# Exposé: An ontology for machine learning experimentation

Joaquin Vanschoren, K.U.Leuven (Belgium), U. Leiden (The Netherlands)
Larisa Soldatova, University of Aberystwyth (UK)

**DM Ontology Jamboree 2010**

# Exposé: An ontology for machine learning experimentation

Joaquin Vanschoren, K.U.Leuven (Belgium), U. Leiden (The Netherlands)
Larisa Soldatova, University of Aberystwyth (UK)

**DM Ontology Jamboree 2010**

# Overview

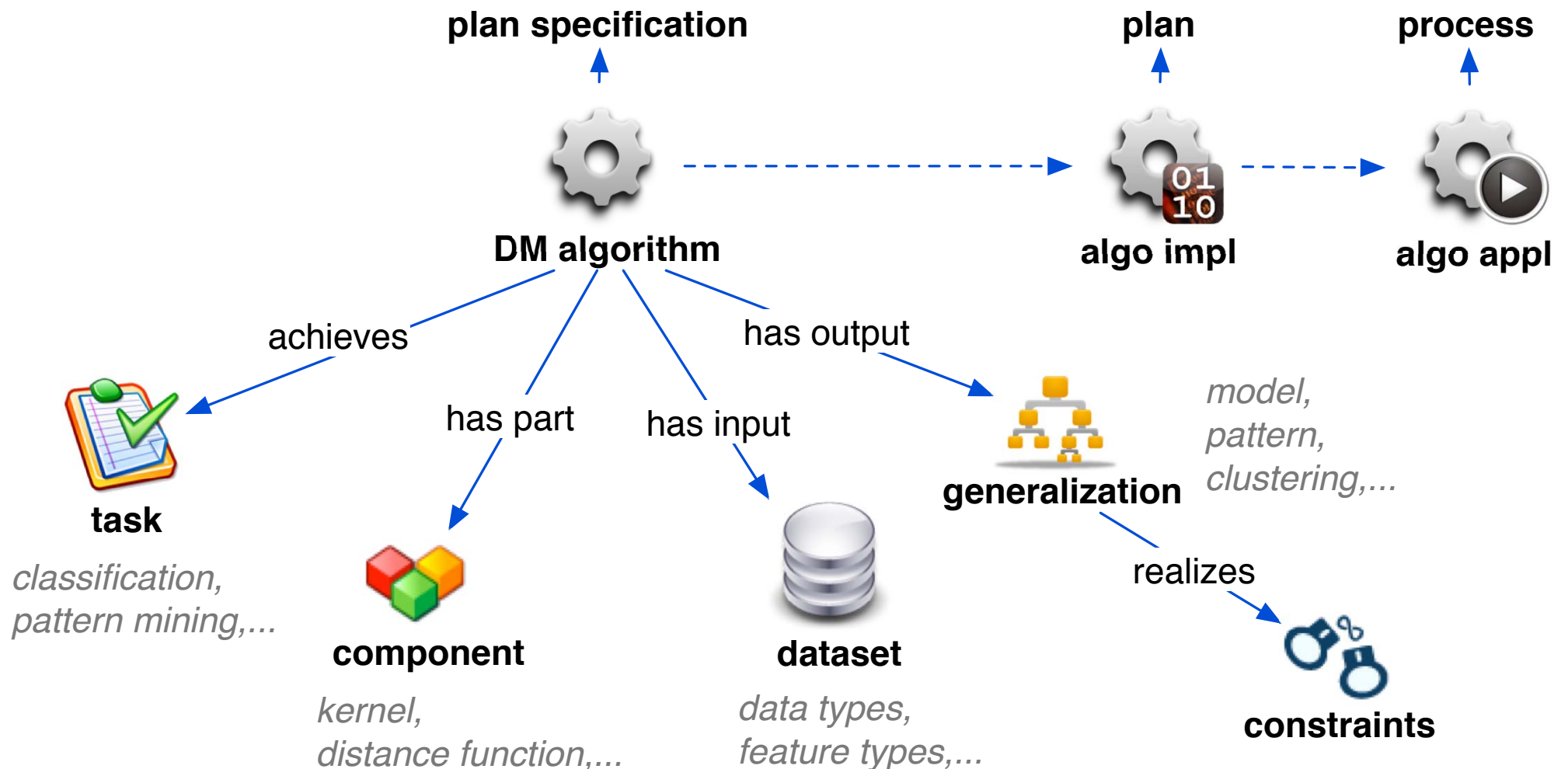Ontology lessons
Exposé ontology
Use cases

# Ontology lessons

What did we learn from other ontologies

# Ontology design

- Start from accepted classes & properties (top-level ontologies, e.g. OBI, RO)

- If possible, reuse prior ontologies to build on their knowledge/consensus

- Use ontology design patterns: reusable patterns for recurrent problems
  - http://ontologydesignpatterns.org

- Check clarity, consistency, extensibility, minimal commitment

# Ontology recap:
# OntoDM (Panov et al., '09,'10)

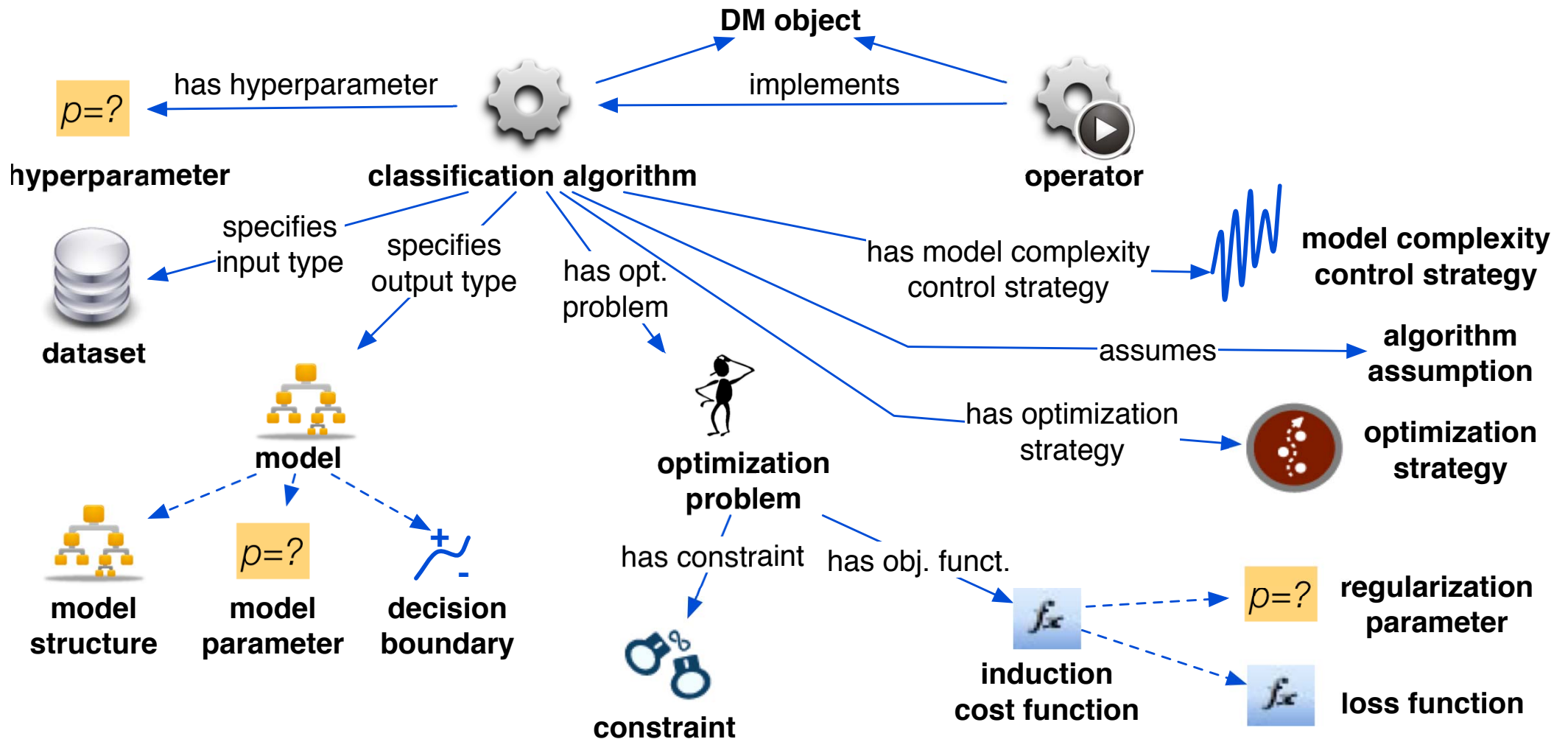- Aim: unified framework for DM research, builds on BFO



**plan specification**

**plan**

**process**

**DM algorithm**

**algo impl**

**algo appl**

achieves

has part

has input

has output

**task**

*classification, pattern mining,...*

**component**

*kernel, distance function,...*

**dataset**

*data types, feature types,...*

**generalization**

*model, pattern, clustering,...*

realizes

**constraints**

# Ontology recap:
# OntoDM (Panov et al., '09,'10)

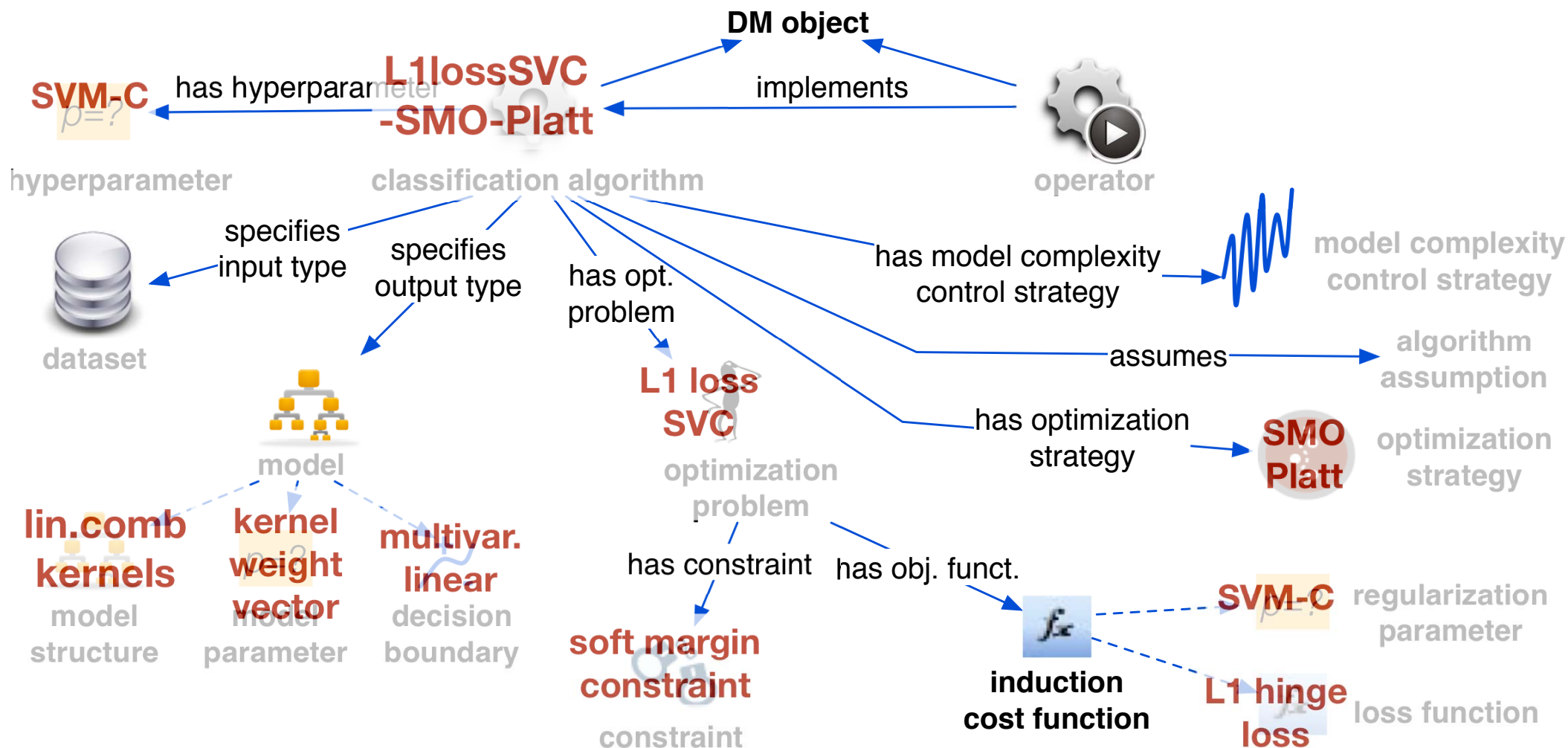- Aim: unified framework for DM research, builds on BFO

# Ontology recap:
# DMOP (Hilario et al., '09)
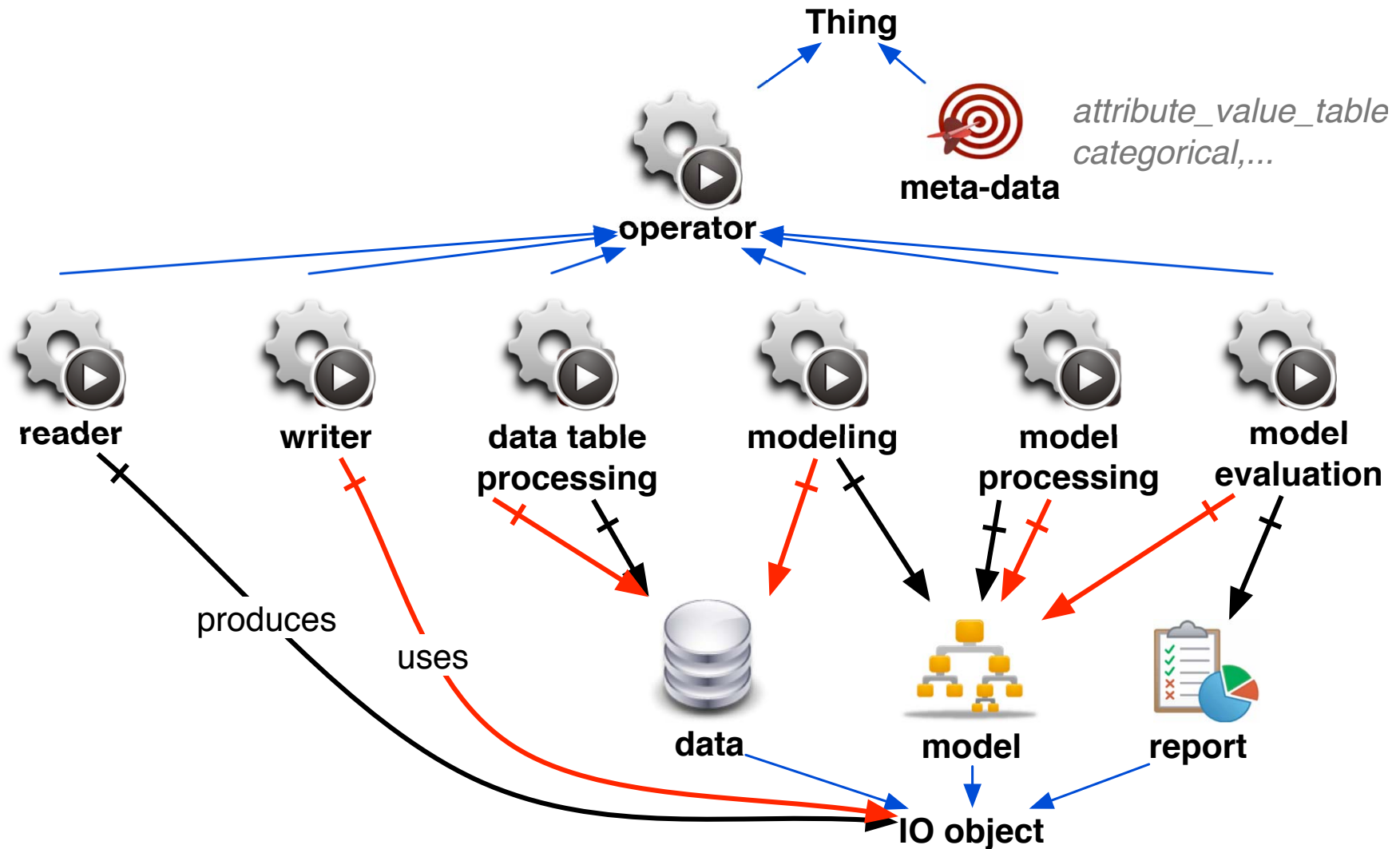
- Model internal structure of learning algorithms

# Ontology recap:
# DMOP (Hilario et al., '09)

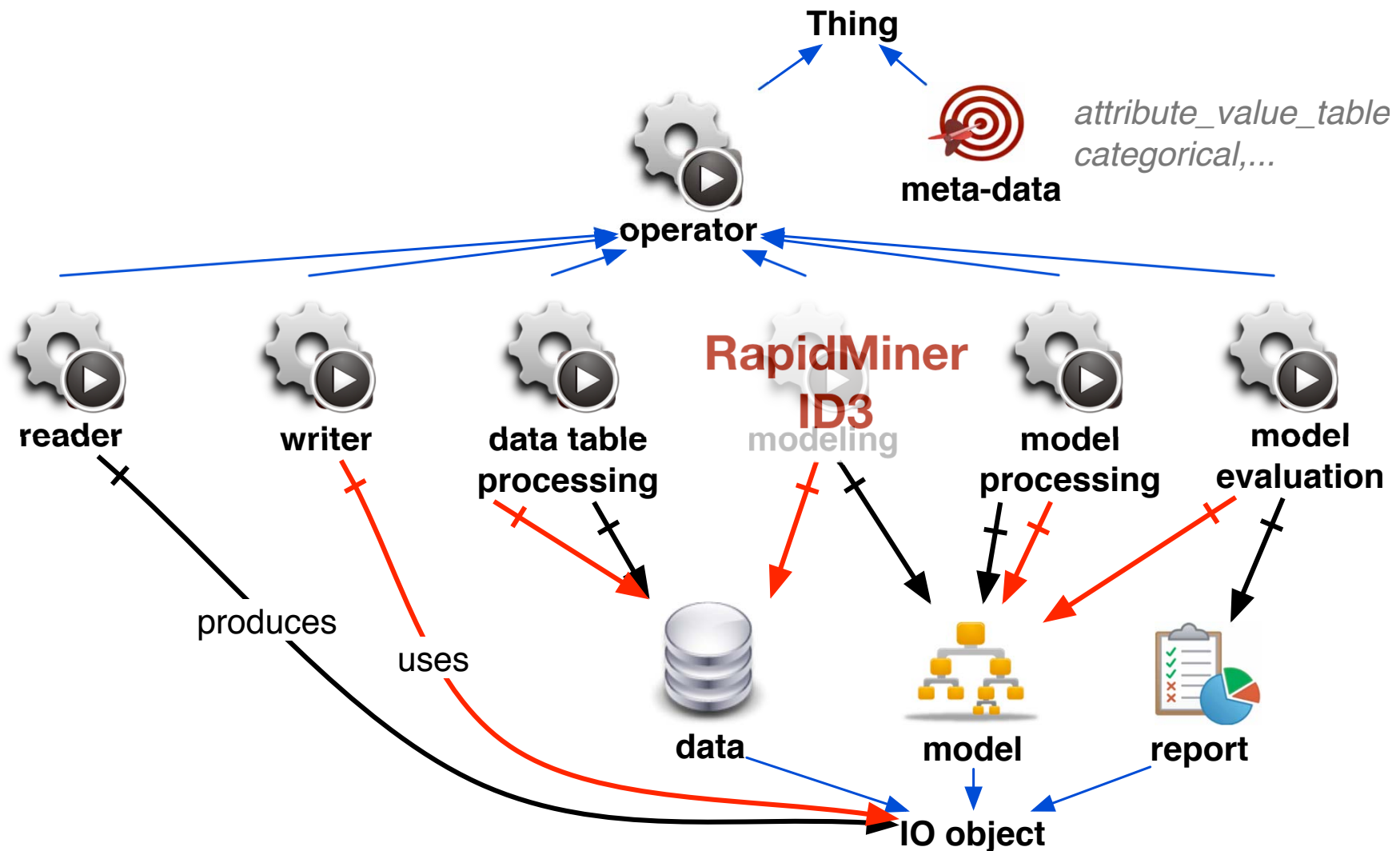- Model internal structure of learning algorithms

# Ontology recap:
# DMWF (Kietz et al., '09)

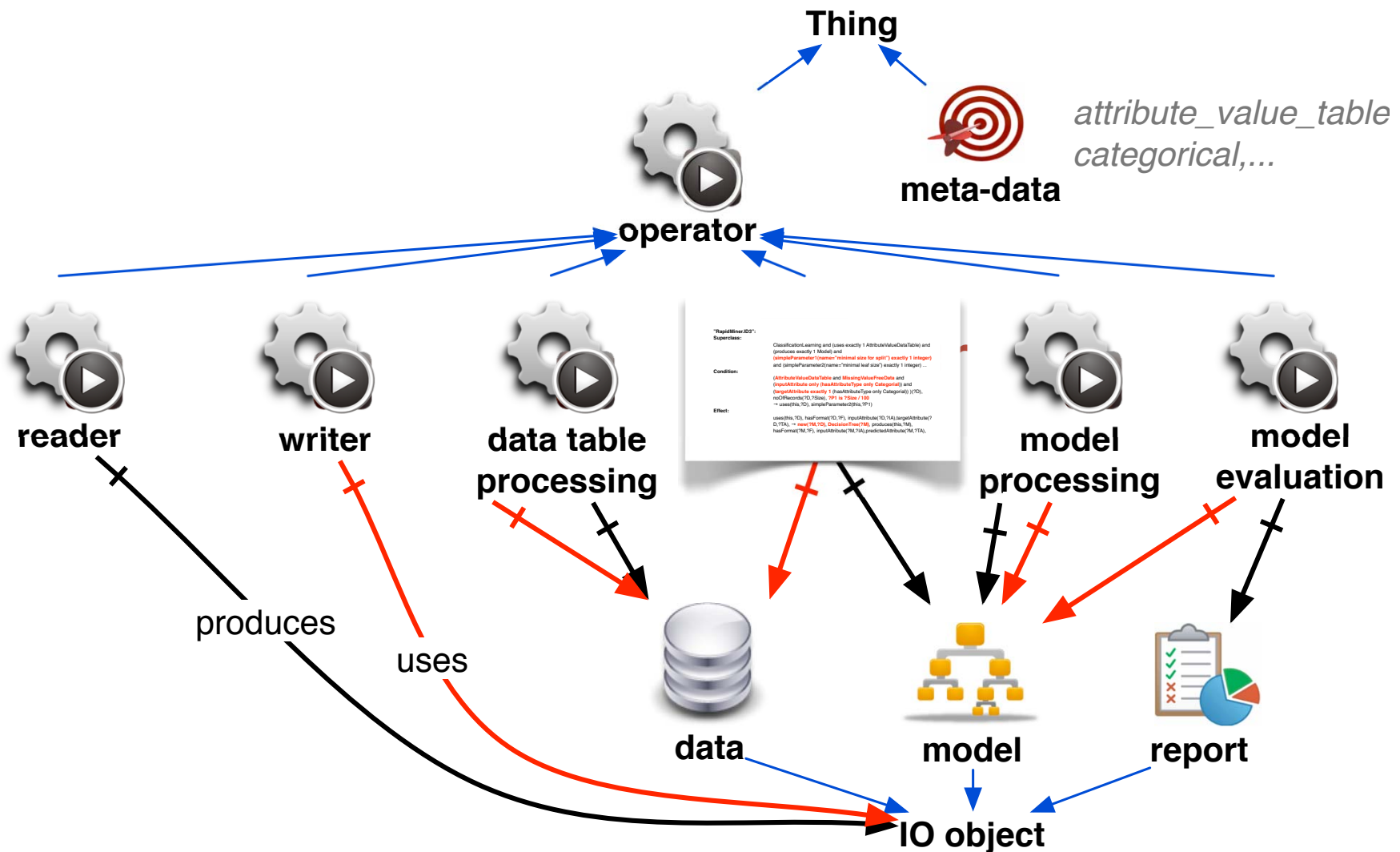- Reason about KD operators: in/outputs, conditions/effects (SWRL rules)

# Ontology recap:
# DMWF (Kietz et al., '09)

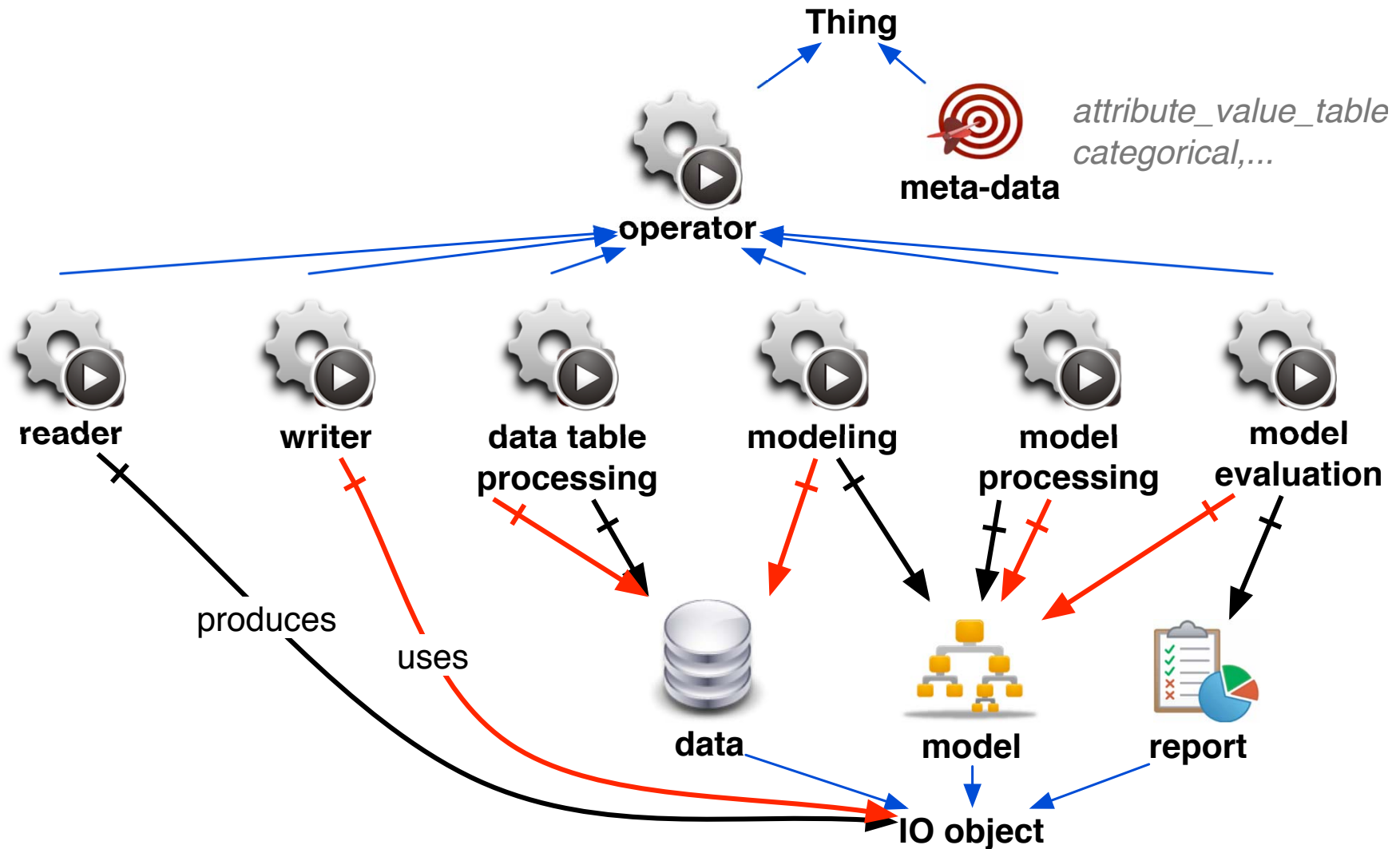- Reason about KD operators: in/outputs, conditions/effects (SWRL rules)

# Ontology recap:
# DMWF (Kietz et al., '09)

- Reason about KD operators: in/outputs, conditions/effects (SWRL rules)

# Ontology recap:
# DMWF (Kietz et al., '09)

- Reason about KD operators: in/outputs, conditions/effects (SWRL rules)

# Ontology recap:
# DMWF (Kietz et al., '09)

- Reason about KD operators: in/outputs, conditions/effects (SWRL rules)

**"RapidMiner.ID3":**
**Superclass:**
    ClassificationLearning and (uses exactly 1 AttributeValueDataTable) and
    (produces exactly 1 Model) and
    **(simpleParameter1(name="minimal size for split") exactly 1 integer)** and
    (simpleParameter2(name="minimal leaf size") exactly 1 integer) ...
**Condition:**
    (**AttributeValueDataTable** and **MissingValueFreeData** and
    (**inputAttribute only (hasAttributeType only Categorial**)) and
    (**targetAttribute exactly 1** (hasAttributeType only Categorial)) )(?D),
    noOfRecords(?D,?Size), **?P1 is ?Size / 100**
    → uses(this,?D), simpleParameter2(this,?P1)
**Effect:**
    uses(this,?D), hasFormat(?D,?F), inputAttribute(?D,?IA),targetAttribute(?D,?TA),
    → **new(?M,?D), DecisionTree(?M)**, produces(this,?M), hasFormat(?M,?F),
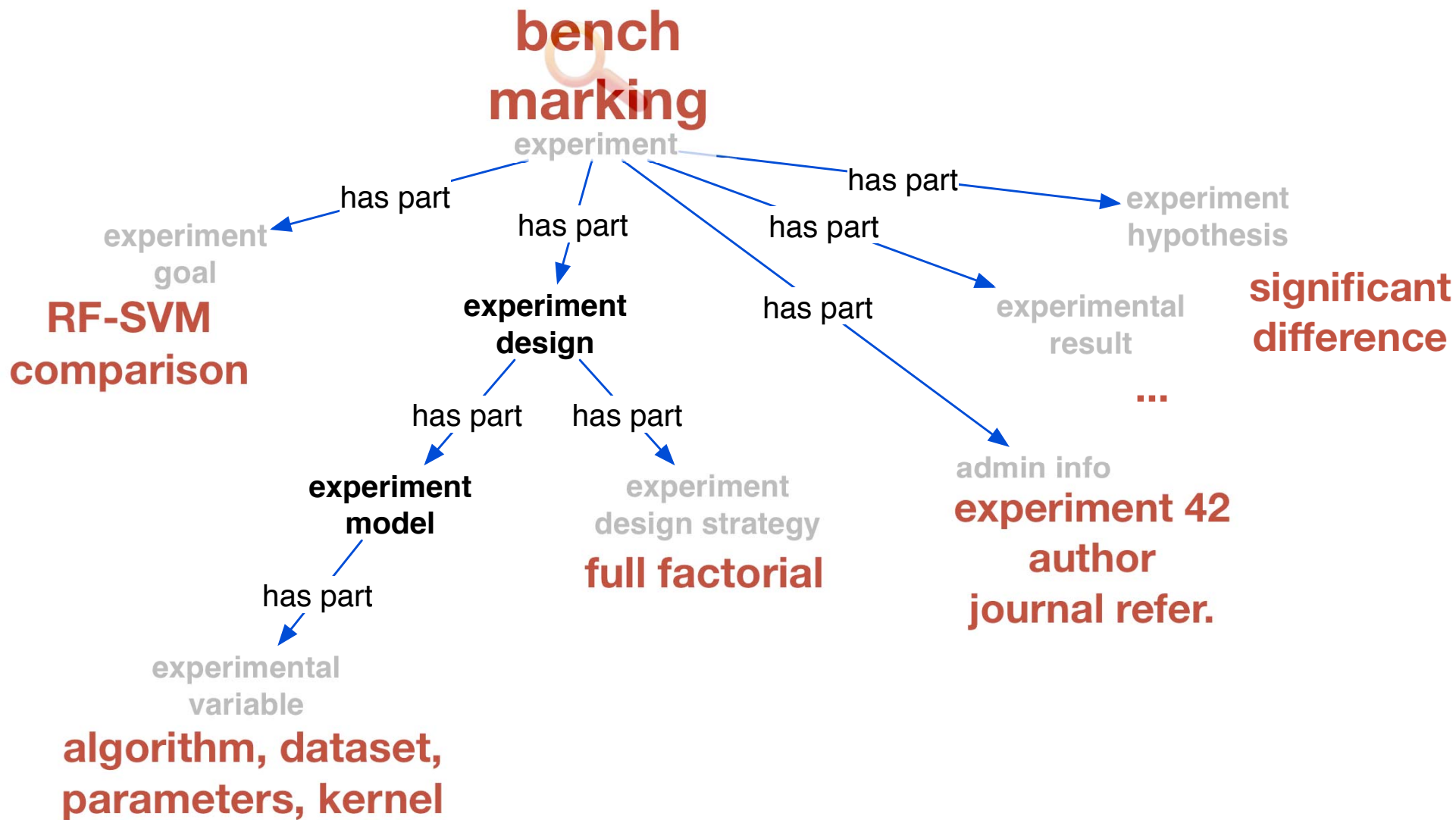    inputAttribute(?M,?IA),predictedAttribute(?M,?TA),

# Ontology recap:
# EXPO (Soldatova and King, '06)

- Make goal and structure of scientific experiments more explicit

# Ontology recap:
# EXPO (Soldatova and King, '06)

- Make goal and structure of scientific experiments more explicit

# Exposé

an ontology for data mining experimentation

# Context

- Giant, public database(s) of data mining experiments

- We need:

  - Common language to share experiments (through DM tools)

  - Intuitive ways to store and query experimental results

- We want:

  - Interoperable ontology: OntoDM for top-level, DMOP for detailed properties of learning algorithms

  - Driven by actual experiments submitted to database

    - New algorithms -> ideally, described by author

    - Instances automatically extracted from database

# Problem 1: Experiments

What is a machine learning experiment?
What do we need to know about it?

# Exposé: Experiments

hp: has participant
hd: has description



**experiment workflow**

**KD workflow**

**Workflow:
has inputs, outputs,
operators (participants)**

# Exposé: Experiments

hp: has participant
hd: has description

**composite experiment** — hp → **experimental design** — hp → **experimental variable**

**composite experiment** — hp → **singular experiment**

**singular experiment** — is exec. on → **machine**

**EXPO**

**experiment workflow**

**KD workflow**

**Workflow:**
**has inputs, outputs, operators (participants)**

# Exposé: Experiments

hp: has participant
hd: has description

EXPO

composite experiment —hp→ experimental design —hp→ experimental variable

composite experiment —hp→ singular experiment

singular experiment —is exec. on→ machine

machine ←hd— model evaluation result

experiment workflow

KD workflow

**Workflow:**
**has inputs, outputs,**
**operators (participants)**

performance estimation

learner evaluation —hp→ performance estimation

learner evaluation —has output→ model evaluation result

model evaluation result —has part→ evaluation

learner evaluation —hp→ model evaluation function

learner evaluation —hp→ learning algorithm

learner evaluation —has input→ dataset

model evaluation function —has input→ model

model —hd→ model evaluation result

learning algorithm —has output→ model

model —has output→ prediction result

dataset

learning algorithm

model evaluation function

model

prediction result

# Problem 2: Algorithms

When talking about an algorithm, what is meant?

General algorithm?
Specific implementation? Which version?
When run, which parameters, components?

# Exposé: Algorithms
*Specification, implementation, application*

hp: has participant
hd: has description
ico: is concretization of

Similar to OntoDM



algorithm specif.

ico

algo impl

hd — name, version, url,...

has part → p=? param impl

has quality → algo quality

hp

algo appl

hp → operator

function appl

hp → p=? param setting

# Same for functions and parameters



ico = is concretization of
hp = has participant

# Problem 3:
# Algorithm composition

plug-in functions, kernels, other algorithms

such components play different roles
-> Agent-role pattern

# Exposé: Algorithms

hp: has participant
hd: has description

# Exposé: Algorithms

hp: has participant
hd: has description

# Exposé: Algorithms

hp: has participant
hd: has description



kernelized algorithm

learning algorithm

algorithm specif.

baselearner

search

algorithm role

data preprocessor

algorithm component role

role

function role

kernel

distance functions

ico

name, version, url,...

hd

algo impl

has part

p=?

param impl

has quality

hp

hp

algo quality

operator

algo appl

realizes

function appl

hp

p=?

param setting

# Problem 4: Workflows

Inputs, outputs, operators
Hierarchical: workflows within workflows
Reuse, parameterize common workflows, e.g. k-fold CV

# Exposé: workflows

# Problem 5: Reuse

How can we make maximal use of existing ontologies?

OBI: top-level
OntoDM: top-level DM concepts
DMO: operators, learning mechanisms

# BFO: accepted top-level classes

# BFO: accepted top-level classes

# OntoDM: top-level DM concepts



ico = is concretization of
 hp = has participant

# DMO: operators, learning mechanisms



ico = is concretization of
hp = has participant

# Exposé: top level classes



**BFO**

thing

quality — realizable entity — material entity — planned process — digital entity — information content entity

objective

data quality - algorithm quality

role

machine

KD workflow — *hp* → operator

implemen-tation

dataset — *has descr.* → dataset spec

model — *ico* → model spec

**OntoDM**

*executed on*

function appl. — *hp* → function impl. — *ico* → function

algo appl — *hp* → algo impl — *ico* → algorithm specif.

*hp*

*has part*

*has part*

$p=?$ param setting — *hp* → $p=?$ param impl — *ico* → $p=?$ parameter

**DMOP**

$p=?$

ico = is concretization of
hp = has participant

25

# Other aspects

# Datasets



has input

has output

**dataset**

**data processing
appl**

hp

**data processing
workflow**

**KD
workflow**

# Datasets

evaluation
function appl.

hp

hp

learner
evaluation

evaluation
function impl.

ico

evaluation
function

# Evaluation

evaluation function appl.

evaluation function impl.

evaluation function

learner evaluation

hp

hp

ico

has description

name

version

association evaluation measure

support

confidence

leverage

conviction

frequency

lift

clustering evaluation measure

density-based clustering measure

integrated squared error

distance-based clustering measure

inter-cluster similarity

intra-cluster variance

probabilistic distribution evaluation measure

integrated average squared error

probability distribution scoring function

distribution likelihood

distribution log-likelihood

computational evaluation measure

probabilistic model distance measure

Kullback-Leibner divergence

likelihood ratio

build cpu time

build memory consumption

predictive model evaluation measure

class prediction evaluation measure

AUPRC

AUROC

single point AUROC

f- measure

derived measure

binary prediction evaluation function

recall

precision

specificity

confusion matrix

has participant

averaged binary prediction measure

multi-class prediction evaluation measure

probability error-based measure

numeric prediction evaluation measure

has participant

kappa statistic

predictive accuracy

class RMSE

correlation coefficient

error-based evaluation measure

information criterion

graphical evaluation measure

ROC_curve

cost curve

lift chart

RMSE

MAD

RSS

RRSE

MAPE

AIC

BIC

precision-recall curve

PRgraph point

has part

has part

# Experiment context

**singular experiment**

hp

**composite experiment**

# Experiment context

# Exposé: final notes

- In total 860 classes, 32 properties (from RO + DMOP)

- Individuals: all algorithms, preprocessors, evaluation from WEKA

  - actually stored in experiment database

  - should be programmatically added (and updated)

- Written in OWL-DL, using Protégé 4.0

- Can be browsed at:

  - http://expdb.cs.kuleuven.be/expdb/expose.owl

  - http://www.e-lico.eu/OWLBrowser2/manage/

# Use Cases

# Goal: Collaborative experimentation
# Now: small-scale, not repeatable, not reusable

# Goal: Collaborative experimentation
# Now: small-scale, not repeatable, not reusable

---

⚙️ !

**new algorithm**

# Goal: Collaborative experimentation
# Now: small-scale, not repeatable, not reusable

datasets

!

# Goal: Collaborative experimentation
# Now: small-scale, not repeatable, not reusable

preprocessing
workflows

# Goal: Collaborative experimentation
# Now: small-scale, not repeatable, not reusable



evaluation
procedures

# Goal: Collaborative experimentation
# Now: small-scale, not repeatable, not reusable

# Goal: Collaborative experimentation
# Now: small-scale, not repeatable, not reusable

# Goal: Collaborative experimentation
# Now: small-scale, not repeatable, not reusable

# Goal: Collaborative experimentation
# Now: small-scale, not repeatable, not reusable

# Goal: Collaborative experimentation
# Now: small-scale, not repeatable, not reusable

- A lot of work, limits depth

- Results cannot be reused by others (have to be repeated)

- Hard to repeat experiments from descriptions in papers!

# Goal: Collaborative experimentation
# Now: small-scale, not repeatable, not reusable



- A lot of work, limits depth
- Results cannot be reused by others (have to be repeated)

...iments from ...rs!

*The journal system is perhaps the most open system for the transmission of knowledge that could be built ... with 17th century media.* Nielsen (APS Physics 2008)

# Data mining as an e-science
# Ontologies: experiments shared, run automatically

# Data mining as an e-science
# Ontologies: experiments shared, run automatically

- Share experiments
  - Internet = large, collaborative workspace

# Data mining as an e-science
# Ontologies: experiments shared, run automatically

- Share experiments
  - Internet = large, collaborative workspace

- Store them in ***experiment databases***
  - Ensure reproducibility
  - Reuse millions of prior experiments
  - Use all info on algorithms, datasets
  - Results universally accessible and useful

# e-Sciences
# Astrophysics: Virtual Observatories

# e-Sciences
# Bio-informatics: Micro-array Databases

# e-Sciences
## Bio-informatics: Micro-array Databases

# Collaborative Experimentation
# Why?

# Collaborative Experimentation
# Why?

**Reproducibility**

**Good science**

# Collaborative Experimentation
# Why?

**Reproducibility**

**Good science**

**Quick, easy analysis**

**Querying: Answer questions**

**Test hypotheses**

# Collaborative Experimentation
# Why?

**Reproducibility**

**Good science**

**Quick, easy analysis**

**Querying: Answer questions**

**Test hypotheses**

**Reuse**

**Save time & energy
(e.g. benchmarking)**

# Collaborative Experimentation
# Why?



**Reproducibility**
    **Good science**

**Quick, easy analysis**
    **Querying: Answer questions**
    **Test hypotheses**

**Reuse**
    **Save time & energy**
    **(e.g. benchmarking)**

**Generalizability:**
    **Plug into prior results: larger studies**

# Collaborative Experimentation
# Why?

**Reproducibility**

**Good science**

**Quick, easy analysis**

**Querying: Answer questions**

**Test hypotheses**

**Integration**

**Data mining tools**

**import/export**

**Reuse**

**Save time & energy**

**(e.g. benchmarking)**

**Generalizability:**

**Plug into prior results: larger studies**

# Collaborative Experimentation
# Why?

**Reproducibility**
**Good science**

**Quick, easy analysis**
**Querying: Answer questions**
**Test hypotheses**

**Integration**
**Data mining tools**
**import/export**

**Reuse**
**Save time & energy**
**(e.g. benchmarking)**

**Reference**
**'Map' of known approaches**
**Compare to state-of-the-art**
**Includes negative results**

**Generalizability:**
**Plug into prior results: larger studies**

# Use Case 1

Describe experiments in a common language
-> sharing or running experiments on grid

# Use Exposé to define common language: ExpML

# Use Exposé to define common language: ExpML



in → learner evaluation

# Use Exposé to define common language: ExpML



in → learner evaluation → out

# Use Exposé to define common language: ExpML



in → learner evaluation → out

has participant

performance estimation

evaluation function

# Use Exposé to define common language: ExpML



appl

has participant

in → learner evaluation → out

has participant

performance estimation

evaluation function

38

# Use Exposé to define common language: ExpML



parameter setting

$p=?$

operator (component)

has participant

appl

has participant

impl

has participant

in

learner evaluation

out

has participant

performance estimation

evaluation function

# ExpML: a markup language for DM experiments



- Share DM experiments, XML-based

dataset

learner evaluation

 appl

 impl

 *p=?* param.sett.

 operator

 perform. estim. appl.

 appl

 eval. function appl.

 appl

 model evaluation

# ExpML: a markup language for DM experiments

└ [database icon] dataset

└ ( learner evaluation )

    └ [gear icon] appl

        └ [gear icon] impl

        └ *p=?* param.sett.

        └ [gear icon] operator

    └ [gear icon] perform. estim. appl.

    └ [fx icon] eval. function appl.

        └ [gear icon] appl

        └ [fx icon] appl

└ [bar chart icon] model evaluation

- Share DM experiments, XML-based

```
<expml>
    <dataset id='d1'>
    <learner evaluation id='e1' input_data='d1'>
        <learner_appl>
            <learner_impl name=... version=...>
            <parameter_setting name='P' value='100'/>
            <learner_appl role= 'base-learner'>
                ...
        </learner_appl>
        <performance_estimation_appl>
        ...
        <model_evaluation_function_appl>
        ...
    </learner_evaluation>
    <model_evaluation_result output_of='e1'>
        <evaluation name='accuracy' value= '0.99'>
        ...
```
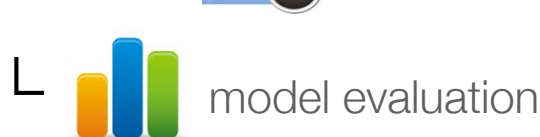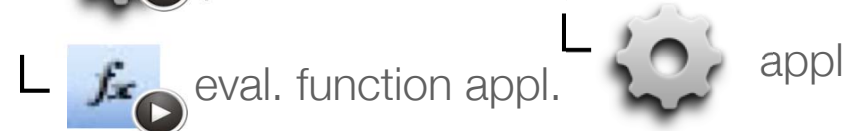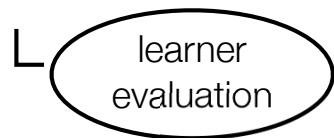
# ExpML: a markup language for DM experiments



L  dataset

L  learner evaluation

L  appl

| | |
|---|---|
| ontology | XML |
| has-part, has-participant | XML subelement |
| | (with role attribute) |
| has-description | (required) attribute |
| has-quality | `property' subelement |
| is-concretization-of | implementation_of attr. |
| part-of | attributes |
| has-specific-input | input_data attribute |
| has-specified-output | output_of attribute |

- Share DM experiments, XML-based

```
<expml>
    <dataset id='d1'>
    <learner evaluation id='e1' input_data='d1'>
        <learner_appl>
            <learner_impl name=... version=...>
            <parameter_setting name='P' value='100'/>
            <learner_appl role= 'base-learner'>
                ...
        </learner_appl>
        <performance_estimation_appl>
        ...
        <model_evaluation_function_appl>
        ...
    </learner_evaluation>
    <model_evaluation_result output_of='e1'>
        <evaluation name='accuracy' value= '0.99'>
        ...
```
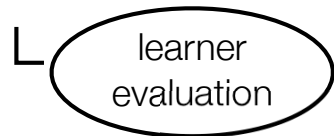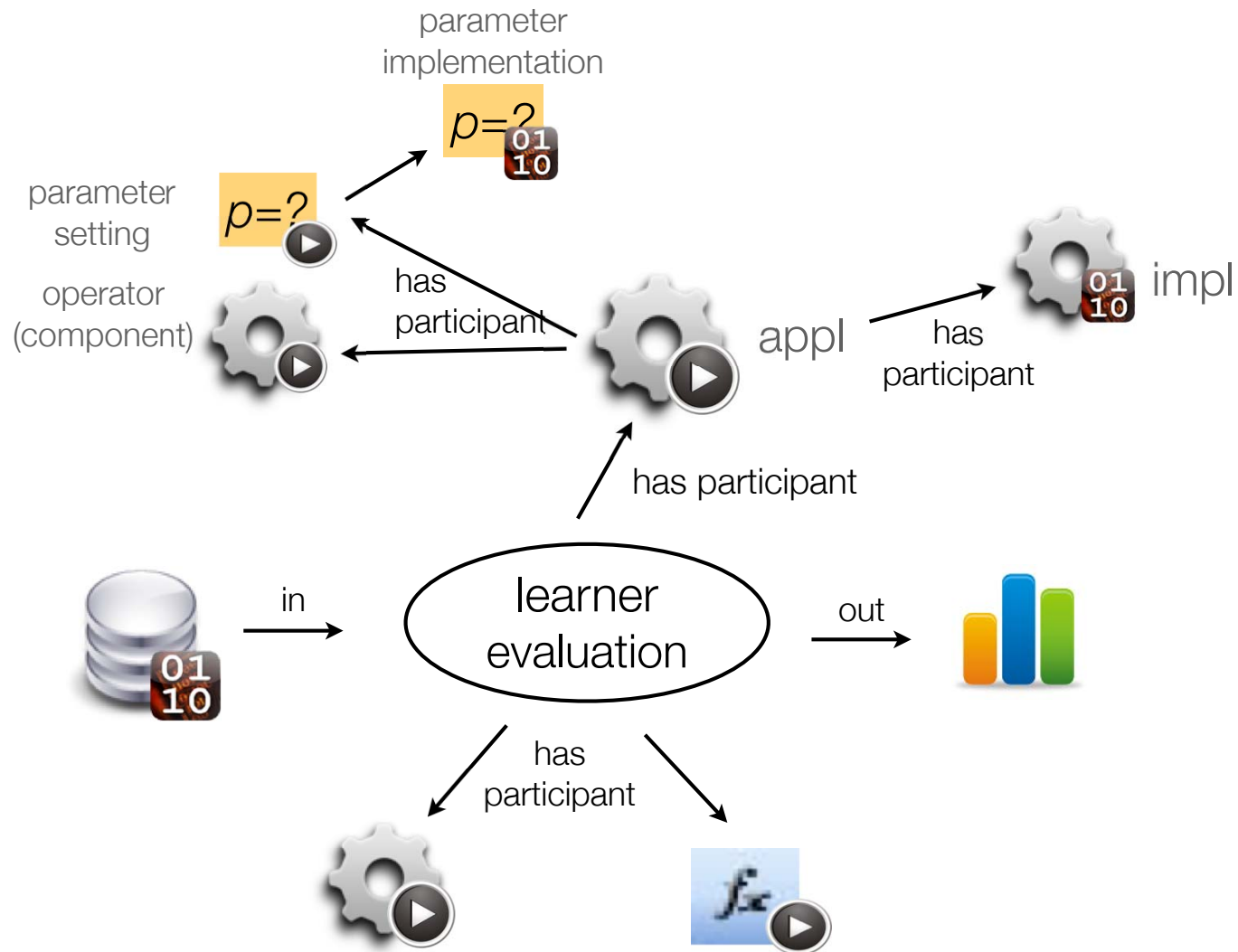
# Use Case 2

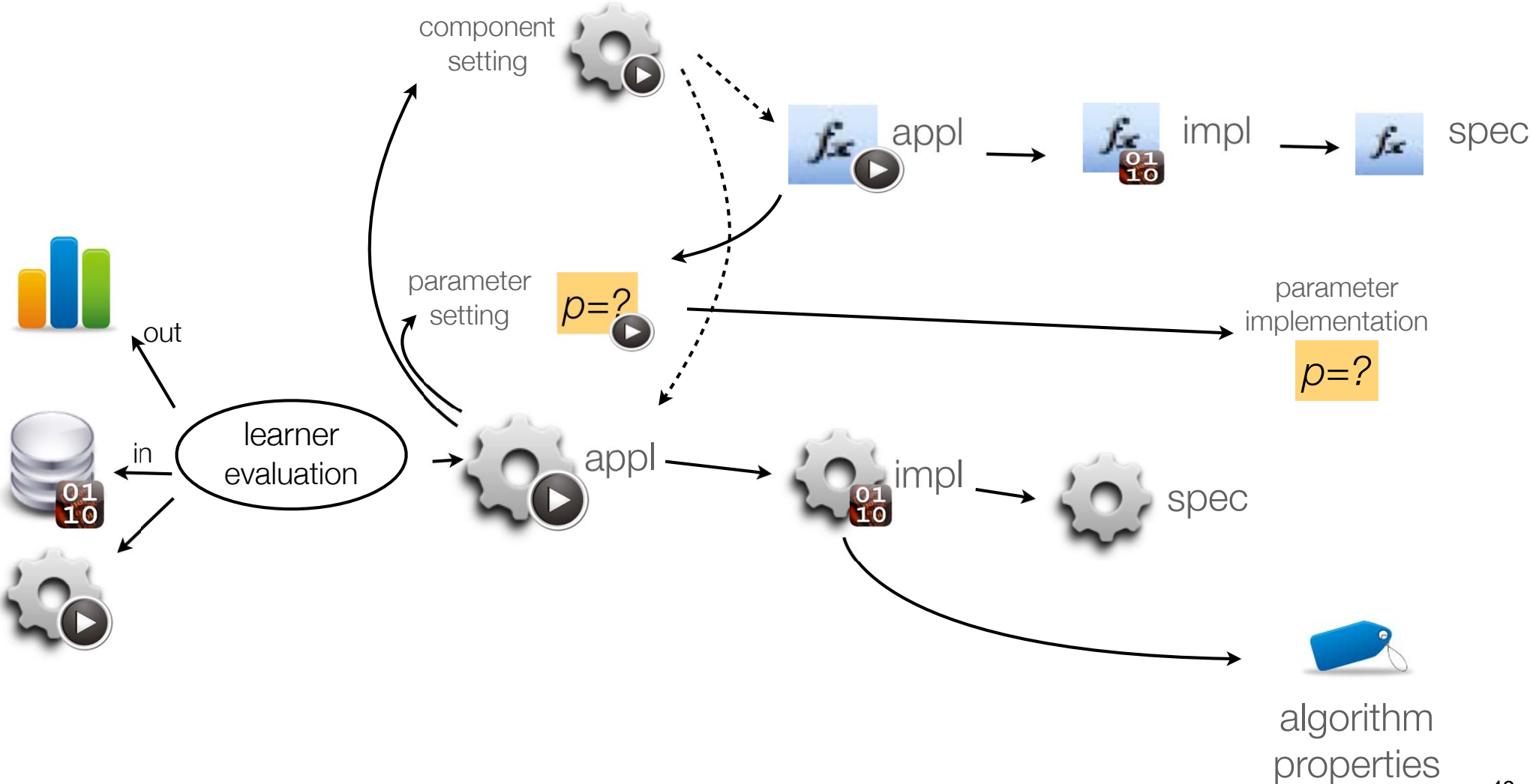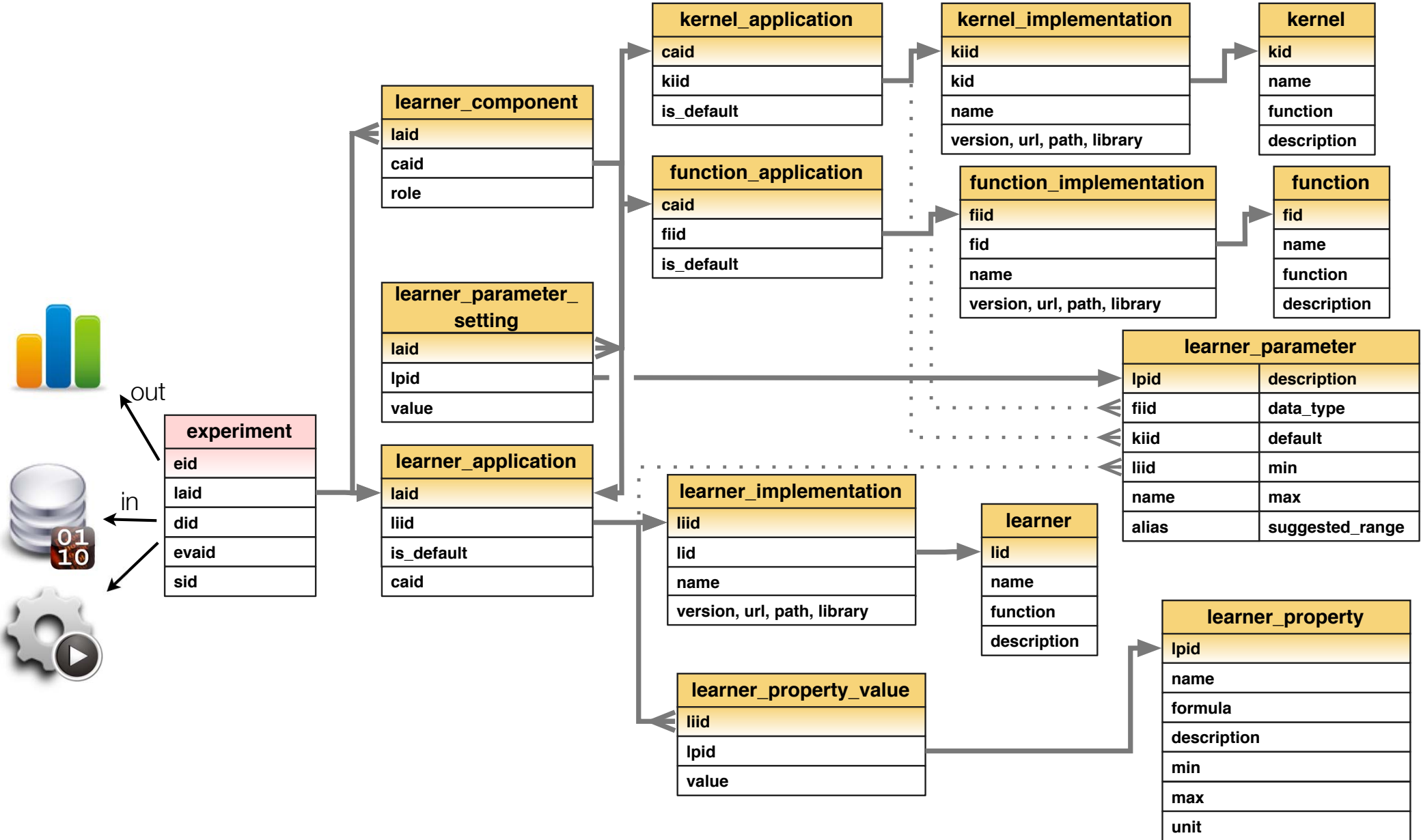Collect experiments in a database
to query all empirical results

# ExpDB: a database to share experiments

# Experiment Database

>650,000 experiments, 54 algorithms,
>87 datasets, 45 evaluation measures,
2 data processors, bias-variance analysis



component setting

appl

impl

spec

parameter setting

$p=?$

parameter implementation

$p=?$

out

in

learner evaluation

appl

impl

spec

algorithm properties

# Experiment Database

>650,000 experiments, 54 algorithms,
>87 datasets, 45 evaluation measures,
2 data processors, bias-variance analysis

# Use Case 3

Intuitive querying

# Query Interface (YouTube "experiment database")

http://expdb.cs.kuleuven.be

# The way ahead

- 3rd generation of tools could make data mining into e-science

  - Experiments shared, reused, run worldwide

  - Repeatable, generalizable, reusable

- Cooperation on a standardized ontology for data mining?

- Automatic ontology extraction: DM paper -> ontology extension

- RDF experiment databases?

- Open problems:

  - Queriable models, auto-population (active meta-learning), quality control

Hvala

Danke

Thanks

Xie Xie

Diolch

Toda

Merci

Grazie

Spasiba

Efharisto

Obrigado

Arigato

Köszönöm

Tesekkurler

Dank U

Dhanyavaad

Gracias

http://expdb.cs.kuleuven.be