# Dr.Aid: an automated formal framework to support data-governance rule compliance for decentralized collaboration

Rui Zhao

Supervised by: Petros Papapanagiotou, Malcolm Atkinson, Jacques Fleuriot

University of Edinburgh

2021-05-24

# Rules – heterogeneity, similarity and redundancy

I agree to restrict my use of CORDEX model output for non-commercial research and educational purposes only. [1]

In publications that rely on the CORDEX model output, I will appropriately credit the data providers by an acknowledgement similar to the following: "We acknowledge..." [1]

You may extract, download, and make copies of the data contained in the Datasets, and you may share that data with third parties according to these terms of use. [2]

When sharing or facilitating access to the Datasets, you agree to include the same acknowledgment requirement in any sub-licenses of the data that you grant, and a requirement that any sub-licensees do the same. [2]

Data is non-transferrable (other than as permitted in the licence) and confidential in nature. [3]

Data is not to be used to identify, contact or target patients or general

# Problem to solve

- Large collection of ... data
  - personal
  - sensitive
  - valuable

# Problem to solve

- Large collection of ... data
  - personal
  - sensitive
  - valuable
- Issues of ...
  - data privacy
  - ownership
  - governance

# Problem to solve

- Large collection of ... data
    - personal
    - sensitive
    - valuable
- Issues of ...
    - data privacy
    - ownership
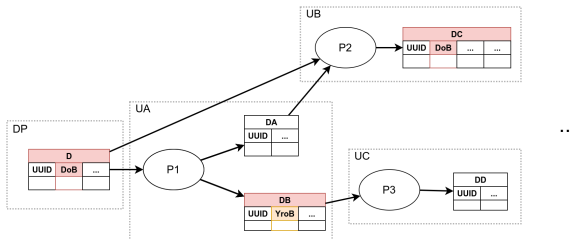    - governance

*Compliance checking* remains ...
- manual
- time consuming
- error prone

# Example use case

**Data governance rules associated with dataset D by data provider DP**

Users should properly acknowledge the data provider in their publications, in a form similar to "This work ...".
Field DoB (Date of Birth) in the data is potentially *sensitive* and any use of it should be reported to the data provider, through the URL report.example.ac.
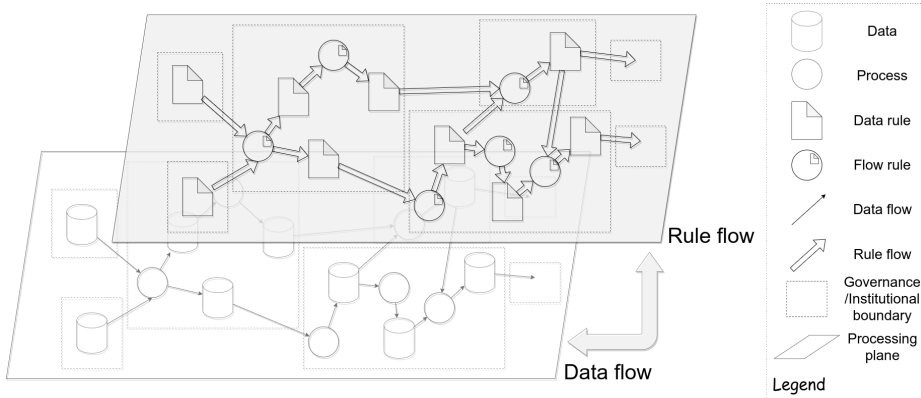
# Issues identified

1. (Personnel) **Scattering**: data processing is *multi-institutional* (providers don't work together with users);
2. (Rule) **Propagation**: derived data can be used as input data further;
3. (Rule) **Diversity**: policies can be more than access control, e.g. *obligations*;
4. **Dynamic** (rule) **application**: processes can change the policies applied to data;
5. (Rule) **Combination** and **separation**: processes are multi-input-multi-output (MIMO).

# Related research

| Framework | Scattering | Propagation | Diversity | Dynamic application | Combination | Separation |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| E-P3P[38] | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ |
| Thoth[29] | ✗ | ✓ | ✗ | ✓[2] | ✓ | ✗ |
| DAPRECO[26, 57] | ? | ✗ | ✓ | ✗ | ✗ | ✗ |
| Smart object[59] | ✓ | ✓ | ✓ | ✗ | ✓ | ✗ |
| CamFlow[53] | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ |
| Meta-code[36] | ✓[3] | ✓ | ✓[4] | ✓ | ✗ | ✗ |
| Dr.Aid (our work) | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

# From **data flow** to **rule flow**

# What is Dr.Aid

Dr.Aid (Data Rule Aid) is a *logic-based* framework for *automated compliance checking* of data governance rules over data flow *graphs*.

# What is Dr.Aid

Dr.Aid (Data Rule Aid) is a *logic-based* framework for *automated compliance checking* of data governance rules over data flow *graphs*.

Highlights:
- Logic-based (situation calculus)
- All 5 issues addressed, in particular:
  - Recognises rule diversity
  - Addresses dynamic rule application
  - Supports MIMO

# Model

- Language:
  - **data rule**, for data-governance rules for multi-staged processing, e.g. "*users must properly acknowledge the data providers*"
  - **flow rule**, for the *changes* of data rules in each *process* as a result of data transportation and transformation, e.g. "*column 3 from input 1 is changed to column 2 on output 2*"
- Supports *obligation*: user must do something after using the data
- Retrospective analysis from provenance
  - CWLProv (file-oriented)
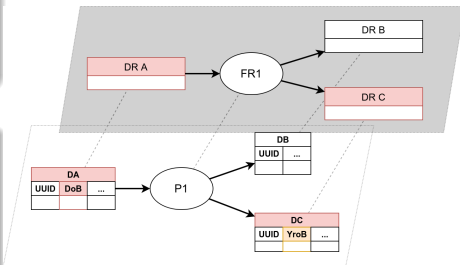  - S-Prov (data-streaming)

# Language – example

## Data-use policy (as *data rule*)

Field DoB (Date of Birth) in the data is potentially *sensitive* and any use of it should be reported to the data provider.

## Process information (as *flow rule*)

- changes column DoB from input 1 to column YroB on output 2;
- removes column DoB from input 1 on output 1;
- propagates the rest of the data unaffected regarding the data-use policy.

# Language – example modelled

This dataset contains potentially sensitive information in column DoB.
Any use of it should be reported to the data providers through link report.example.ac.

Natural-language policy

attribute(pf, column "DoB")
    Name    Type    Value
attribute(ru, url "report.example.ac")
    Name Type        Value
obligation(report ru, [pf], action = *)
    Obligated  Action    Validity    Activation
    action   argument  binding     condition

User notation

attribute(pf, column, "DoB", [input1, pf_1], s0)
                                History (ID)   Situation
attribute(ru, url, "report.example.ac", [input1, ru_1], s0)
                                        History (ID)    Situation
obligation(report, [[input1, ru_1]],
            [[input1, pf_1]], "action = *", input1, s0)
                                            Port    Situation

Situation calculus representation

obligation(report attribute(ru, url "report.example.ac"), [attribute(pf, column "DoB")], action = *)

Nested equivalent view

# Language – example modelled and query

Querying is analogous to the projection task, i.e. querying predicates (fluents) that hold at the targeted final situation:

$$S_f = Do(pr(input1, [output1, output2]) :$$
$$delete(input1, output1, *, column, "DoB") :$$
$$edit(input1, output2, *, column, "DoB", column, "YroB") :$$
$$end([output1, output2]), s0).$$
$$? - attribute(N, T, V, X, S_f).$$

## Language – example modelled and query

Querying is analogous to the projection task, i.e. querying predicates (fluents) that hold at the targeted final situation:

$$S_f = Do(pr(input1, [output1, output2]) :$$
$$delete(input1, output1, *, column, "DoB") :$$
$$edit(input1, output2, *, column, "DoB", column, "YroB") :$$
$$end([output1, output2]), s0).$$
$$? - attribute(N, T, V, X, S_f).$$

resulting in

$$attribute(ru, url, "report@example.ac", [output1, input1, pf_1], S_f)$$
$$attribute(pf, column, "YroB", [output2, input1, pf_1], S_f)$$
$$attribute(ru, url, "report@example.ac", [output2, input1, pf_1], S_f)$$

... and the same procedure for *obligations*.

# Language – Data rule

- *attributes* ($N$, $T$, $V$), describing properties of the data
    - a name $N$
    - a type $T$
    - a value $V$
- *obligations* ($OD$, $VB$, $AC$)
    - an obligation definition $OD$ (the obligated action to perform upon activation)
    - a validity binding $VB$ (describing additional applicability constraints)
    - an activation condition $AC$ (the triggering condition)

# Language – flow rule

- **propagate** $pr(P_{in}, Ps_{out})$
- **edit** $edit(P_{in}, P_{out}, N, T, V, T_{new}, V_{new})$
- **delete** $delete(P_{in}, P_{out}, N, T, V)$

# Language – flow rule

- **propagate** $pr(P_{in}, Ps_{out})$
- **edit** $edit(P_{in}, P_{out}, N, T, V, T_{new}, V_{new})$
- **delete** $delete(P_{in}, P_{out}, N, T, V)$

$$pr(input1, [output1, output2])$$
$$delete(input1, output1, *, column, DoB)$$
$$edit(input1, output2, *, column, DoB, column, YroB)$$

$*$ represents *anything*, used as wildcard.

# Situation calculus formalisation

Situation calculus models things into three components:

- *fluent*: predicate with *situation*, holding different truth values at different *situation*s
- *situation*: state of the world, where the initial situation specifies the initial state of the world
- *action*: behaviour to change the world, applying to a *situation* and resulting in a new *situation*

# Situation calculus formalisation

Situation calculus models things into three components:

- *fluent*: predicate with *situation*, holding different truth values at different *situation*s
- *situation*: state of the world, where the initial situation specifies the initial state of the world
- *action*: behaviour to change the world, applying to a *situation* and resulting in a new *situation*

We model:

- data rule terms as *fluents*
- the initial data rules are specified for the initial *situation*
- the flow rules are modelled as *actions*

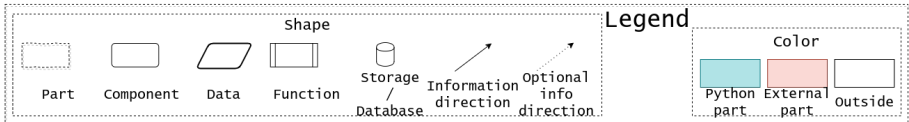# Situation calculus formalisation
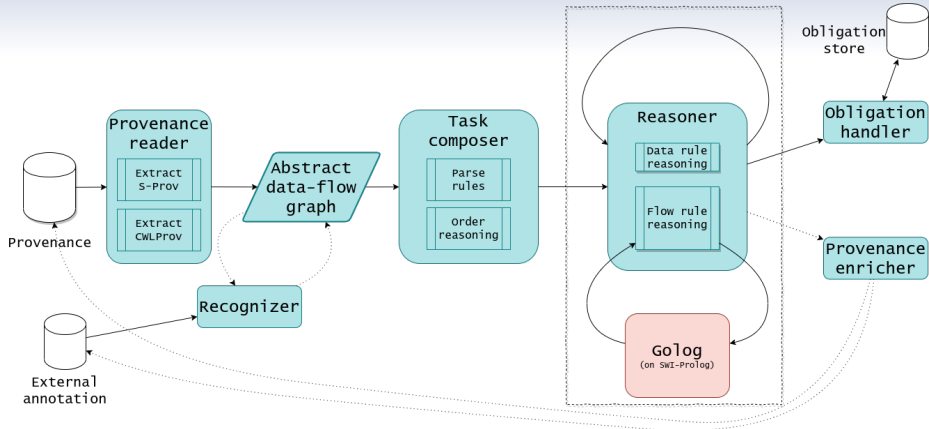
We model:

- data rule terms as *fluents*
- the initial data rules are specified for the initial *situation*
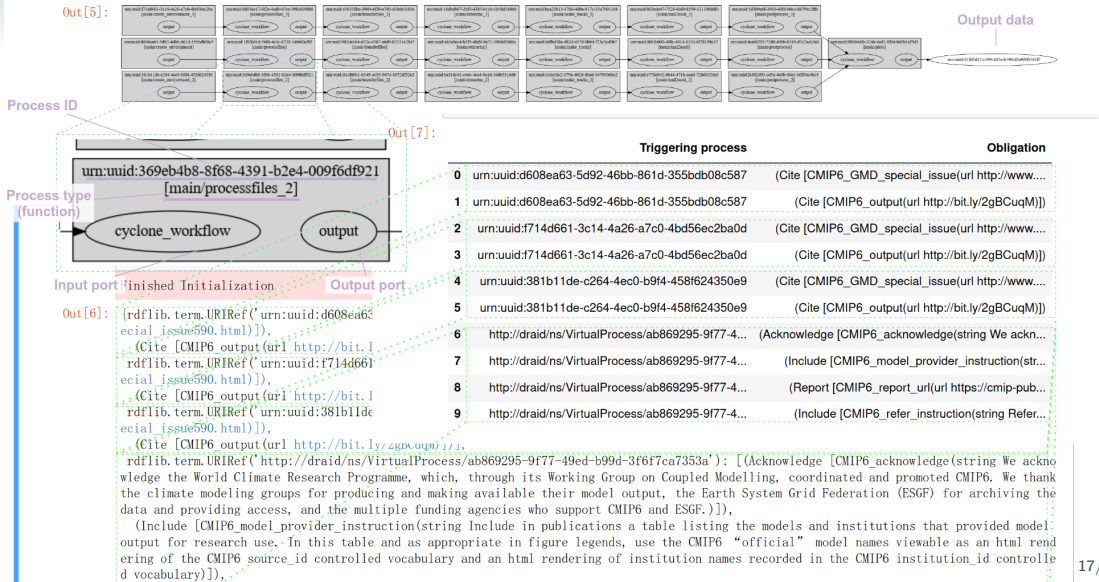- the flow rules are modelled as *actions*

We do:

- flow rules as action sequence
- final situation $S_f$ as the situation after performing all flow rules starting from the initial situation $s_0$
- query all *attribute*() and *obligation*() which hold in the final situation $S_f$

# System structure

# Outcomes from system

# Language coverage

| Policy source | # sentences | # rules | # actioning | # implied | # encoded | actioning coverage | total coverage |
|---|---|---|---|---|---|---|---|
| CMIP6[3] | 35 | 9 | 8 | 1 | 7 | 100% | 89% |
| EIDA[8] | 20 | 5 | 3 | 0 | 3 | 100% | 60% |
| INGV[7] | 2 | 2 | 2 | 0 | 2 | 100% | 100% |
| CC-BY[5] | 12 | 6 | 5 | 2 | 3 | 100% | 83% |
| CMT Catalogue[9] | 15 | 4 | 4 | 0 | 4 | 100% | 100% |
| CORDEX[4] | 22 | 9 | 6 | 0 | 5 | 83% | 55% |
| ISMD[12] | 2 | 1 | 1 | 0 | 1 | 100% | 100% |
| RCMT[10] | 14 | 3 | 3 | 2 | 1 | 100% | 100% |
| MIMIC[13] | 17 | 4 | 4 | 0 | 4 | 100% | 100% |
| CPRD[6] | 21 | 7 | 6 | 0 | 2 | 33% | 29% |
| PIMA[15] | 2 | 1 | 1 | 0 | 1 | 100% | 100% |
| ISC[2] | 21 | 7 | 7 | 0 | 7 | 100% | 100% |
| IRIS[11] | 28 | 10 | 10 | 0 | 10 | 100% | 100% |
| OGL[14] | 30 | 7 | 4 | 3 | 1 | 100% | 57% |
| World Bank[16] | 40 | 12 | 7 | 2 | 3 | 71% | 42% |
| **Total** | 281 | 87 | 71 | 10 | 54 | 90% | 74% |

# Thanks
# for listening!