# Strengths and deficits of CEN Workshop Agreements for data documentation in materials modelling and characterization

**Björn Schembera (UStuttgart)**
**Martin T Horsch (NMBU)**
**Heinz A Preisig (NTNU)**

# Going beyond FAIR

Björn Schembera

Institute for Applied Analysis and Numerical Simulation

University of Stuttgart

# Dark Data - Definition

- Dark Data[1] which is
  - Hidden
  - Unavailable
  - Unstructured
  - Undocumented/unannotated
  - Biased
  - Stemming from abandoned research
- Up to 80% of the global data is dark[2]
- Why should we care?
  - Economical: it costs
  - Ecological: "..annual global [..] footprints resulting from storing dark data might approach 5.26 million tons CO2 [..]."[3]
  - Responsibility problem
  - Epistemic problem (epistemic opaqueness)
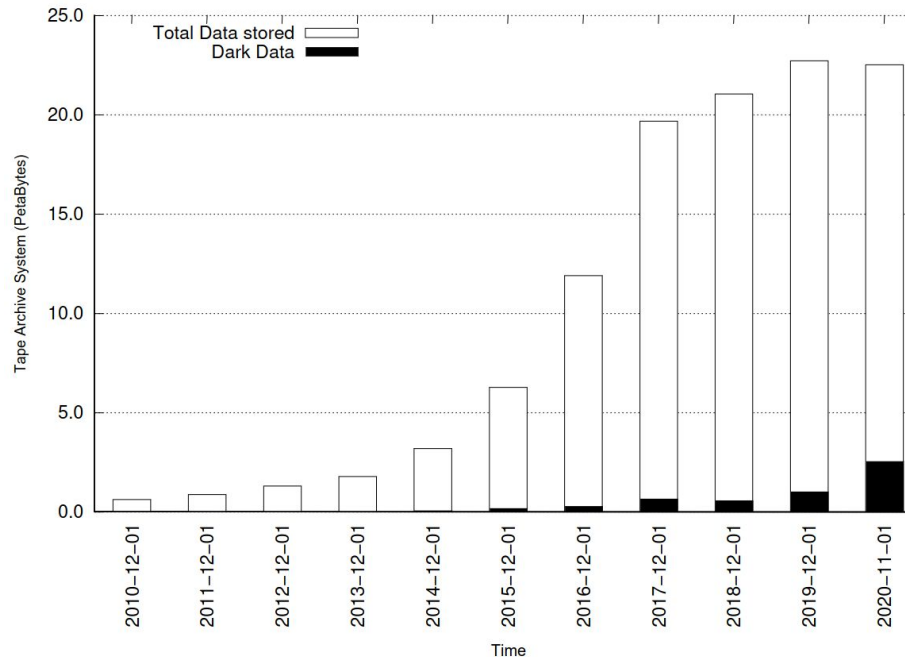  - Legal Implications

[1] Schembera, Björn, and Martin T. Horsch. Dark data and epistemic metadata in molecular modeling. In preparation
[2] Ahmad, Norita, Areeba Hamid, and Vian Ahmed. "Data Science: Hype and Reality."Computer 55.2 (2022): 95-101.
[3] Al Kez, Dlzar, et al. Exploring the sustainability challenges facing digitalization and internet data centers. J. of Cleaner Production 371 (2022).

# Dark Data - Example 1: HPC Center

- Dark data at an HPC Center[4], > 11% dark data by 2020
  - Lots of data is dark due to orphaned accounts or missing metadata



[4] Schembera, Björn, and Juan M. Durán. "Dark data as the new challenge for big data science and the introduction of the scientific data officer."Philosophy & Technology 33.1 (2020): 93-115.
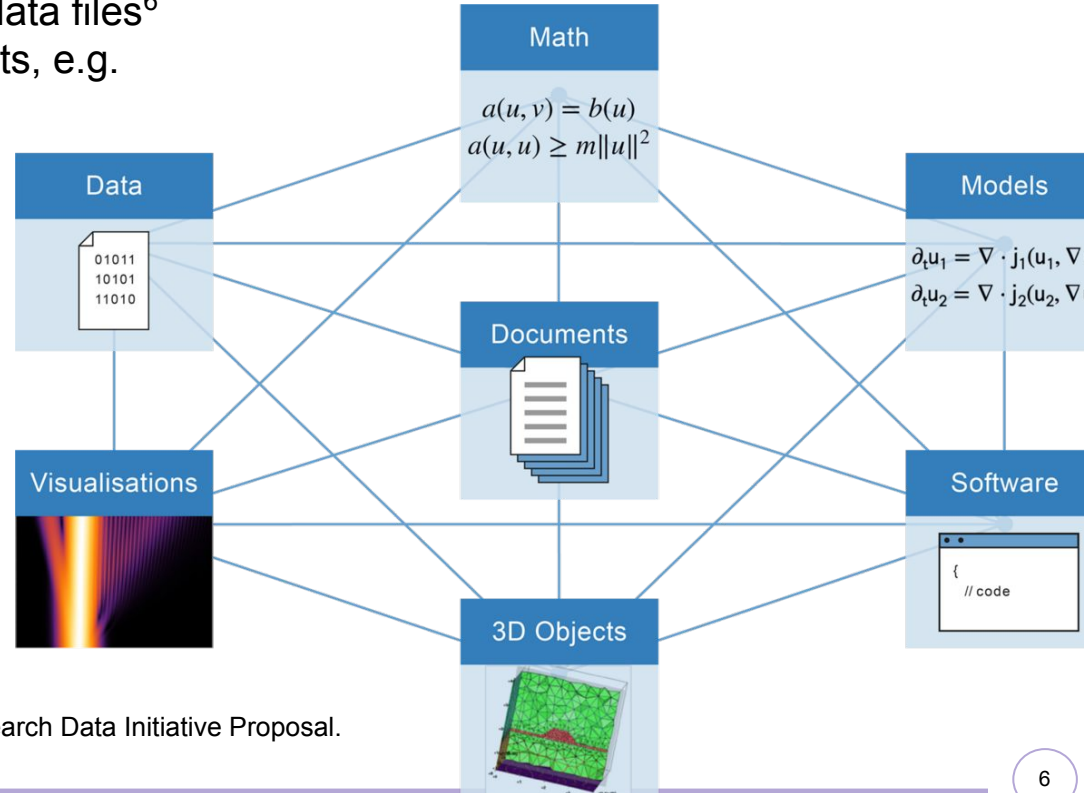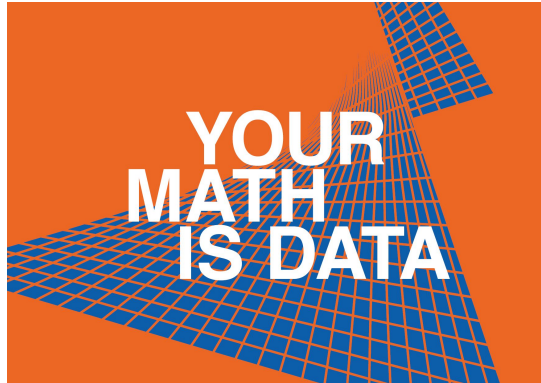
# Dark Data - Example 2: Mathematics

- Another example from the field of mathematics[5]

"[..] a classification of all conditional independence structures on up to four discrete random variables, originally published in a series of papers (Matus and Studeny, 1995; Matus, 1995, 1999). [..] Simcek (2006) digitized this result and left the field after his PhD in 2007. His **research data was deleted in 2021 from his former institute's website** [..]. It was encoded in a packed binary **format which is hard to read, search, and reuse**. Some files supporting the correctness of the classification for binary distributions use an **unspecified, compiler-specific binary serialization format** for floating-point data. The programs used for the creation and inspection of the database were written in a **dialect of the Pascal** programming language which **has not been maintained since 2006**. The **sparse documentation is in Czech**."

[5] Boege, Tobias, et al. Data Management Planning in the German Mathematical Community. arXiv preprint arXiv:2211.12071 (2022).
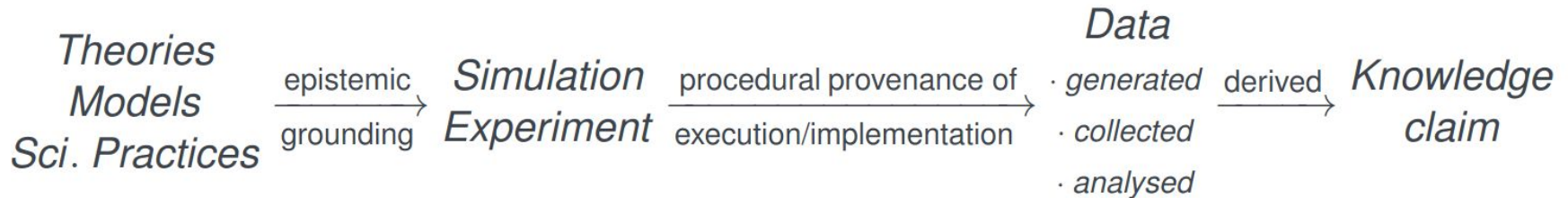
# Your Math is Data!

- Research Data goes beyond pure data files[6]
- It's a spiderweb of information assets, e.g.
  - Models
  - Algorithms
  - …
- MaRDI@NFDI is working on this





[6] The MaRDI consortium. (2022). MaRDI: Mathematical Research Data Initiative Proposal. https://doi.org/10.5281/zenodo.6552436l

# Dark Data vs. FAIR Data

- FAIR has reached a quasi-standard paradigm in research data management
- Dark data is not FAIR -> Diminish dark data to achieve FAIRness
  - Organizational measures
    - Data Management Plans
    - Scientific Data Officers / Data curators
    - Incentives (extrinsic / intrinsic)
  - Technical measures
    - Semantic Technology
    - Data infrastructures
- Usually we limit ourselves to record pure data provenance (origin and genesis) when talking FAIR



*Theories Models Sci. Practices* → epistemic grounding → *Simulation Experiment* → procedural provenance of execution/implementation → *Data · generated · collected · analysed* → derived → *Knowledge claim*

# Reproducibility and FAIR Data

- FAIR does not make statements about reproducibility, just about reusability

Example: A study that analyzed 108 publications for reproducibility:

Data accessibility | Data availability | Reproducibility

Some availability: 41

Some online: 39

Fully Reprod.: 4

Partly Reprod.: 1

108 Publications

Not Reprod.: 34  } This is data FAIR, however not reproducible

Not specified where: 41 → Only in article: 1

Theoretical: 26

Author request: 1

Riedel, Christian, et al. "Including data management in research culture increases the reproducibility of scientific results." *INFORMATIK 2022* (2022).

# …is FAIR really the End of the Line?

- FAIR formulates a minimum standard in data documentation
- Its ultimate goal is to optimize data for reusability with regards to data formats, licenses, retrievability, and provenance
- However it does not make statements about
  - To what extent the data is reproducible
  - What's the scientific ground for the data?
  - Why is it valid
  - Which claims have been formulated
  - Responsibility
- There is much more than just FAIR data
  - RIOT[7]: Reproducible, interpretable, open, transparent
  - CARE[8]: Collective benefit, authority to control, responsibility, ethics
  - XAIR[9]: Explainable AI ready

[8] E. Ganley et al., BMC Res. Notes 15: 51, doi:10.1186/s13104-022-05932-5, 2022
[9] S. Russo Carroll et al., Sci. Data 8: 108, doi:10.1038/s41597-021-00892-0, 2021
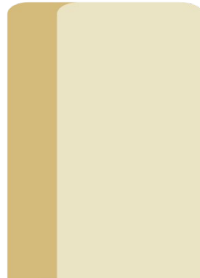[10] Horsch, M., et al. "Epistemic metadata for computational engineering information systems." *Manuscript, to appear in Proc. FOIS 2023 (2023).*
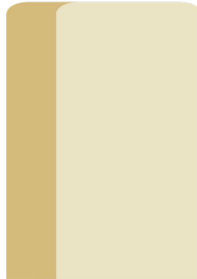
# RIOT

**R** Reproducible
Get the same answer asked of the same or different dataset

**I** Interpretable
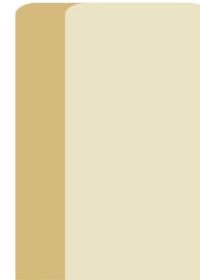Be clear, concise, accessible, and unambiguous

**O** Open
Open, inclusive, diverse, collective effort

**T** Transparent
Whenever possible, make public every part of research

# CARE

# Epistemic Metadata and XAIR

- Goal: Make the knowledge reusable and not only the data
- This can be accomplished by epistemic metadata
  - documenting the knowledge status of data
  - Research data must be stored and exchanged jointly with this metadata
- This makes the data XAIR (**Explainable** <u>AI</u> <u>r</u>eady)
- All XAIR data is FAIR
- Not all FAIR data is XAIR
- Metadata standardization is utterly important



Personification of Episteme in Celsus Library in Ephesus, Source: Wikipedia

# Toward data documentation standardization meeting regulatory and users' requirements

Martin T Horsch

Dept of Data Science

Norwegian University of Life Sciences

# Responsible data documentation

**European AI Act proposal:** "To address the **opacity** that may make certain AI systems incomprehensible to or too complex for natural persons, a certain degree of transparency should be required for high-risk AI systems. [...] High-risk AI systems should therefore be accompanied by **relevant documentation**"

- "High-risk" includes energy safety, water, *etc.*, and all that affects fundamental rights.
- This is not in force yet, negotiations are taking place at least until end of the year.

**Epistemic opacity** (Humphreys, 2011): A cognitive "process is **epistemically opaque** relative to [...] agent $X$ at time $t$ [... if …] X does not know at t all of the epistemically relevant elements"

Tendency: Making data trustworthy through explanations becomes a **legal** requirement.

# Responsible data documentation

**European AI Act proposal:** "To address the **opacity** that may make certain AI systems incomprehensible to or too complex for natural persons, a certain degree of transparency should be required for high-risk AI systems. [...] High-risk AI systems should therefore be accompanied by **relevant documentation**"

- "High-risk" includes energy safety, water, *etc.*, and all that affects fundamental rights.
- This is not in force yet, negotiations are taking place at least until end of the year.

**Epistemic opacity** (Humphreys, 2011): A cognitive "process is **epistemically opaque** relative to [...] agent *X* at time *t* [... if …] X does not know at t all of the epistemically relevant elements"

- The "epistemically relevant elements" from Humphreys are the same as the "relevant documentation" from the AI Act. We call them the **epistemic metadata**.

Tendency: Making data trustworthy through explanations becomes a **legal** requirement.

This means that **explainable-AI-ready** (XAIR) data cannot rely on *informal* metadata standardization. *Formal standardization* going through the official agencies becomes necessary.

# From informal to formal standards

**European AI Act proposal:** "To address the **opacity** that may make certain AI systems incomprehensible to or too complex for natural persons, a certain degree of transparency should be required for high-risk AI systems. [...] High-risk AI systems should therefore be accompanied by **relevant documentation**"

- "High-risk" includes energy safety, water, *etc.*, and all that affects fundamental rights.
- This is not in force yet, negotiations are taking place at least until end of the year.

**Epistemic opacity** (Humphreys, 2011): A cognitive "process is **epistemically opaque** relative to [...] agent $X$ at time $t$ [... if …] X does not know at t all of the epistemically relevant elements"

- The "epistemically relevant elements" from Humphreys are the same as the "relevant documentation" from the AI Act. We call them the **epistemic metadata**.

Beginning with the EC's Battery Regulation, **digital product passports** will be mandatory first for batteries, later textiles, electronics, and successively more and more products.

- Characterizing the **knowledge status** becomes a priority.

# Molecular modelling case study

**Epistemic metadata** and their **documentation** were explored in molecular thermodynamics:

**First stage report (10 cases)**, doi:10.5281/zenodo.7516532, **2023**.

Discussion of five papers each from two research groups (Berlin, London) without involving the papers' authors. Obtained a tentative **taxonomy for epistemic metadata** and explored the patterns of epistemic grounding.

**Second stage report (12 claims)**, doi:10.5281/zenodo.7608074, **2023**.

Discussion of two claims each from six papers, involving the papers' authors, some of whom became co-authors of the present work. **Ontology of epistemic metadata**, except for epistemic grounding, implemented in PIMS-II.

Good data documentation standards give researchers the freedom to say what they want to say. Ontologies should **provide a language, not micromanage** researchers' self-expression.

# Epistemic metadata and reproducibility claims

# Epistemic metadata

Metadata are "descriptive data about an object" (ISO 11179).

Epistemic metadata are metadata that support **characterizing the knowledge status** of data.

Epistemic metadata:

a) "what **knowledge claim** $\varphi$ has been formulated?,"

b) "where do the data and the claim come from?" (**provenance**),

c) "what **validity claim** was made about $\varphi$?,"

d) "why should we accept any of this?" (**grounding**).

These concepts are implemented in the PIMS-II ontology.

# Subject matter of research data

We understand subject matter of a knowledge claim and/or the associated research data as given by the research question that is being answered, or by the "equivalence relation over logical space" with respect to that question.[1]

Assertion: "*A is the factually correct answer to question Q.*"

Subject matter of the assertion: Q.

Equivalence relation: Two states of affairs are equivalent if the answer to Q is the same for both. *Two knowledge bases are equivalent if they return equivalent tables for the respective SPARQL query.*

[1]S. Yablo, *Aboutness*, Princeton Univ. Press (ISBN 978-0-691-14495-5), **2014**.

# Subject matter of research data

We understand subject matter of a knowledge claim and/or the associated research data as given by the research question that is being answered, or by the "equivalence relation over logical space" with respect to that question.[1]

With respect to the research question[2]

$$q_1 = \text{"What is the } \mathbf{D} \text{ matrix of liquid } M \text{ as a function of } \mathbf{x}, p, \text{ and } T?\text{,"}$$

two states of affairs are equivalent if their $\mathbf{D}(\mathbf{x}, p, T)$ dependencies are the same.

[1]S. Yablo, Aboutness, Princeton Univ. Press (ISBN 978-0-691-14495-5), **2014**.

[2]G. Guevara Carrión *et al.*, *J. Phys. Chem. B* **124**(22): 4527–4535, doi:10.1021/acs.jpcb.0c01625, **2020**.

# Reproducibility claims

Common formulation and schema for reproducibility claims (RCs):

«Whenever research process $\kappa$" is carried out, it must lead to the outcome $\varphi$".»

1. Researcher $a$ did $\kappa$ and found $\varphi$.

2. Researcher $b$ did $\gamma$, somehow similar to $\kappa$, and found something that is inconsistent with $\varphi$.

3. Now we think that $\varphi$ has not been reproduced successfully, maybe it is "falsified." But why?

# Reproducibility claims

Common formulation and schema for reproducibility claims (RCs):

«Whenever research process $\kappa$" is carried out, it must lead to the outcome $\varphi$".»

1. Researcher $a$ did $\kappa$ and found $\varphi$.

   Here, $a$ also made the **positive reproducibility claim** $\psi = \Box(\varphi" \mid \kappa")$.

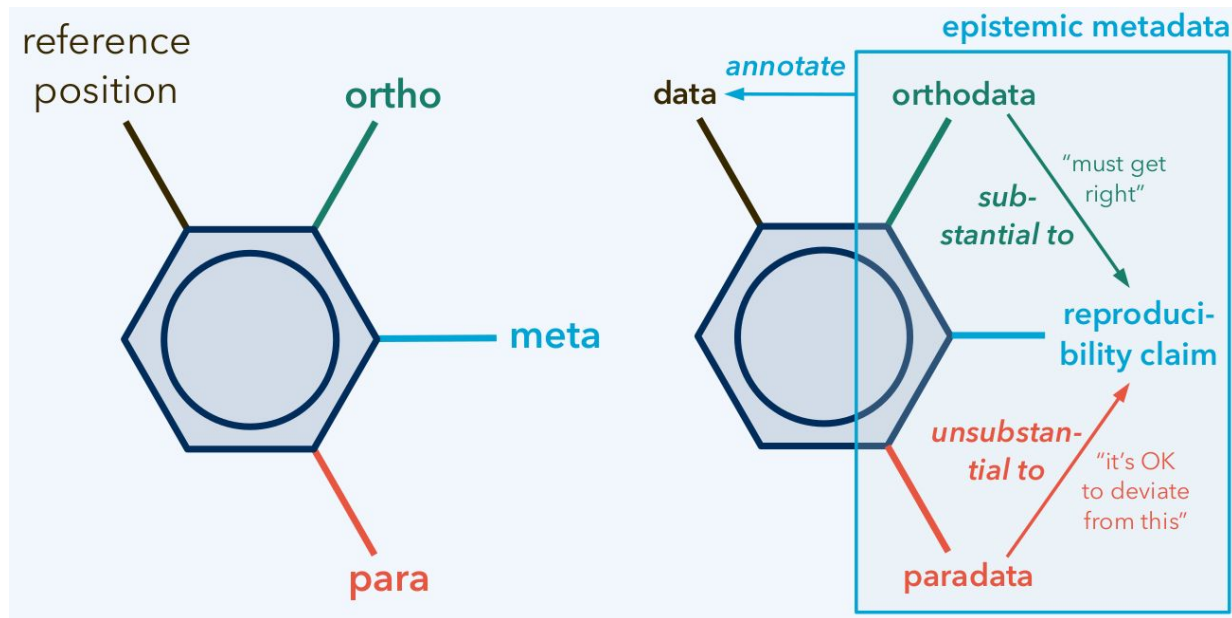   - Read $\Box(\varphi" \mid \kappa")$ as «given that $\kappa$" is done, necessarily a result consistent with $\varphi$" is obtained.»

2. Researcher $b$ did $\gamma$, consistent with $\kappa$", and found something that is inconsistent with $\varphi$".

   Here, $b$ made the **negative reproducibility claim** $\Diamond(\neg\varphi" \mid \kappa") \equiv \neg\Box(\varphi" \mid \kappa") \equiv \neg\psi$.

   - Read $\Diamond(\neg\varphi" \mid \kappa")$ as «if $\kappa$" is done, it can happen that a result consistent with $\neg\varphi$" is obtained.»

3. What is relevant there is the contradiction between $\psi$ and $\neg\psi$.
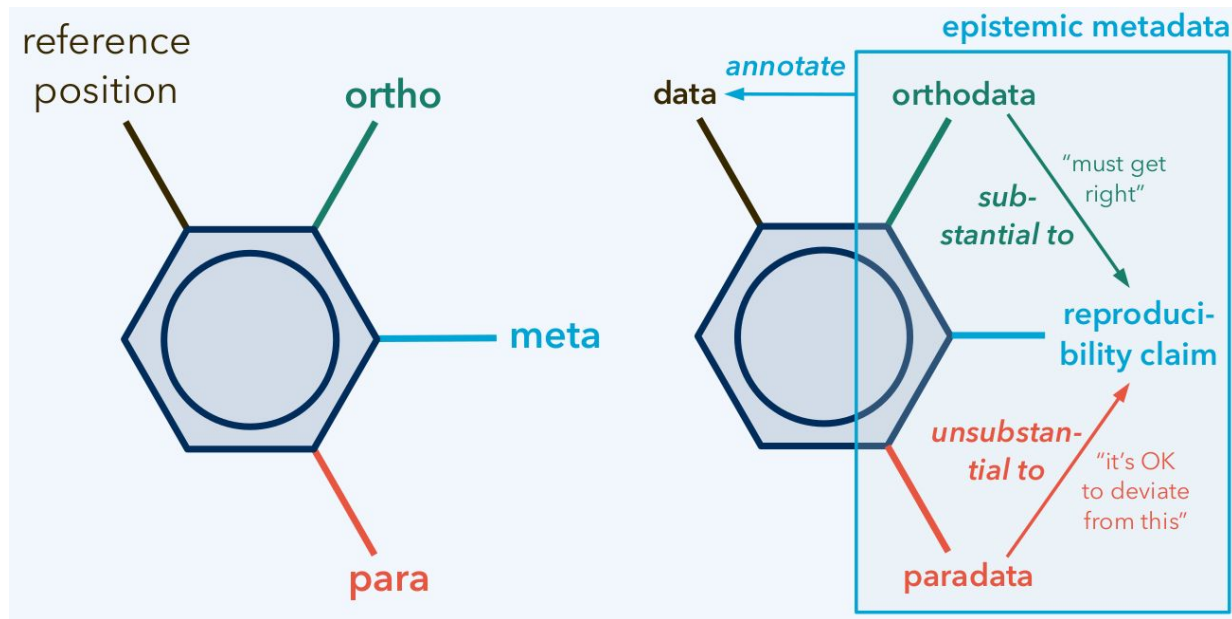
# Reproducibility claims



provenance metadata $\kappa$
provenance paradata $\kappa'$

output metadata $\varphi$
output paradata $\varphi'$

«repeat $\kappa$, but no need to retain $\kappa'$»

«obtain $\varphi$ again, except for $\varphi'$ maybe»

# Paradata and logical subtraction



«repeat $\kappa$, but no need to retain $\kappa$'»

«obtain $\varphi$ again, except for $\varphi$' maybe»

provenance metadata $\kappa$

provenance paradata $\kappa$'

$\kappa$'' = $\kappa$ − $\kappa$'

(provenance orthodata )

output metadata $\varphi$

output paradata $\varphi$'

$\varphi$'' = $\varphi$ − $\varphi$'

(output orthodata )

# Paradata and logical subtraction

**Logical subtraction** is a concept from analytic philosophy.

Its formalization is closely connected to the theory of **subject matter**.

Example from Yablo (*Aboutness*, **2014**): Someone who rejects ontological commitment to the existence of numbers is asked how many prime numbers there are greater than ten. "Infinitely many, of course, except that numbers don't exist."

Reproducibility in computational engineering:

> Could you try to replicate my old simulation result? Just do the same as I did.
>
> Except that you of course log in with your user account, not mine.
>
> Your result was off by 0,5%? Don't worry, that is totally normal.

# Connection to NFDI4Ing work



**DORIS**

- … will develop **standards for reproducibility**.
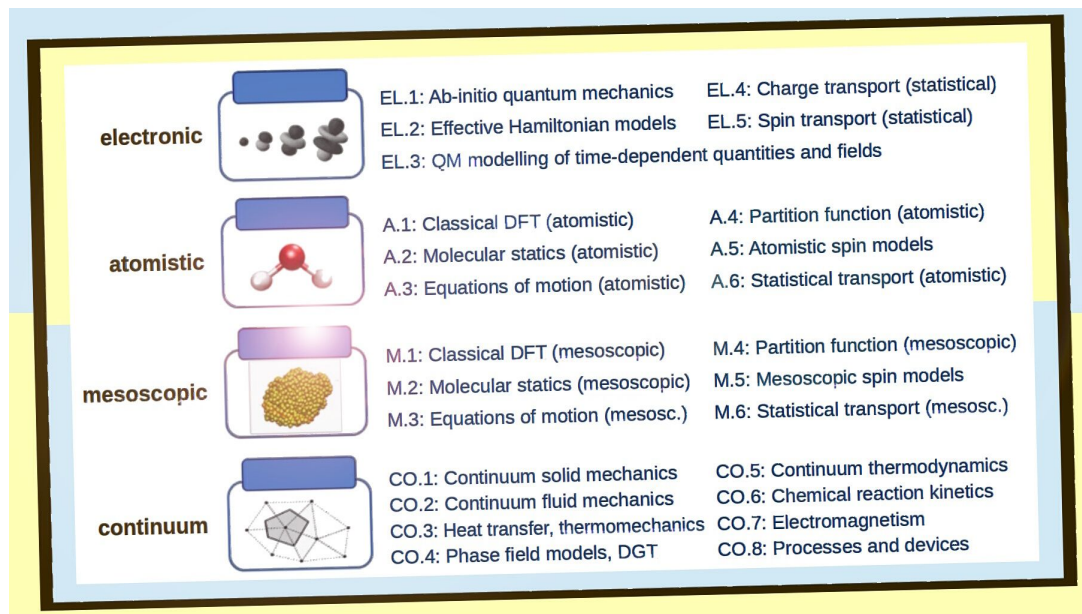- … and **best-practice guidelines** for reproducibility.

# Previous work at the European level

# CEN Workshop Agreements (CWAs)

As an attempt at metadata standardization, **MODA** resulted in a closed epistemic space with a rigid categorization of modelling and simulation methodologies.
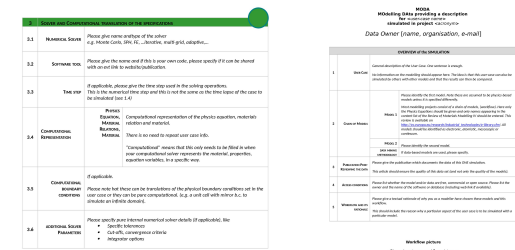


**MODA:** "Model Data"

**CWA 17284**:2018 E

# Criticism of the old CWAs

As an attempt at metadata standardization, MODA resulted in a closed epistemic space with a rigid categorization of modelling and simulation methodologies. Both **MODA** and **CHADA** documentations are **hard to create** and **hard to use** by humans, but **not machine-actionable**.



**MODA:** "Model Data"
**CWA 17284**:2018 E

**CHADA:** "Characterization Data"
**CWA 17815**:2021 E

# Criticism of the old CWAs

Priorities (**DORIC principles**) following doi:10.5281/zenodo.4571052

**D** diversify technology
**O** observe practices
**R** have **realistic** objectives
**I** incentivize open data
**C** co-design data and workflows

- **MODA** was a **closed epistemic space**, modelling methods had to be chosen from a small list.

- **MODA** imposed a **given level of detail** in data documentation; namely, **unrealistically** detailed.

- **MODA** documentations are **complicated**, and of **limited use** to all, including to humans.

# Toward meeting user requirements



Priorities (**DORIC principles**) following doi:10.5281/zenodo.4571052

**D** diversify technology
**O** observe practices
**R** have **realistic** objectives
**I** incentivize open data
**C** co-design data and workflows

- MODA was a closed epistemic space, modelling methods had to be chosen from a small list. **ModGra** gives the user a **highly expressive** graph **language** to describe their method.
- MODA imposed a given level of detail in data documentation; namely, unrealistically detailed. **ModGra** gives the user the choice to document the model at a **flexible level of detail**.
- MODA documentations are complicated, and of limited use to all, including to humans. **ModGra** specifies semantics at the level of physics and is **actionable** through ProMo tools.

# CWA 17960:2022
# ModGra: A graphical representation of physical process models

Heinz A Preisig

Dept of Chemical Engineering

Norwegian University of Science and Technology

# Contents

1. How to go about generating a CWA
2. ModGra
   a. motivation
   b. approach
   c. basic components
   d. examples

# CWA

# CEN workshop agreement

step 1: decide on what shall be standardise -- aim at something minimal

step 2: contact people who may be interested

step 2: draft a working plan

step 3: get in contact with a standardisation organisation in a European country

step 4: provide workshop proposal to standardisation organisation

( form https://boss.cen.eu/media/BOSS%20CEN/formtemp/ws_proposal.docx )

step 5: form a committee -- aim at a wide spread in terms of expertise

step 6: standardisation organisation submits proposal to CEN and announce 30 days in advance a kick-off meeting

step 7: keep on meeting until standard is established -- standardisation organisation provides the secretary

step 8: the CWA is submitted to CEN

Instructions can be found on : https://boss.cen.eu/developingdeliverables/cwa/pages/ or google for *how to generate CEN workshop agreement*

# Practicalities

- Aim at a small document
- Simple over complex
- A main body of a CWA document is structured like a contract
  - terms are defined
  - terms can only be used once defined or defined as commonly known
- Adding examples helps
- Simple language helps

Experience: terminology requires a lot of effort

# ModGra

# Approach: Reductionism

Break a process recursively down into smaller and smaller entities

How long?

- until a level of granulation is achieved that captures the essentials of the process
- the granules can be viewed as **simple systems, control volumes** characterised by
  - a time scale
  - a distribution property
    - intensive properties are a function of the location → **distributed** systems
    - intensive properties are NOT a function of the location → **lumped** systems

# Approach

View a process model as a directed graph with:

- nodes being capacities, control volumes containing conserved extensive quantities
- arcs transferring conserved extensive quantities

Add information processing as directed graph with main application of control:

- nodes being input/output functions
- arcs transferring information -- signals

Abstraction: directed graph with tokens living in there, giving the graph context
similar to Petri's thinking

# Foundation

A physical process contains

- mass,
- energy,
- momentum,
- charge

$\rightarrow$ conserved quantities.

There are

- "holders"/"accumulators" and
- "transfers"
  of conserved quantities

Abstract to

**tokens**

being stored/accumulated and move

- **capacities** for tokens
- **transfer** of tokens

# Foundation

Abstract to

**tokens**

being stored/accumulated and move

- **capacities** for tokens
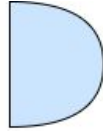- **transfer** of tokens

Abstract to

**directed graph**

with

- nodes being the capacities
- arcs the transport of tokens between capacities

# Capacities

| | | | |
|---|---|---|---|
| constant | | physical : | temporal constant, spatial: uniform & infinitely large |
| | | information: | temporal constant |
| dynamic lumped | | physical : | temporal dynamic, spatial: uniform (0D) & finite size |
| | | information: | temporal dynamic |
| dynamic distributed | 1-3D | physical : temporal dynamic, spatial: not uniform (nD, n := 1,2,3) & finite size | |
| dynamic point | ● | physical : | temporal event-dynamic, spatial: infinite small |
| | | information: | temporal event-dynamic |
| event-dynamic distributed | | physical : | temporal event-dynamic, spatial: not uniform & finite size |
| | | information: | temporal event-dynamic |

# Extensions

| | | |
|---|---|---|
| composite | | physical :    composite entity -- a subgraph<br>information:    ditto |
| surrogate | | physical :    surrogate replacing typically a composite entity -- a subgraph<br>information:    ditto |

Enables

- the construction of large systems
- empirical models

# Transfers

| | | |
|---|---|---|
| ——— label ———→ | physical : | mass transfer |
| - - - · label - - - → | physical : | diffusional mass transfer |
| - - - · label - - - → (red) | physical : | conductive heat transfer |
| - - - · label - - - → (blue) | physical : | work flow |
| ········· label ········→ (olive) | information : | signal |

Directionality defines a reference coordinate system for the respective flow

# A simple example

A
B

tank reactor plant with reaction
A + B -> C

A,B,C

A
B

A+B -> C

A,B,C

Simplest model
- reservoirs for supply and product
- lumped liquid volume in the tank
- no control

# More complex mixing model



A+B -> C

Only hydraulic shown
Would need definition
- where reaction takes place
- where A and B flows in
- where product is drawn

top
loop

bottom
loop

internal
loop

# Adding control

# Multi-scale model -> workflow

**Melting process**

three staged melting process

**Molecular modelling**

n : species mass
V : volume
p : pressure
T : temperature
E : energy
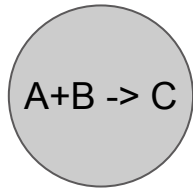
d : distribution
s : control signal

thermostat
adjust temperature
$T_m - T < \varepsilon$

barostat
$p_m - p < \varepsilon$

$T_m - T < \varepsilon$
$p_m - p < \varepsilon$

equipartition --> T

Boltzmann dist

adjust volume

velocity distribution

switch @ start

equilibrium condition
V & E converged

macroscopic level

switch @ start

unit cell

time-windowed average volume

viral theorem averaging --> p

time-windowed average energy

molecular level

macroscopic level

# Model simplification

# Stirred tank



schema of a standard stirred tank



a key assumption : contents of the tank, the liquid, is ideally stirred --> modelled as a lumped entity

**Left diagram labels:**

room
gas phase / air
breating pipe
lid
reactant source
reactant
feed pipe
flansh top
condensed liquid
h/c hot
jacket return
flansh bottom
condensate
h/c cold
jacket inlet
jacket contents
gas contents
wall
shell
wall
vessel liquid contents
contents
jacket
reactor
outflow
product sink

@ assumption: shell ideally insulated --> no heat loss to the room from the shell
@ assumption: flansh on the top is insulated from the flansh on the lower piece

**Right diagram labels:**

room
reactant source
reactant
feed pipe
flansh bottom
h/c hot
jacket return
wall
h/c cold
jacket inlet
jacket contents
wall
vessel liquid contents
contents
jacket
reactor
outflow
product sink

@ assumption: temperature in contents approximately room temperature

- insignificant evapouration from the fluid and negligible heat losses to the room
- outer construction is not actively storing energy
- gas capacity is not active

54

**Left diagram:**

room

reactant source

reactant → feed pipe

h/c hot ← jacket return

h/c cold → jacket inlet

jacket contents

vessel liquid contents

contents

jacket

outflow

reactor

product sink

@ assumption: fast heat transfer through the wall
@ assumption: wall capacity negligible
@ assumption: jacket uniform

**Right diagram:**

room

reactant source

reactant

h/c hot ← jacket return

h/c cold → jacket inlet

jacket contents

vessel liquid contents

contents

jacket

reactor

product sink

@ assumption: inflow and outflow without deadtime

# Use in model-design software

**ProMo_Sandbox8**

**extracting_reactor**

## main

### Navigation

- ○ D - delete
- ○ E - explore
- ◉ Esc - explore
- ○ G - explode
- ○ I - insert

**connector**

**name node** -

**controls**

**topology**

### nodes
- ◉ 0 :: constant|infinity
- ○ 1 :: dynamic|lumped
- ○ 2 :: event|lumped

### networks
- ○ 0 :: control
- ○ 1 :: gas
- ◉ 2 :: liquid
- ○ 3 :: material
- ○ 4 :: reactions
- ○ 5 :: solid

### named_networks
- ○ 0 :: A-liquid
- ◉ 1 :: B-liquid

### edit
- add
- edit
- delete
- colour

### token
- ◉ 0 :: charge
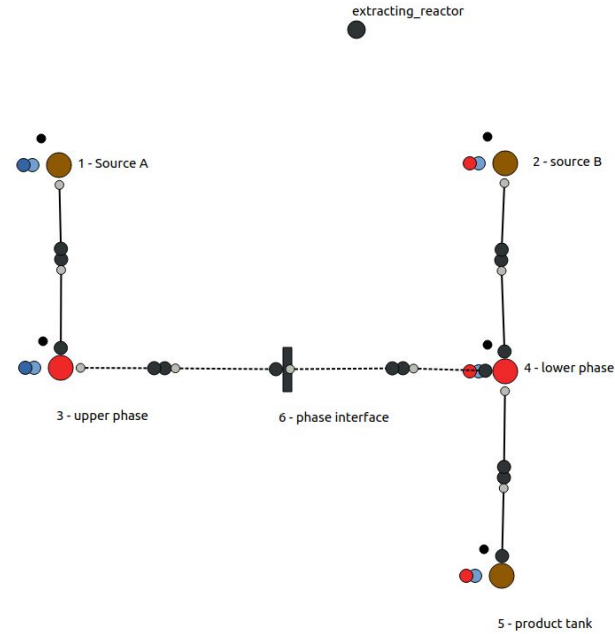- ○ 1 :: energy
- ○ 2 :: mass

### mechanism
- ◉ 0 :: conduction

### nature
- ◉ 0 :: lumped

### variants

**node variant**        **arc variant**

extracting_reactor

1 - Source A

2 - source B

3 - upper phase

6 - phase interface

4 - lower phase

5 - product tank

57

# ProMo -- Process Modeller

- expert section defines primitive blocks
- translator builds models using the primitive blocks
  and generate higher-level models again as a building block
- all generated building blocks can be reused
- automatic code generation

Software suite ProMo written in Python using pyqt and deployed with ABCdesktop as browser application on https://promo-abcloudtop.io/ as beta release. First release soon. Will be announced on my webpage https://folk.ntnu.no/preisig/

# Conclusions

- ModGra is a powerful model design tool
    - as a discussion tool
    - model reduction
    - multi-scale workflow generation tool
    - generation of a model library -- model dissemination