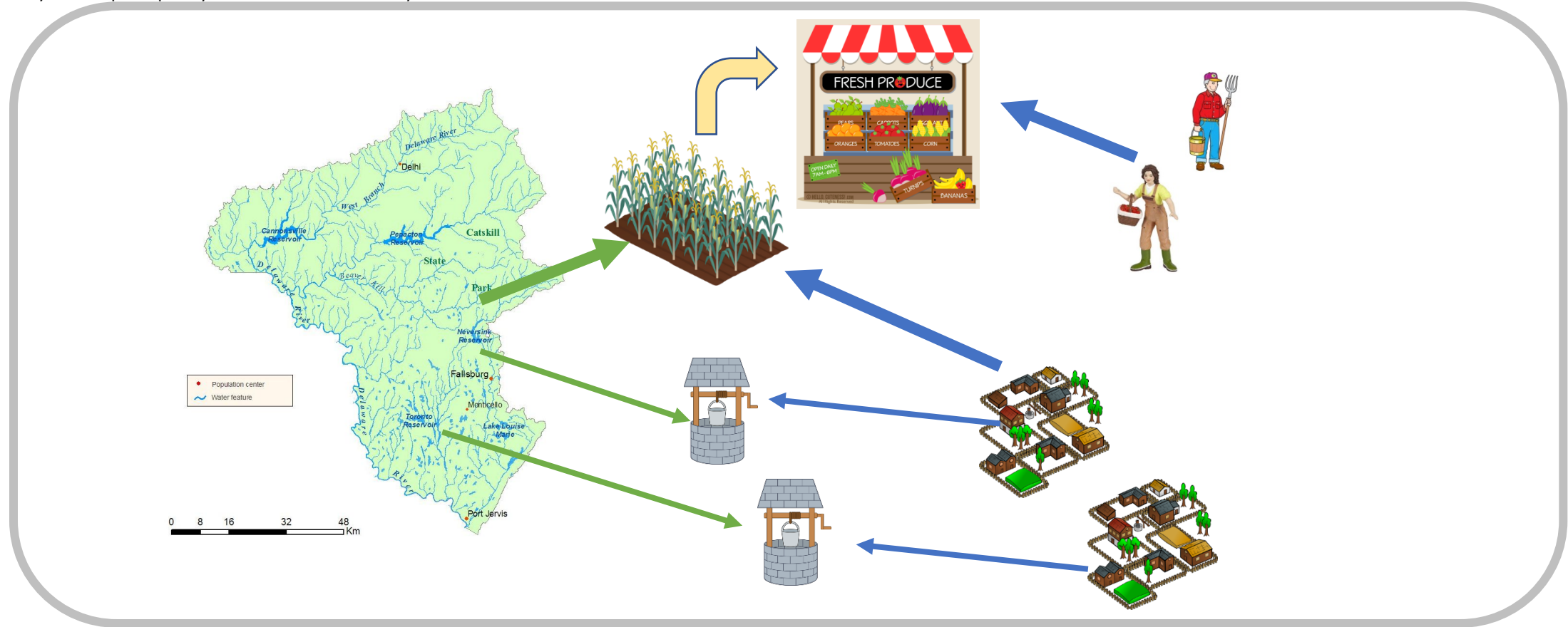
Providers (e.g. forests, watersheds): where valuable ecosystem function happens

Transactors (e.g. wells, crops, atmosphere): where natural value is generated

Beneficiaries (e.g. farmers, coastal dwellers): demand agents for natural value

Example: building an eco-social flow network

Triggered by a simple query: "observe social dynamics of water in watershed X"



The model for the system (e.g. forests and watersheds) are first identified and built by the AI engine. Providers (e.g. forests and watersheds) are first identified and built by the AI engine. This ontology defines types of Transactors (e.g. wells, crops, and spheres), starting with provision (provider->transactor)... and following with use (beneficiary <- transactor), building a (potentially) differently scaled model for each flow. Beneficiaries (e.g. farmers, coastal dwellers) are identified last. Intermediate transactors (e.g. markets) are brought in according to the ontologies. They can be local or remote.

Resolution of models based on semantics

Model statements are stored in a distributed database. Each dependency is stated conceptually and resolved contextually.

