



Norges miljø- og  
biovitenskapelige  
universitet



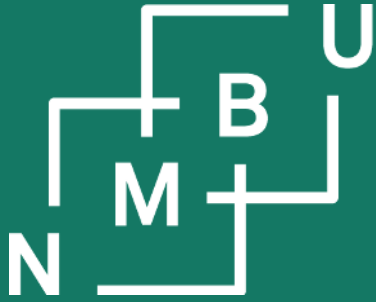
# Case study on epistemic metadata in molecular modelling

Martin Horsch,<sup>1,2</sup> Silvia Chiacchiera,<sup>2</sup> and Björn Schembera<sup>3</sup>

<sup>1</sup>Norwegian Univ. Life Sciences, Faculty of Science and Technology, Ås, Norway

<sup>2</sup>UKRI STFC Daresbury Laboratory, Scientific Computing Department, Daresbury, UK

<sup>3</sup>Univ. Stuttgart, Institute of Applied Analysis and Numerical Simulation, Stuttgart, Germany



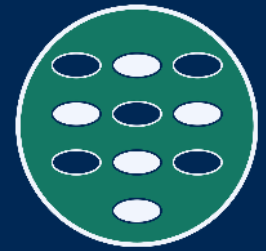
Noregs miljø- og  
biovitenskapelige  
universitet

# Domain and background

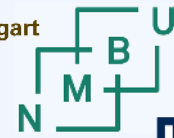
First stage and an ontology

Second stage and grounding

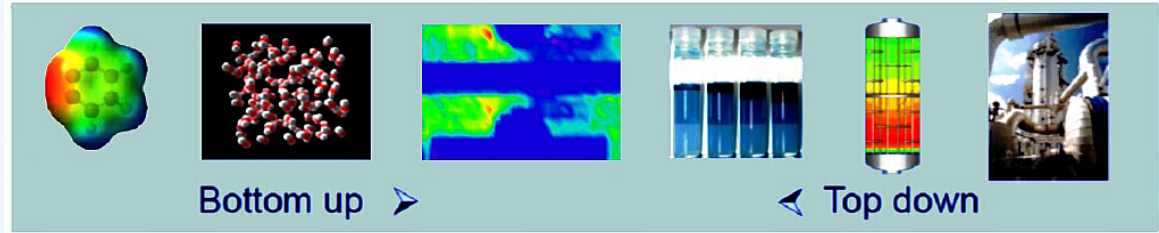
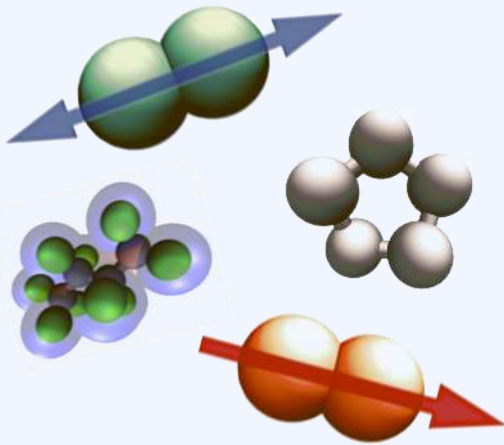
Materialteori og -informatikk



Digitalisering på Ås



# Molecular simulation in engineering



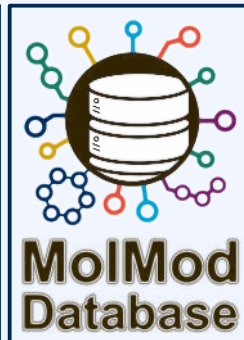
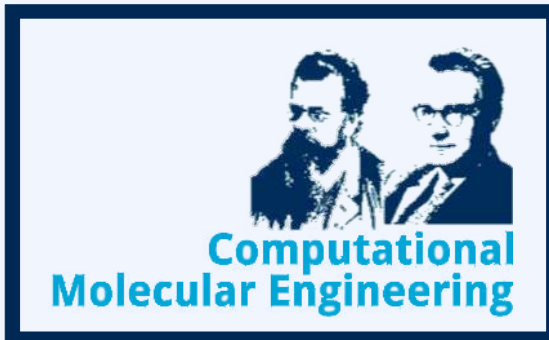
physics-driven aspects  
(qualitative validity)

data-driven aspects  
(quantitative reliability)

- Realistic representation of underlying physical features

- Models with parameters that can be adjusted to data

- Reliable interpolation, extrapolation, prediction



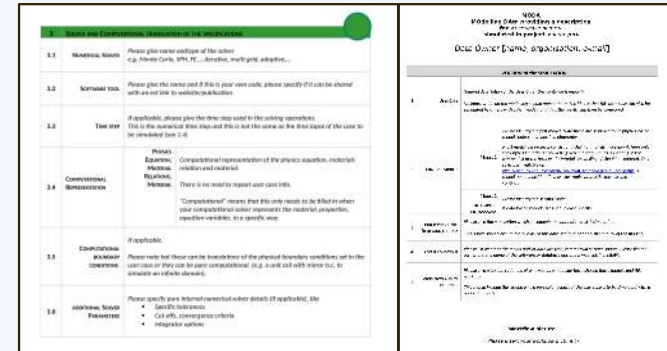
- Process Engineering
- Thermodynamics
- Scientific Computing



# The EMMC's first attempts

## EMMC: European Materials Modelling Council

The EMMC's documentation standard MODA ("model data") **failed to meet researchers' needs** by making too much annotation mandatory, much more than needed in practice.



- MODA was a closed semantic and epistemic space: Modelling methods had to be chosen from a small list.<sup>1, 2</sup>
- MODA imposed a **given level of detail** in workflow documentation.<sup>1</sup>
- MODA documentations were **complicated**.<sup>3</sup>

**CWA 17284:**  
2018 E  
«MODA»



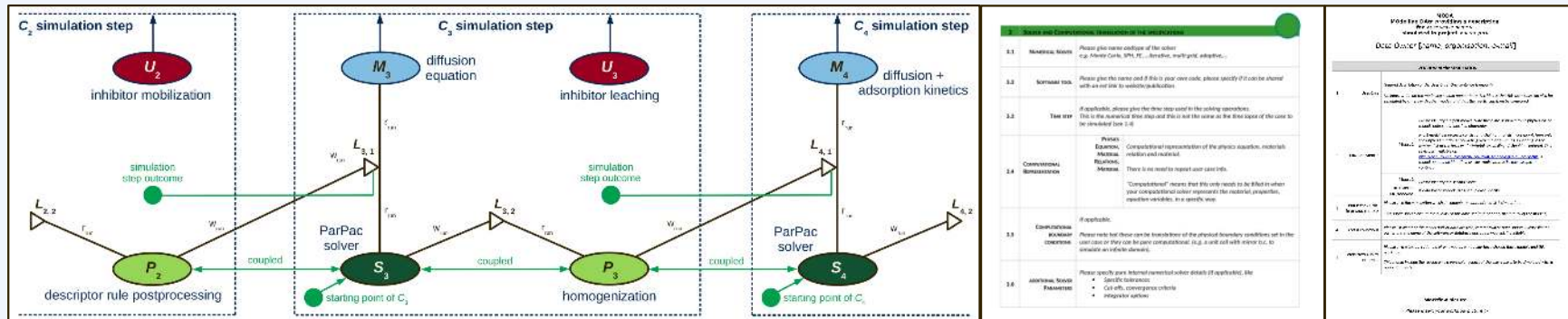
<sup>1</sup>CEN Workshop Agreement 17284:2018 E, «Materials modelling: Terminology, classification and metadata», 2018.

<sup>2</sup>A. F. de Baas, *What Makes a Material Function? Let Me Compute the Ways*, EU Publications, doi:10.2777/417118, 2017.

<sup>3</sup>ReaxPro project deliverable D2.1, «ReaxPro MODA diagrams», 2020.

# The EMMC's first attempts

## Focus on provenance documentation



- MODA was a **closed semantic and epistemic space**: Modelling methods had to be chosen from a small list.
- MODA imposed a **given level of detail** in workflow documentation; namely, **unrealistically detailed**.
- MODA documentations were **complicated** and **of limited use** to all, including to humans.<sup>1, 2</sup>

**CWA 17284:**  
2018 E  
«MODA»



<sup>1</sup>ReaxPro project deliverable D2.1, «ReaxPro MODA diagrams», 2020.

<sup>2</sup>«European standardization efforts from FAIR toward explainable-AI-ready data documentation in materials modelling», in *Proc. ICAPAI 2023*, doi:10.1109/icapai58366.2023.10193944, IEEE, 2023.

# Case study in molecular modelling

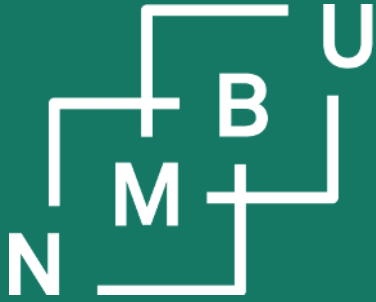
**Epistemic metadata** and their documentation were explored for the domain of molecular modelling and simulation within engineering thermodynamics:

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

Discussion of *five papers each* from *two research groups* (London, Berlin) without involving the papers' authors. Obtained a tentative **taxonomy for epistemic metadata**, later implemented into the PIMS-II ontology.

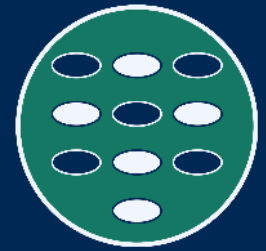
**Second stage report (12 claims)**, doi:10.5281/zenodo.7608074, **2023**.

Discussion of *two claims each* from *six papers*, with two papers each from three research groups (London, Berlin, Kaiserslautern), involving the papers' authors. Discussed aspects such as the **grounding of knowledge claims** with authors.



Noregs miljø- og  
biovitenskaplege  
universitet

Materialteori og -informatikk



Digitalisering på Ås

# Domain and background

## First stage and an ontology

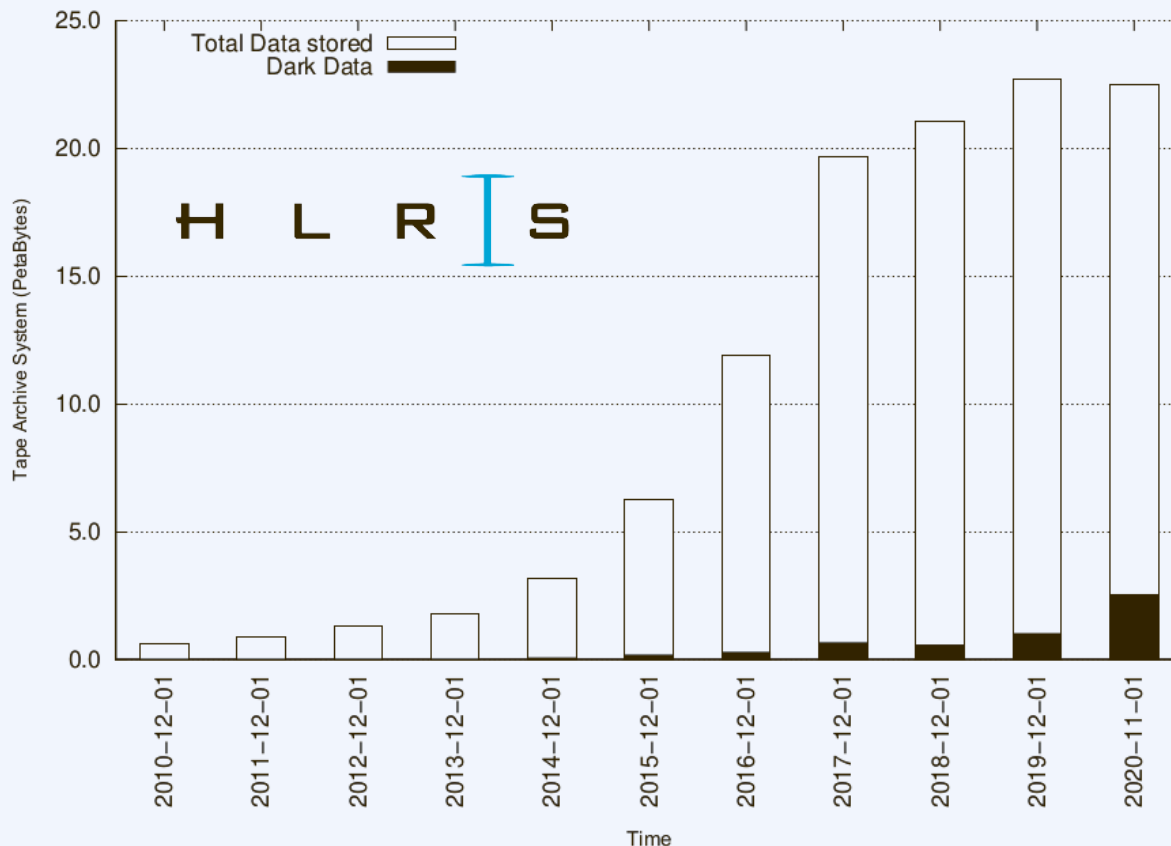
### Second stage and grounding

Key epistemic metadata items are the **knowledge claims** made based on data, their **provenance**, **validation and reproducibility**, and **epistemic grounding**.

# Problem and idea behind epistemic metadata

**Dark data** are data with an uncharacterized knowledge status.

In other words: *We don't know what we know from and about the data.*



**Flood of dark data:**  
*More and more data are accumulated, but are dark - and useless.*

Discussed in work by Björn Schembera and Juan Durán.<sup>1,2</sup>

<sup>1</sup>Figure from Björn Schembera's doctoral thesis, doi:10.18419/opus-11028, **2019**.

<sup>2</sup>B. Schembera, J. Durán, *Philos. Technol.* **33**: 93-115, doi:10.1007/s13347-019-00346-x, **2019**.



# Problem and idea behind epistemic metadata

**Epistemic metadata** are the information that **establishes the knowledge status** of data or digital objects.<sup>1</sup>

**Questions we must answer to establish the knowledge status:**

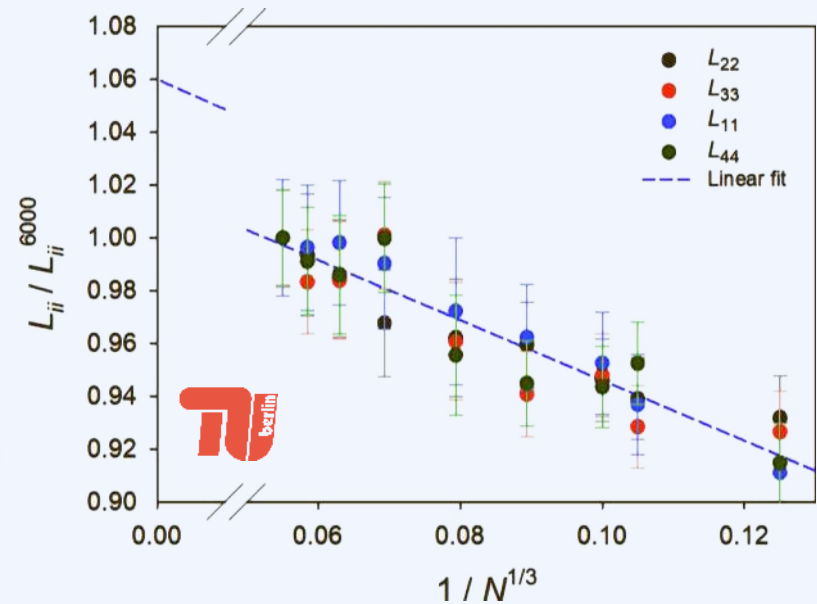
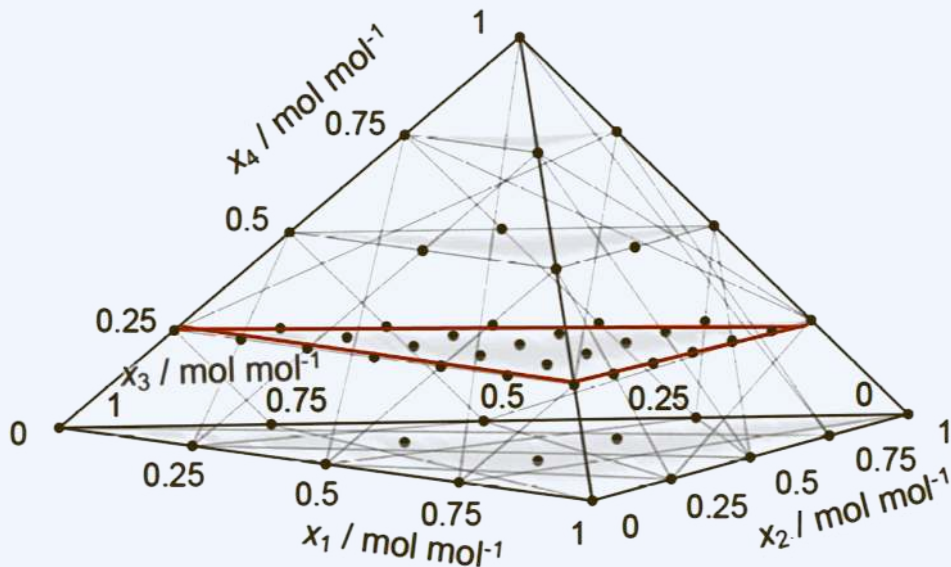
- a) "what **knowledge claim**  $\varphi$  has been formulated?,"
- b) "where do the data and the claim come from?" (**provenance**),
- c) "what **validity claim** was made about  $\varphi$ ?,"
- d) "why should we accept any of this?" (**grounding**).

Key epistemic metadata items are the **knowledge claims** made based on data, their **provenance**, **validation and reproducibility**, and **epistemic grounding**.

<sup>1</sup>«Documentation of epistemic metadata by a mid-level ontology of cognitive processes», in *Proc. JOWO 2022*, CEUR vol. **3249**: p. 2 (CAOS), CEUR-WS, **2022**.

# The first stage of the case study

**Example:** The work by Guevara *et al.*<sup>1</sup> (2020) was considered at both stages.<sup>2,3</sup>



<sup>1</sup>G. Guevara Carrión, R. Fingerhut, J. Vrabec, «Fick diffusion coefficient matrix of a quaternary liquid mixture by molecular dynamics», *J. Phys. Chem. B* **124**(22): 4527–4535, doi:10.1021/acs.jpccb.0c01625, **2020**.

<sup>2</sup>M. T. Horsch, B. Schembera, «Epistemic metadata in molecular modelling: First-stage case-study report (10 cases)», Inprodat technical report 2023-A, doi:10.5281/zenodo.7516532, **2023**.

<sup>3</sup>M. Horsch, S. Chiacchiera, G. Guevara, M. Kohns, E. Müller, D. Šarić, S. Simon, I. Todorov, J. Vrabec, B. Schembera, «Epistemic metadata in molecular modelling: Second-stage case-study report (12 claims)», Inprodat technical report 2023-B, doi:10.5281/zenodo.7610237, **2023**.

# Guevara et al. (2020) paper:<sup>1</sup> First-stage analysis<sup>2</sup>

**Question:** What is a good methodology for obtaining Fick diffusion coefficients in multicomponent mixtures by [equilibrium molecular dynamics] simulation?

**Object of research:** The object of research is the Fick diffusion coefficient matrix as such.

**Knowledge claim:** [...] methodology [...] first, the explicit inclusion of a finite-size correction, where it is specifically novel that this correction is applied to the Onsager coefficients, and second, obtaining the Darken correction from [Kirkwood-Buff] integrals.

**Grounding:** KB part [...] validated against “the Wilson excess Gibbs energy model [...]” [...] not clear what should make us accept the finite-size methodology [...]. It yields a correction of 6% [...] whereas the “[...] following Yeh and Hummer would have led to corrections of around 15%.” It is based on a linear regression in  $N^{-1/3}$  [...] ad hoc fit.

<sup>1</sup>G. Guevara Carrión, R. Fingerhut, J. Vrabec, «Fick diffusion coefficient matrix of a quaternary liquid mixture by molecular dynamics», *J. Phys. Chem. B* **124**(22): 4527–4535, doi:10.1021/acs.jpccb.0c01625, **2020**.

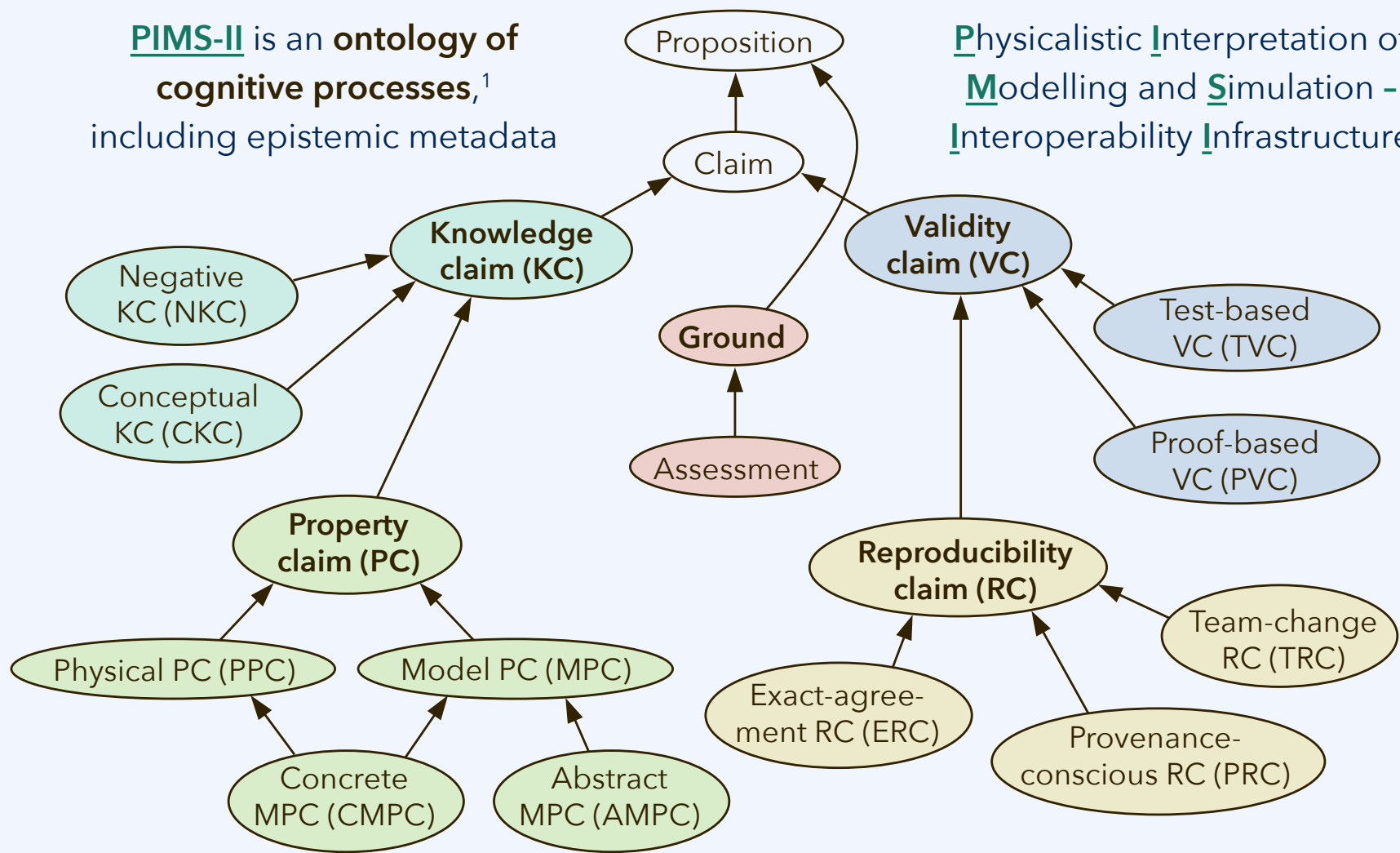
<sup>2</sup>M. T. Horsch, B. Schembera, «Epistemic metadata in molecular modelling: First-stage case-study report (10 cases)», Inprodat technical report 2023-A, doi:10.5281/zenodo.7516532, **2023**.

<sup>3</sup>M. Horsch, S. Chiacchiera, G. Guevara, M. Kohns, E. Müller, D. Šarić, S. Simon, I. Todorov, J. Vrabec, B. Schembera, «Epistemic metadata in molecular modelling: Second-stage case-study report (12 claims)», Inprodat technical report 2023-B, doi:10.5281/zenodo.7610237, **2023**.

# Ontology of epistemic metadata<sup>1</sup>

**PIMS-II** is an **ontology of cognitive processes**,<sup>1</sup> including epistemic metadata

**Physicalistic Interpretation of Modelling and Simulation - Interoperability Infrastructure**



<sup>1</sup>OWL implementation under <http://www.molmod.info/semantics/pims-ii.ttl>

# Ontology of epistemic metadata<sup>1, 2</sup>

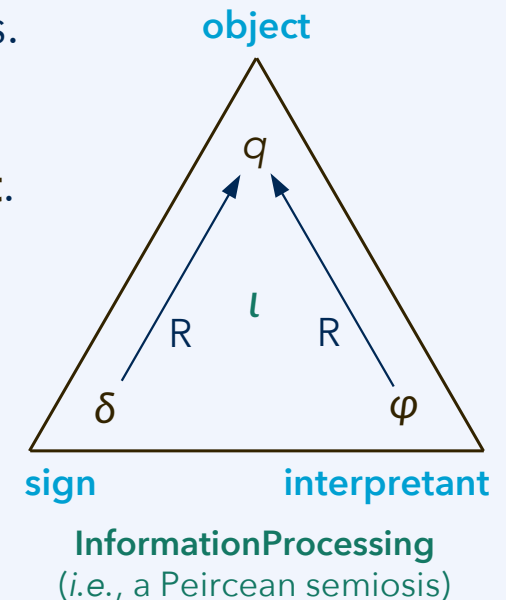
Peircean semiotics is applied to the description of cognitive processes, e.g., consider a process in which dataset  $\delta$  is analysed, yielding knowledge claim  $\varphi$ :

- The data  $\delta$  are about some research question  $q$ .  
So  $\delta$  is a representamen for  $q$ ; it has the role of the **sign**.
- The research question  $q$  is the **object** of the semiosis.
- As an outcome of the semiosis, claim  $\varphi$  is obtained, which is a new representamen for  $q$ , the **interpretant**.

The part of the PIMS-II ontology that deals with Peircean semiotics is also axiomatized in first-order logic,<sup>1</sup> in addition to the OWL implementation.<sup>2</sup>

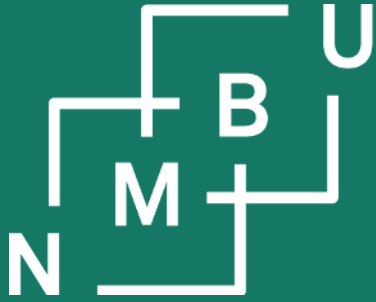


C. S. Peirce



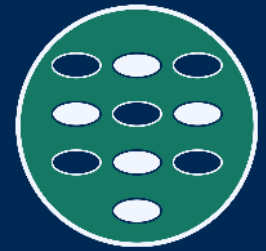
<sup>1</sup>«Mereosemiotics: Parts and signs», in *Proc. JOWO 2021 (FOUST)*, 2021.

<sup>2</sup>OWL implementation under <http://www.molmod.info/semantics/pims-ii.ttl>



Noregs miljø- og  
biovitenskaplege  
universitet

Materialteori og -informatikk



Digitalisering på Ås

# Domain and background

## First stage and an ontology

## Second stage and grounding

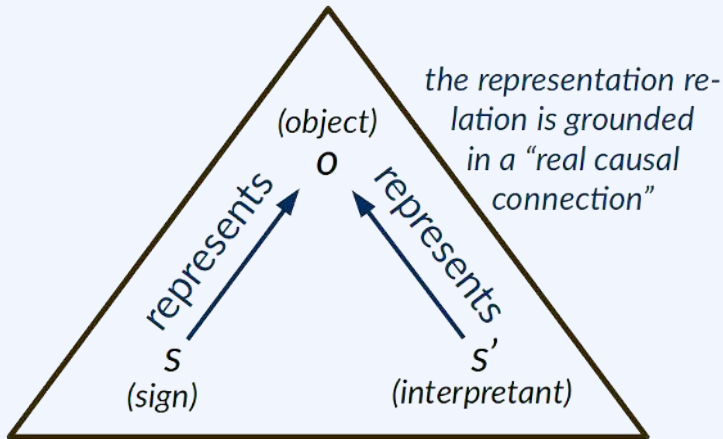
**Knowledge claim (KC), including the provenance**

«Researcher  $a$  did  $\kappa$  and found  $\varphi$  (and thus claims to know  $\varphi$ ).»

→ Therefore, when research process  $\kappa$  is carried out, it can lead to outcome  $\varphi$ .

# Grounding as implemented in PIMS-II (so far)

## Peircean semiotics

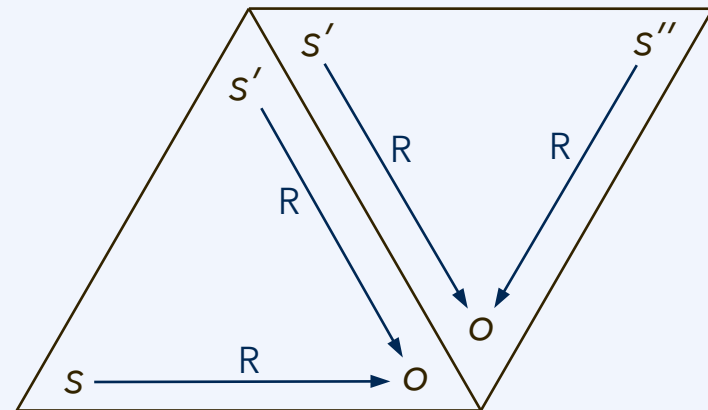


the semiosis, a process by which a new representamen, the interpretant, is created

## Cognitive process (example):

- First, experimental data  $s$  for material  $o$  are used to parameterize a model, obtaining model  $s'$ .
- Then, a simulation is done using model  $s'$ , yielding the simulation result  $s''$  (which also represents  $o$ ).

## Research workflows as cognitive processes:<sup>1</sup>



cognitive process  $K$

Each cognitive step starts from one representation relation, e.g.,  $R_{so}$ , and creates a new one,  $R_{s'o}$ .

The successor step reuses  $R_{s'o}$  and creates the next relation,  $R_{s''o}$ .

<sup>1</sup>«Mereosemiotics: Parts and signs», in *Proc. JOWO 2021 (FOUST)*, 2021.

# Guevara *et al.* (2020) claims:<sup>1</sup> Second-stage analysis<sup>2</sup>

**Interviews** were done with the authors; e.g., on 24<sup>th</sup> January 2023, two papers, among them Guevara *et al.*<sup>1</sup> (2020) were discussed in a 70-minutes meeting. Two of the three authors participated (Gabriela Guevara and Jadran Vrabec).

## «Why is it knowledge?

- Yeh & Hummer instead use a semiempirical correlation, relying on all sorts of properties, working with the end result which  $D$ .
- The new method is formally much simpler, relying only on  $N$ , and it works with the underlying quantity  $L$  which is more fundamental, rather than with the end outcome  $D$ .
- Also, linear behaviour of  $D$  in  $1/N^3$  was already claimed before by others, and not only for  $D$ , it is something like "community shared understanding". In particular, Yeh-Hummer also has  $1/N^3$  in it.

## Validation:

- Is it better than Yeh-Hummer? Really such a validation still needs to be done.»

<sup>1</sup>G. Guevara Carrión, R. Fingerhut, J. Vrabec, «Fick diffusion coefficient matrix of a quaternary liquid mixture by molecular dynamics», *J. Phys. Chem. B* **124**(22): 4527–4535, doi:10.1021/acs.jpccb.0c01625, **2020**.

<sup>2</sup>M. Horsch, S. Chiacchiera, G. Guevara, M. Kohns, E. Müller, D. Šarić, S. Simon, I. Todorov, J. Vrabec, B. Schembera, «Epistemic metadata in molecular modelling: Second-stage case-study report (12 claims)», Inprodat technical report 2023-B, doi:10.5281/zenodo.7610237, **2023**.



# Guevara *et al.* (2020) claims:<sup>1</sup> Second-stage analysis<sup>2</sup>

Interviews summarized in the second-stage report,<sup>2</sup> with two claims per paper.

Selected knowledge claims from the paper:

1. A novel finite-size correction methodology for the phenomenological diffusion coefficient matrix  $\mathbf{L}$  based on linear extrapolation over  $1/N^3$  to the limit  $1/N^3 \rightarrow 0$  is proposed and successfully used to calculate  $\mathbf{D}$ .
2. The Fick diffusion coefficient matrix  $\mathbf{D}$  of the considered mixture has the values given in Table 1 of the paper under the conditions specified there.

d The novel method looks preferable or more plausible as it exhibits what is typically seen in the community as theoretical virtues: *First*, Yeh and Hummer [21] use a semiempirical correlation relying on multiple properties, while the novel method is formally much simpler, relying only on  $N$ . *Second*, the Yeh-Hummer method operates on the end result  $\mathbf{D}$ , whereas the novel method operates on the intermediate result  $\mathbf{L}$  that directly experiences the finite-size limitation in the molecular simulation.

<sup>1</sup>G. Guevara Carrión, R. Fingerhut, J. Vrabec, «Fick diffusion coefficient matrix of a quaternary liquid mixture by molecular dynamics», *J. Phys. Chem. B* **124**(22): 4527–4535, doi:10.1021/acs.jpcc.0c01625, **2020**.

<sup>2</sup>M. Horsch, S. Chiacchiera, G. Guevara, M. Kohns, E. Müller, D. Šarić, S. Simon, I. Todorov, J. Vrabec, B. Schembera, «Epistemic metadata in molecular modelling: Second-stage case-study report (12 claims)», Inprodat technical report 2023-B, doi:10.5281/zenodo.7610237, **2023**.

# Type-1 and Type-2 grounding

Type-1 and Type-2 notions inspired by Marr.<sup>1, 2</sup>

|  |   |
|--|---|
| <b>Type-1</b>  |   |
| <b>Type-2</b><br>The <b>provenance</b> of the results tells that they are valid. | <i>Case study example: Šarić et al. argue:</i> <ul style="list-style-type: none"><li>– We are using ion models that worked accurately before.</li><li>– It was shown before that ion models designed for one water model still perform accurately for another water model.</li><li>– Therefore we can carry over the ion models designed for SPC/E water to another water model, namely, TIP4P/ε.</li></ul> |

<sup>1</sup>D. Marr, *Artificial Intelligence* 9(1): 37–48, doi:10.1016/0004-3702(77)90013-3, 1977.

<sup>2</sup>«Documentation of epistemic metadata by a mid-level ontology of cognitive processes», in *Proc. JOWO 2022*, CEUR vol. 3249: p. 2 (CAOS), CEUR-WS, 2022.

# Type-1 and Type-2 grounding

Type-1 and Type-2 notions inspired by Marr.<sup>1, 2</sup>

|  |   |
|--|---|
| <p><b>Type-1</b></p> <p>The results' validity is <u>not</u> grounded in the way the results were obtained.</p> | <p><i>Typical:</i> <b>Mathematical proof</b> in statistical mechanics for a theoretical framework with widely accepted definitions and axioms.</p> <p><i>Case study example:</i> Fingerhut <i>et al.</i> introduce a method based on Kirkwood-Buff integration (building on previous theoretical work by Ben Naim)</p>  |
| <p><b>Type-2</b></p> <p>The <b>provenance</b> of the results tells that they are valid.</p>                    | <p><i>Case study example:</i> Šarić <i>et al.</i> argue:</p> <ul style="list-style-type: none"><li>– We are using ion models that worked accurately before.</li><li>– It was shown before that ion models designed for one water model still perform accurately for another water model.</li><li>– Therefore we can carry over the ion models designed for SPC/E water to another water model, namely, TIP4P/ε.</li></ul> |

<sup>1</sup>D. Marr, *Artificial Intelligence* 9(1): 37–48, doi:10.1016/0004-3702(77)90013-3, 1977.

<sup>2</sup>«Documentation of epistemic metadata by a mid-level ontology of cognitive processes», in *Proc. JOWO 2022*, CEUR vol. 3249: p. 2 (CAOS), CEUR-WS, 2022.

# Normative grounds and reliabilism

|   |  |
|---|--|
| <p><b>Type-1</b></p> <p>The results' validity is not grounded in the way the results were obtained.</p> | <p><i>Typical:</i> <b>Mathematical proof</b> in statistical mechanics for a theoretical framework with widely accepted definitions and axioms.</p>   |
| <p><b>Type-2</b></p> <p>The <b>provenance</b> of the results tells that they are valid.</p>             | <p><b>Reliability</b> of process <math>m</math> means that «If <math>S</math>'s believing <math>p</math> at <math>t</math> results from <math>m</math>, then <math>S</math>'s belief in <math>p</math> at <math>t</math> is <b>justified</b>».<sup>1</sup></p> <p><i>Typical:</i> We used a <b>model</b>, <b>method</b>, and <b>simulation code</b> validated in the past and - usually - very accurate.</p> <p>(<b>process reliabilism</b>)</p> |

# Normative grounds and reliabilism

|   | «evidence supporting trustworthiness cannot be complete» <sup>1</sup><br><b>trust</b>  | «reliance is compatible with - ideally - the complete evidence» <sup>1</sup><br><b>reliance</b>   |
|---|--|---|
| <p><b>Type-1</b></p> <p>The results' validity is not grounded in the way the results were obtained.</p> | <p><i>Typical:</i> Mathematical proof in statistical mechanics for a theoretical framework with widely accepted definitions and axioms.</p>  | <p><b>Case study example:</b><br/>Guevara <i>et al.</i> argue:<br/>Our new finite-size correction is better because it is more simple and because it is applied to the directly computed quantity <b>L</b>.</p> |
| <p><b>Type-2</b></p> <p>The <b>provenance</b> of the results tells that they are valid.</p>             | <p><b>Case study example:</b><br/>Chatwell and Vrabec argue:<br/>It is OK to use a cutoff radius of <math>5.5\sigma</math> for the LJ potential, since this was done in three cited works from the literature.</p> | <p><i>Typical:</i> We used a model, method, and simulation code validated in the past and - usually - very accurate.<br/><b>(process reliabilism)</b></p>   |





See **reference ontology of trust** (ROT) by Baratella *et al.*<sup>1</sup>

<sup>1</sup>Baratella *et al.*, «The many facets of trust», to appear in *Proceedings of FOIS 2023*.

# Normative grounds and reliabilism

|   | trust  | reliance   |
|---|--|--|
| <p><b>Type-1</b></p> <p>The results' validity is not grounded in the way the results were obtained.</p> | <p><i>Typical:</i> Mathematical proof in statistical mechanics for a theoretical framework with widely accepted definitions and axioms.</p>  | <p><i>Schema:</i> A new theory is more reliable because it is simpler, covers more phenomena, or represents underlying physics. <b>(theoretical virtues)</b></p> |
| <p><b>Type-2</b></p> <p>The <b>provenance</b> of the results tells that they are valid.</p>             | <p><i>Case study example:</i><br/>Chatwell and Vrabec argue: It is OK to use a cutoff radius of <math>5.5\sigma</math> for the LJ potential, since this was done in three cited works from the literature.</p> | <p><i>Typical:</i> We used a model, method, and simulation code validated in the past and - usually - very accurate. <b>(process reliabilism)</b></p>            |

# Schema: Grounding of knowledge claims

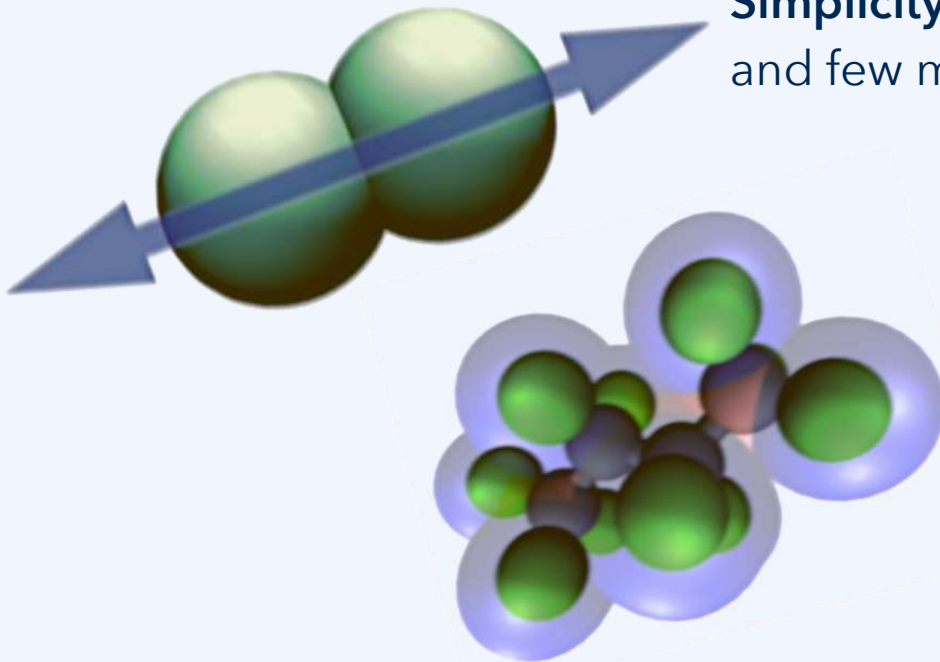
|        | trust   | reliance   |
|--------|---|--|
| Type-1 | rare<br>             | very frequent<br> |
| Type-2 | very frequent<br> | frequent<br>    |

**Sample:** Bowskill *et al.*, Chatwell & Vrabec, Fingerhut *et al.*, Guevara *et al.*, Stephan & Hasse (long-range).

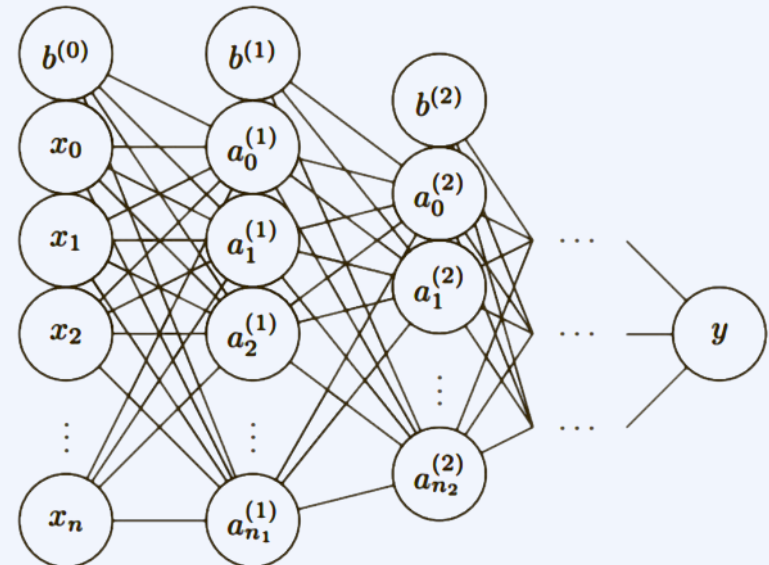
# Theoretical virtues

Theoretical virtues often oppose each other.

**Simplicity** favours few elements  
and few model parameters.



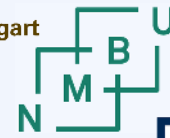
**Representing** all the physics  
**realistically** requires the opposite.



Neural networks are lacking in  
virtue, yet some people use them.

<sup>1</sup>Figure sources: MolMod DB (link) and Anna Jenul's doctoral thesis (link).

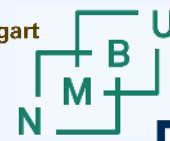




# Theoretical virtues

Theoretical virtues found in the case study:

- **Alignment** of representations, qualitative reflection of underlying physics
  - Chatwell & Vrabec: Relaxation model based on an exponential decay as deduced theoretically; functional form of the EOS made plausible by theory.
  - Zhu and Müller note that ML models are devoid of any physics-based insights.
- **Coverage** or good sampling of the phenomenon's state space
  - Bowskill *et al.*: Test cases are representative of real world problems.
  - Stephan and Hasse: Selected mixtures and conditions are representative.
- **Mechanism** or explainability of dependencies
  - Guevara *et al.*: Finite-size effect on immediate result (**L**), not end result (**D**).
- **Simplicity**
  - Guevara *et al.*: Finite-size correction depends on one quantity only (the size).



# Challenges at documenting grounding

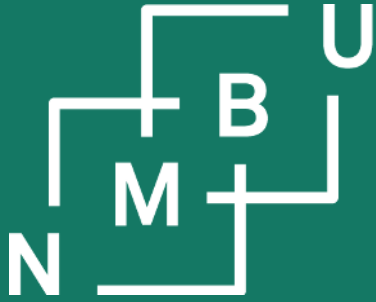
It is hard to differentiate between grounding in virtues or in process reliabilism.

- Model  $s'$  is better than  $s$ . It is equally accurate with fewer parameters,
  - ... so we should prefer it because it is **more simple**.
  - ... so, **from experience**, its extrapolations are **more reliable**.
    - » Nobody writes either of the above explicitly.

When digitalizing research data, we should respect that:

- Research is a **social process** among humans;
- scientific communication is **human communication**;
- it can rely on **pragmatics** – no need say every small thing explicitly;
- **epistemic grounding** is **usually not spelled out** in detail (or at all); we would often need to impute an interpretation onto the authors.

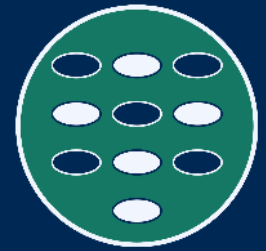
The ontology's aim is: Help people make more explicit statements **if they want**.



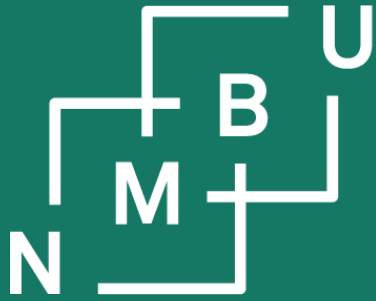
Noregs miljø- og  
biovitenskaplege  
universitet

Domain and background  
First stage and an ontology  
Second stage and grounding

Materialteori og -informatikk



Digitalisering på Ås



Norges miljø- og  
biovitenskapelige  
universitet



# Case study on epistemic metadata in molecular modelling

Martin Horsch,<sup>1,2</sup> Silvia Chiacchiera,<sup>2</sup> and Björn Schembera<sup>3</sup>

<sup>1</sup>Norwegian Univ. Life Sciences, Faculty of Science and Technology, Ås, Norway

<sup>2</sup>UKRI STFC Daresbury Laboratory, Scientific Computing Department, Daresbury, UK

<sup>3</sup>Univ. Stuttgart, Institute of Applied Analysis and Numerical Simulation, Stuttgart, Germany