

The Materials Design Ontology

Patrick Lambrix^{a,b,c,*}, Rickard Armiento^{c,d}, Huanyu Li^{a,c}, Olaf Hartig^a, Mina Abd Nikooie Pour^{a,c} and Ying Li^{a,c}

^a *Department of Computer Science, Linköping University, Sweden*

E-mails: patrick.lambrix@liu.se, huanyu.li@liu.se, olaf.hartig@liu.se, mina.abd.nikooie.pour@liu.se, ying.li@liu.se

^b *Department of Building Engineering, Energy Systems and Sustainability Science, University of Gävle, Sweden*

E-mail: patrick.lambrix@liu.se

^c *The Swedish e-Science Research Centre, Linköping University, Sweden*

E-mails: patrick.lambrix@liu.se, rickard.armiento@liu.se, huanyu.li@liu.se, mina.abd.nikooie.pour@liu.se, ying.li@liu.se

^d *Department of Physics, Chemistry and Biology, Linköping University, Sweden*

E-mail: rickard.armiento@liu.se

Abstract. In the materials design domain, much of the data from materials calculations is stored in different heterogeneous databases with different data and access models. Therefore, accessing and integrating data from different sources is challenging. As ontology-based access and integration alleviates these issues, in this paper we address data access and interoperability for computational materials databases by developing the Materials Design Ontology. This ontology is inspired by and guided by the OPTIMADE effort that aims to make materials databases interoperable and includes many of the data providers in computational materials science. In this paper, first, we describe the development and the content of the Materials Design Ontology. Then, we use a topic model-based approach to propose additional candidate concepts for the ontology. Finally, we show the use of the Materials Design Ontology by a proof-of-concept implementation of a data access and integration system for materials databases based on the ontology.¹

Keywords: Ontology, Ontology development, Data access, Data integration, Materials science, Materials Design Ontology

1. Introduction

Materials design and materials informatics is central for technological progress, not the least in the green engineering domain. Many traditional materials contain toxic or critical raw materials, whose use should be avoided or eliminated. Also, there is an urgent need for new environmentally friendly energy technologies. The design of viable materials with the right properties is a key component for enabling such technologies [3]. Computational materials design has contributed to recent progress in fields relevant to the move to eco-friendly solutions such as battery technologies and solar cells; other relevant examples of materials design for novel technologies include thermoelectrics and magnetic transport [4–6].

The space of potentially useful materials yet to be discovered — the so-called ‘chemical white space’ — is immense. The possible combinations of, say, up to six different elements, constitute many billions. The space is further extended by possibilities of different phases, low-dimensional systems, nanostructuring, and so forth, which

*Corresponding author. E-mail: patrick.lambrix@liu.se.

¹This paper is an extension of [1] with results from [2] and currently unpublished results regarding an application using the ontology.

adds several orders of magnitude. This space was traditionally explored by experimental techniques, i.e., materials synthesis and subsequent experimental characterization. Parsing and searching the full space of possibilities this way is, however, hardly practical. Recent advances in condensed matter theory and materials modeling make it possible to generate reliable materials data by means of computer simulations based on quantum mechanics [7]. High-throughput simulations combined with machine learning can speed up progress significantly and also help to break out of local optima in composition space to reveal unexpected solutions and new chemistries [8]. The progress brought by the combination of machine learning models and databases of materials data, is now so rapid that it can be discussed as a lead-up to a *singularity* for the field of materials design [9].

This development has led to several global efforts to assemble and curate databases that combine experimentally known and computationally predicted materials properties, along with a desire to make them interoperable [10]. These efforts have collectively been referred to as the Materials Genome Initiative (<https://www.mgi.gov/>). A central idea is that materials design challenges can be addressed by searching these databases for entries with desired combinations of properties. Nevertheless, these data sources also open up for *materials informatics*, i.e., the use of big data methodology and data mining techniques to discover new physics from the data itself. A workflow for such a discovery process can be based on a typical data mining process, where key factors are identified, reduced and extracted from heterogeneous databases, similar materials are identified by modeling and relationship mining and properties are predicted through evaluation and understanding of the results from the data mining techniques [11]. The use of the data in such a workflow requires addressing problems in data integration, provenance, and semantics.

Even when a new material has been invented and synthesized in a lab, much work remains before it can be deployed. Production methods allowing manufacturing the material at large scale in a cost effective manner need to be developed, and integration of the material into the production must be realized. Furthermore, life-cycle aspects of the material need to be assessed. Today, this post-invention process takes typically about two decades [5, 12]. Shortening this time is in itself an important strategic goal, which could be realized with the help of an integrated informatics approach [5].

It is clear that materials data, experimental as well as simulated, has the potential to speed up progress significantly in many steps in the chain starting with materials discovery, all the way to marketable product. However, the data needs to be suitably organized and easily accessible, which in practice is highly nontrivial to achieve. It requires a multidisciplinary effort and the various conventions and norms in use need to be integrated. Materials data is highly heterogeneous [12].

In this paper we address the data access and interoperability issue by developing an ontology suitable for the OPTIMADE (Open Databases Integration for Materials Design, <https://www.optimade.org/>) effort. The OPTIMADE consortium aims to make materials databases interoperable by developing a specification for a common REST API [13]. The consortium includes many of the data providers in computational materials science. However, although a first version of a common API has been defined, there is no semantic and integrated access support yet. Therefore, in this paper, we develop the Materials Design Ontology (MDO) that covers the content and is guided and inspired by the databases in OPTIMADE. The current focus is on solid-state physics and condensed matter theory. Further, we provide a proof-of-concept implementation of a data access and integration system that currently covers two of the databases in the OPTIMADE consortium. The used framework is general and in the future other databases as well as the OPTIMADE API will be added to the implementation.

The paper is organized as follows. In Section 3 we describe the development of MDO while the ontology itself is described in Section 4. In Section 5 we propose new concepts for an extension of MDO. Currently, these concepts are under discussion. In Section 6 we show the use of MDO in our MDO proof-of-concept implementation of a data access and integration system for materials science databases. The paper concludes in Section 7. We start with some background in Section 2.

2. Background

2.1. Ontologies in materials science

A number of ontologies in materials science have been developed. To find these ontologies, we used services such as BioPortal (<https://bioportal.bioontology.org/>), MatPortal (<https://matportal.org/>) and Linked Open Vocabularies

Table 1
Characteristics of some materials ontologies

Ontologies	Knowledge Representation Perspective			Materials Science Perspective	
	Ontology Metrics	Language	Modularity	Domain	Application Scenario
MatOnto [14]	78 concepts, 10 relations, 24 instances	OWL	✓	Crystals	Materials discovery
MatOWL [15]	(not available)	OWL		Materials	Semantic querying
Materials Ontology [16]	606 concepts, 31 relations, 488 instances	OWL	✓	Thermal properties	Data exchange, search
ELSSI-EMD ontology [17]	35 concepts, 37 relations, 33 instances	OWL	✓	Materials testing	Standardization
NanoParticle Ontology [18]	1904 concepts, 81 relations	OWL		Nanotechnology	Data integration, search
eNanoMapper [19]	12781 concepts, 5 relations 464 instances	OWL	✓	Nanotechnology	Data integration
MMOY [20]	2325 concepts, 9 relations, 1738 instances	OWL		Metals	Knowledge extraction
MAMBO [21]	26 concepts, 33 relations	OWL	✓	Molecules-based materials	Knowledge representation
Dislocation Ontology [22]	18 concepts, 16 relations	OWL	✓	Crystalline Materials	Knowledge representation
PMD [23]	13 concepts, 7 relations	OWL	✓	Materials experiments	Knowledge representation, Data curation
MDO	37 concepts, 64 relations	OWL	✓	Materials design	Semantic querying over multiple databases

(<https://lov.linkeddata.es/dataset/lov/>) and search engines such as Google. We then conducted a literature review regarding these ontologies to find out about their characteristics.

There are a number of top-level ontologies that are interesting for conceptualization in the materials science domain. For instance, these top-level ontologies commonly contain definitions relevant to *continuants* and *occurents*. The former can represent materials objects that endure over time, and undergo a variety of changes, while the latter can represent events that unfold over time, and manifest as changes in continuants [24]. Examples of such top-level ontologies include EMMO (earlier known as European Materials & Modelling Ontology, and recently renamed Elementary Multiperspective Material Ontology, <https://github.com/emmo-repo/EMMO>) an ontology based on physics and analytical philosophy, in particular mereotopology and semiotics, the Basic Formal Ontology (BFO) [25], the Descriptive Ontology for Linguistic and Cognitive Engineering (DOLCE) [26], and the General Formal Ontology (GFO) [27].

Most ontologies, however, are domain ontologies, for which we show some characteristics from the knowledge representation and the materials science perspectives in Table 1. These ontologies focus on specific sub-domains of the materials field (Domain in Table 1) and have been developed with a specific use in mind (Application Scenario in Table 1). MatOnto [14], based on DOLCE, aims to represent structured knowledge, properties and processing steps relevant to materials for data exchange, reuse and integration. MatOWL [15] is extracted from MatML schema data to enable ontology-based data access. MatML ([28], <https://www.matml.org>), is an extensible markup language (XML) for exchanging materials information. The Materials Ontology in [16] was designed for data exchange among thermal property databases, particularly focusing on representing knowledge relevant to material processing, measurement methods and manufacturing processes. The NanoParticle Ontology [18], based on BFO, and the eNanoMapper ontology [19] are two ontologies in the nanotechnology domain. The former represents properties of nanoparticles to design new nanoparticles, while the latter focuses on assessing risks caused by the use of nanomaterials in engineering. Extensions to these ontologies are computed in [29]. The MMOY ontology [20] captures metal

materials knowledge from Yago ([30], <https://yago-knowledge.org>), a large knowledge base on many topics including materials and properties. The Materials and Molecules Basic Ontology (MAMBO, [21]) focuses on concepts and relations emerging on materials where the relationship between individual molecules and molecular aggregation is relevant to the properties of the system, such as in molecular materials and nanomaterials. MAMBO integrates with EMMO, Chemical Entities of Biological Interest (CheBI, [31]) and MDO. The Dislocation Ontology [22] focuses on representing knowledge related to crystalline materials and reuses some concepts in MDO. The Platform MaterialDigital Ontology (PMD, [23]) is a prototype to describe materials science experiments. The Materials Design Ontology (MDO, [1], <https://w3id.org/mdo/>), which is the focus of this paper, is inspired by OPTIMADE, and aims to enable semantic and integrated querying over multiple heterogeneous materials databases such as Materials Project [5], Open Quantum Materials Database (OQMD) [32], Novel Materials Discovery (NOMAD) [33] and Automatic FLOW for Materials Discovery (AFLOW) [34].

From the knowledge representation perspective, the basic terms defined in materials ontologies involve materials, properties, performance, and processing in specific sub-domains. All presented ontologies use OWL as a representation language (Language in Table 1). The number of OWL classes ranges from a few to several thousands (Ontology Metrics in Table 1). Some ontologies have more concepts than relations (e.g., MatOnto, Materials Ontology, NanoParticle Ontology, MMOY and EMMO), while some have many more relations (e.g., MDO). Several ontologies are developed in a modular fashion (Modularity in Table 1).

2.2. Ontology development

In Section 3 we describe the development of the Materials Design Ontology (MDO). Although, we could have used a more modern approach such as the eXtreme Design methodology [35] or its extension that integrates debugging [36], as our initial ontology was expected to be of a smaller size and given our earlier experience with the NeOn methodology for ontology engineering, we decided to use NeOn.

NeOn [37] is a methodology for ontology engineering that proposes nine scenarios which commonly occur, including Scenario 1: From Specification to Implementation, Scenario 2: Reusing and re-engineering non-ontological resources, Scenario 3: Reusing ontological resources, Scenario 4: Reusing and re-engineering ontological resources, Scenario 5: Reusing and merging ontological resources, Scenario 6: Reusing, merging, and re-engineering ontological resources, Scenario 7: Reusing ontology design patterns (ODPs), Scenario 8: Restructuring ontological resources, and Scenario 9: Localizing ontological resources. Depending on different background knowledge resources and purposes of the ontology, developers can make use of different scenarios or combinations of the scenarios. Scenario 1 is necessary in any ontology development and should always be included. The detailed use of NeOn for the development of MDO is described in Section 3.

Further, we also used two tools for detecting defects in the ontology during the development. The first tool, Ontology Pitfall Scanner! (OOPS!, [38]), helps to detect some of the most common pitfalls appearing within ontology development. The second tool, Repairing Ontological Structure Environment (RepOSE, [39]), allows to debug an ontology and proposes additional knowledge that could be interesting to add to the ontology.

2.3. Ontology extension

In Section 5 we describe work on generating new concepts that may be added to MDO. The new concepts are, however, not yet included in the public version of MDO as discussions regarding the scope and the use of the extension are ongoing.

We used the phrase-based topic model generation approach we presented in [29], shown in Figure 1. A topic model is a statistical model for discovering the abstract "topics" that occur in a collection of documents. The topics are often represented as lists of words or phrases. Given a corpus of documents related to the domain of interest and the number of requested topics, a phrase-based topic model is created using an extended version of the ToPMine [40] system as presented in [2].

First, frequent contiguous phrases are mined, which consists of collecting aggregate counts for all contiguous words satisfying a user-defined minimum support threshold. Given a minimum support threshold *min_support*, we say that phrases that occur at least *min_support* times are *frequent phrases*. The ToPMine system identifies frequent

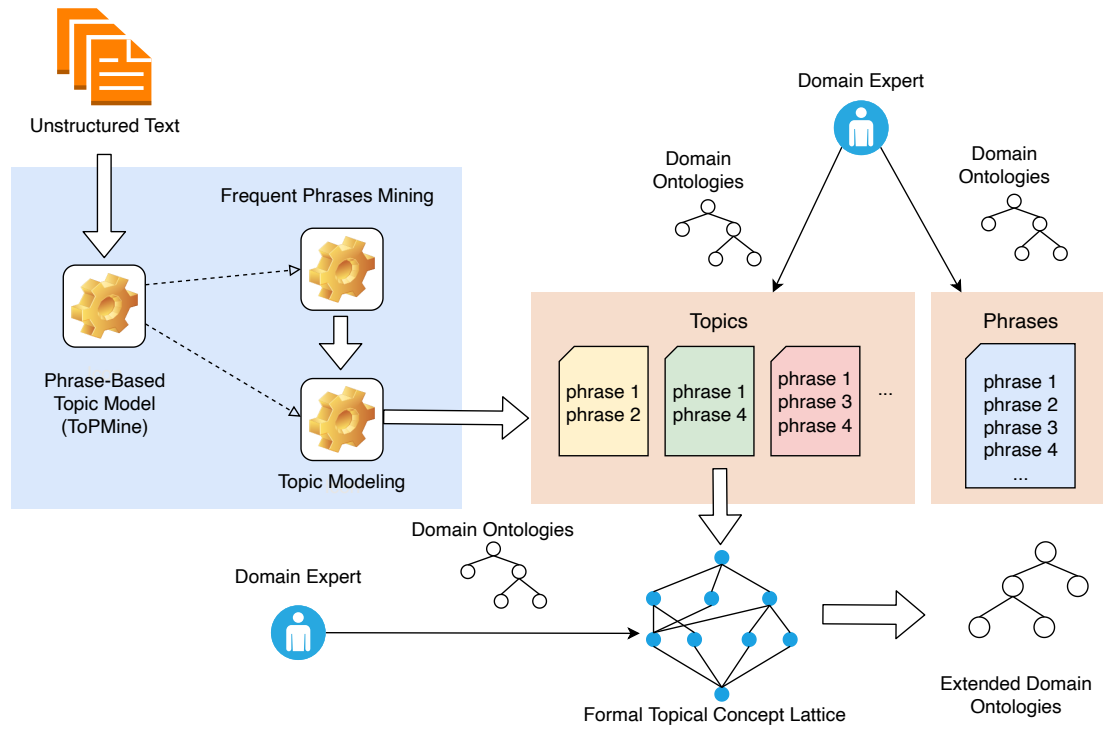


Figure 1. Approach: The upper part of the figure shows the creation of a phrase-based topic model with unstructured text as input and phrases and topics as output. The lower part shows the formal topical concept analysis with as input topics and as output a topical concept lattice. In both parts a domain expert validates and interprets the results [29].

phrases of a length up to a maximum length that is given as an input parameter. Note that ToPMine does not generate all frequent phrases but uses a method based on partitioning documents and using a significance score for deciding which words likely belong together, to produce high-quality frequent phrases. The use of the significance score has, for instance, the effect that frequent sub-phrases of frequent phrases are not counted. In our extended version we also introduce a user-defined maximum threshold for word occurrences to, if so desired, remove very general words. Then the documents are segmented based on the frequent phrases. Further, an agglomerative phrase construction algorithm merges the frequent phrases guided by a significance score.

After this phrase mining, the system performs topic modeling by computing representations of latent topics in the documents. Topics are generated using a variant of Latent Dirichlet Allocation (LDA) [41], called PhraseLDA, that deals with phrases, rather than words. Essentially, topics can be seen as a probability distribution over words or phrases.

The phrases as well as the topics are suggestions that a domain expert should validate or interpret and relate to concepts in the ontology. Based on the validations and interpretations of the domain expert, concepts and axioms are added to the ontology. To help a domain expert with the validation, we implemented a tool of which an early version is described in [42]. The current tool deals with phrases, but not yet with topics. It is available at <https://github.com/LiUsemWeb/phrase2onto>. The tool follows the iterative workflow of validation supporting the different phases. In the *Phrases* phase, the domain expert can look for sub- and super-phrases in the list as well as use string matching to find existing concepts in the ontology with similar names (Figure 2). This is useful to obtain an overview of the phrases as well as to find initial connections to the concepts in the ontology. In the *From Phrase to Concept* phase, the domain expert processes the frequent phrases and decides for each phrase whether one or more concepts related to the phrase can be defined in the ontology. The tool supports the addition of these concepts (Figure 3). In the *Concept* phase new concepts (also unrelated to the phrases) can be added (Figure 4). In our experiments the domain experts requested such functionality as they often thought of related concepts during the validation process. Finally, in the *From Concept to Axiom* phase, support for adding axioms related to the concepts is provided (Figure 5).

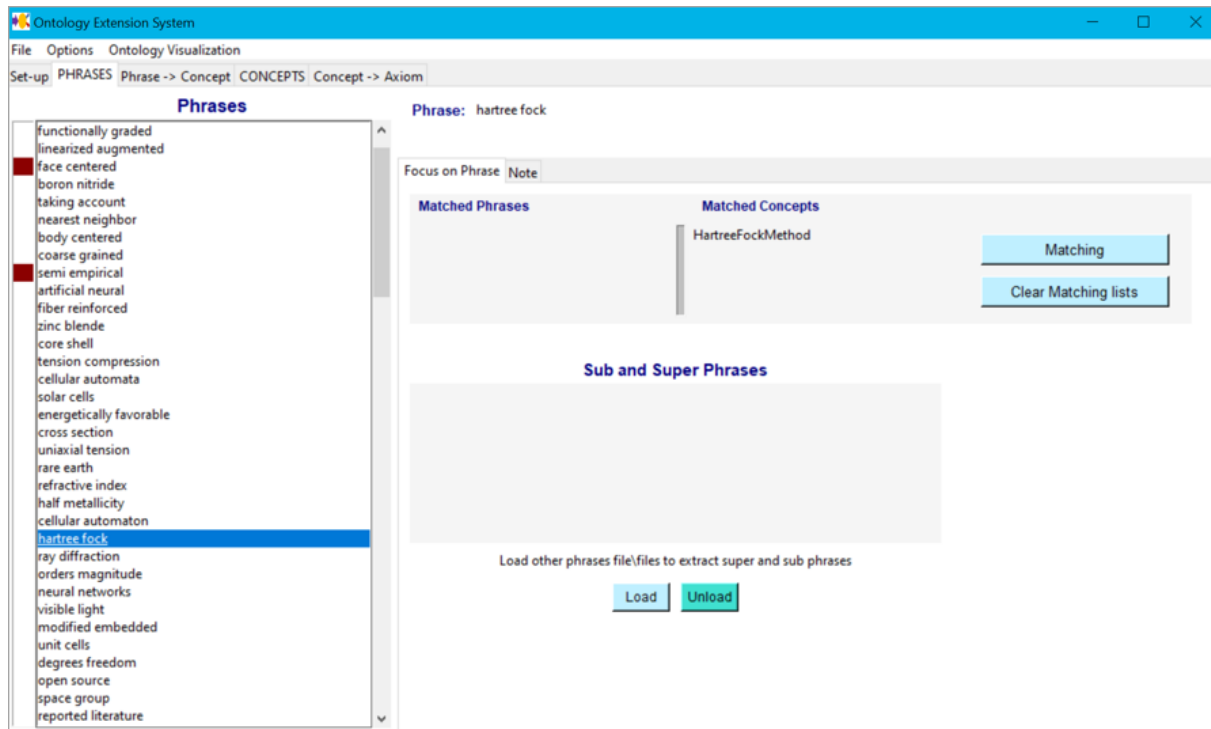


Figure 2. Tool - Phrases tab.

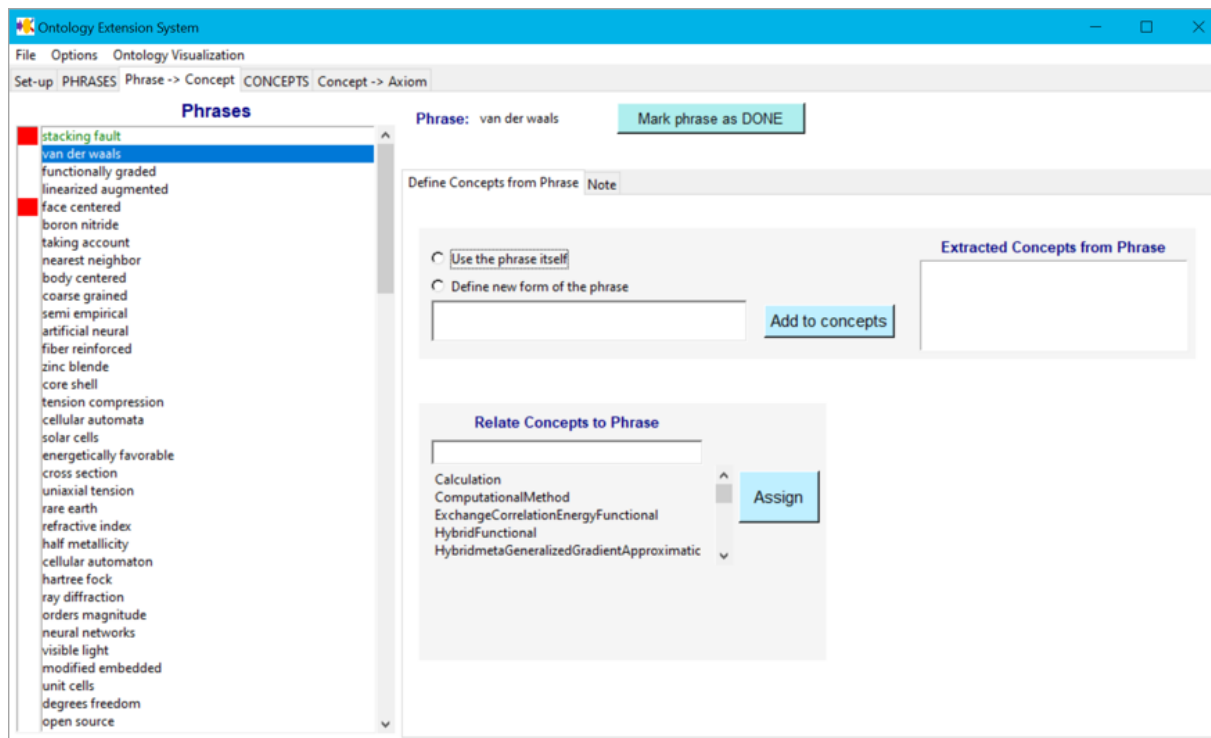


Figure 3. Tool - From Phrase to Concept tab.

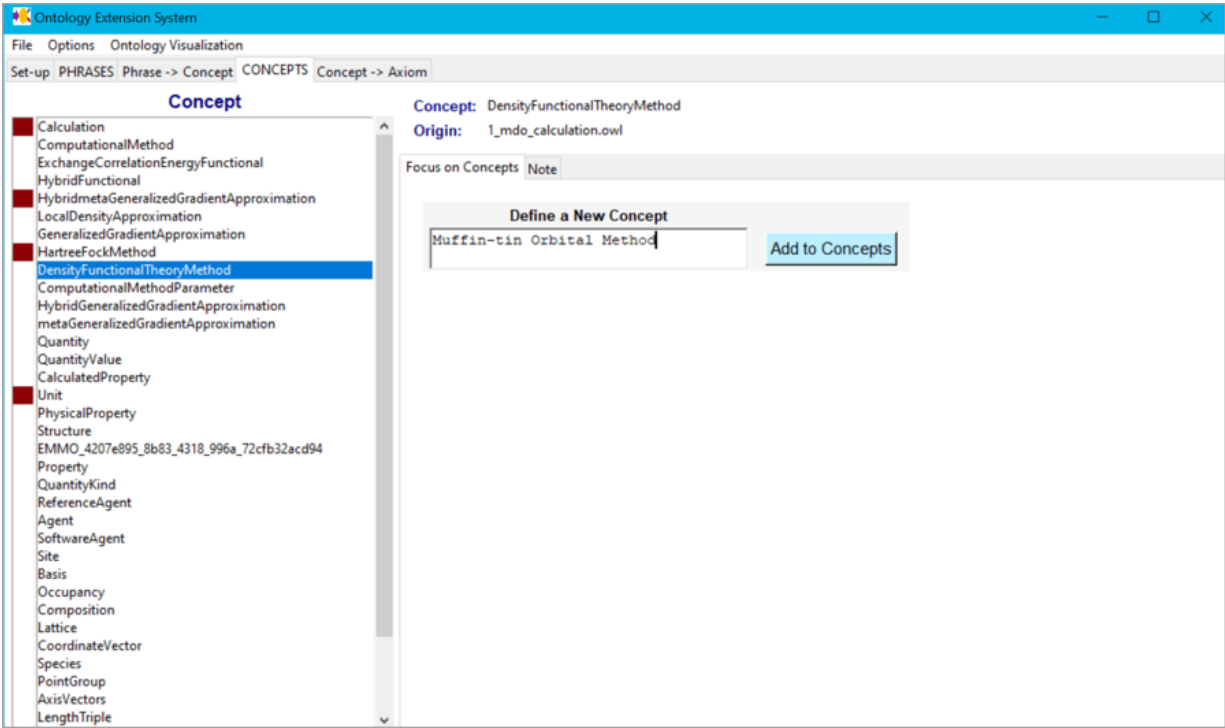


Figure 4. Tool - Concepts tab.

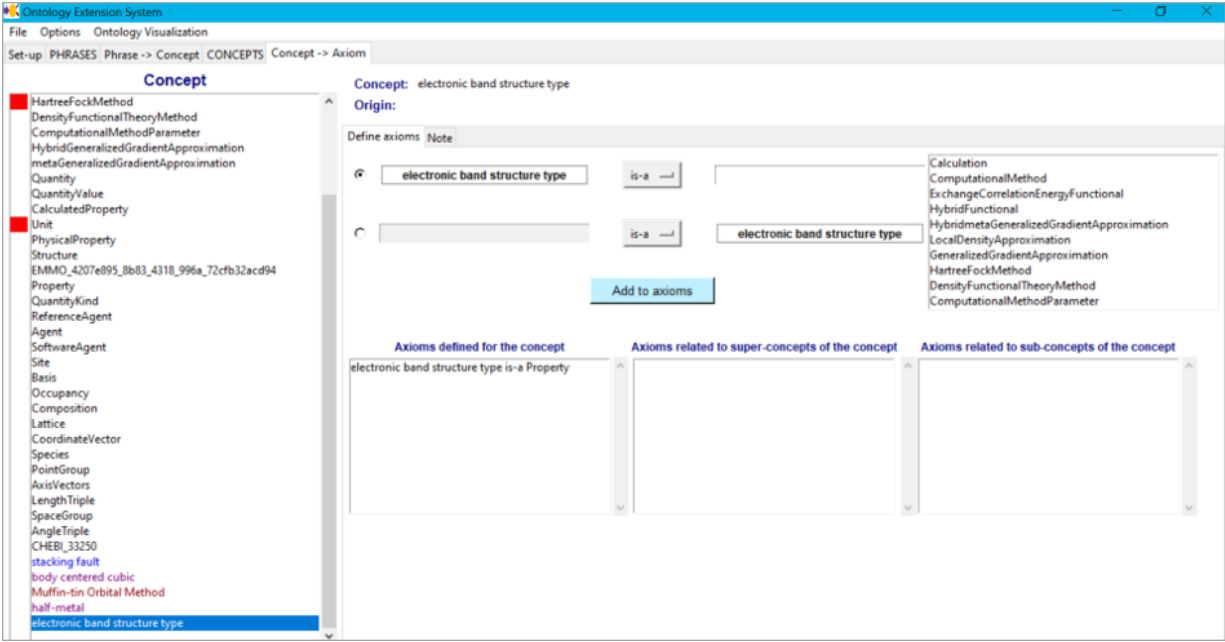


Figure 5. Tool - From Concept to Axiom tab.

3. Developing MDO

The development of MDO followed the NeOn ontology engineering methodology [37]. We focused on applying scenario 1 (*From Specification to Implementation*), scenario 2 (*Reusing and re-engineering non-ontological resources*), scenario 3 (*Reusing ontological resources*) and scenario 8 (*Restructuring ontological resources*). We did not re-engineer or merge ontological resources, i.e., scenarios 4-6 (but did reuse some concepts from other ontologies) and did not translate MDO into another natural language (scenario 9). Further, we did not use existing ontology design patterns (scenario 7), as the only one we are aware of in the materials science field is about materials transformation [43] that is not covered by MDO.

We used OWL2 DL as the representation language for MDO. During the whole process, two knowledge engineers, and one domain expert from the materials design domain were involved. In the remainder of this section, we introduce the key aspects of the development of MDO.

3.1. Requirements analysis

During this step, we clarified the requirements by proposing Use Cases (UC), Competency Questions (CQ) and additional restrictions (AR).

The use cases, which were identified through literature study and discussion between the domain expert and the knowledge engineers based on experience with the development of OPTIMADE and the use of materials science databases, are listed below.

- UC1: MDO will be used for representing knowledge in basic materials science such as solid-state physics and condensed matter theory.
- UC2: MDO will be used for representing materials calculation and standardizing the publication of the materials calculation data.
- UC3: MDO will be used as a standard to improve the interoperability among heterogeneous databases in the materials design domain.
- UC4: MDO will be mapped to OPTIMADE's schema to improve OPTIMADE's search functionality.

The competency questions are based on discussions with domain experts and contain questions that the databases currently can answer as well as questions that experts would want to ask the databases. For instance, CQ1, CQ2, CQ6, CQ7, CQ8 and CQ9 cannot be asked explicitly through the database APIs, although the original downloadable data contains the answers.

- CQ1: What are the calculated properties and their values produced by a calculation?
- CQ2: What are the input and output structures of a materials calculation?
- CQ3: What is the space group type of a structure?
- CQ4: What is the lattice type of a structure?
- CQ5: What is the chemical formula of a structure?
- CQ6: For a series of calculations, what are the compositions of materials with a specific range of a calculated property (e.g., band gap)?
- CQ7: For a specific material and a given range of a calculated property (e.g., band gap), what is the lattice type of the structure?
- CQ8: For a specific material and an expected lattice type of output structure, what are the values of calculated properties of the calculations?
- CQ9: What is the computational method used in a materials calculation?
- CQ10: What is the value for a specific parameter (e.g., cutoff energy) of the method used for the calculation?
- CQ11: Which software produced the result of a calculation?
- CQ12: Who are the authors of the calculation?
- CQ13: When was the calculation data published to the database?

Further, we proposed a list of additional restrictions that help in defining concepts.

- AR1: A property can relate to a structure.

- AR2: A calculation has exactly one corresponding computational method.
- AR3: A structure corresponds to one specific space group.
- AR4: A calculation is performed by some software programs or codes.
- AR5: A structure is a part of some materials.
- AR6: A structure and a property can be published by references which could be databases or publications.
- AR7: A calculation can take some structures as input.
- AR8: A calculation can take some properties as input.

3.2. Reusing and re-engineering non-ontological resources

To obtain the knowledge for building the ontology, we followed two steps: (1) the collection and analysis of non-ontological resources that are relevant to the materials design domain, and (2) discussions with the domain expert regarding the concepts and relationships to be modeled in the ontology. The collection of non-ontological resources comes from: (1) the dictionaries of the Crystallographic Information Framework (CIF, <https://www.iucr.org/resources/cif>) and International Tables for Crystallography (<https://it.iucr.org/>); (2) the APIs from different databases (e.g., Materials Project, AFLOW, OQMD) and OPTIMADE.

3.3. Connection and integration of existing ontologies

We reuse the concepts ‘Agent’ and ‘SoftwareAgent’ from PROV-O [44]. In terms of representation of units we reuse the ‘Quantity’, ‘QuantityValue’, ‘QuantityKind’ and ‘Unit’ concepts from QUDT (Quantities, Units, Dimensions and Data Types Ontologies) [45]. We use the metadata terms from the Dublin Core Metadata Initiative (DCMI, <http://purl.org/dc/terms/>) to represent the metadata of MDO.

4. Description of MDO

MDO consists of one basic module, *Core*, and two domain-specific modules, *Structure* and *Calculation*, importing the *Core* module. In addition, the *Provenance* module, which also imports *Core*, models provenance information. In total, the OWL2 DL representation of the ontology contains 37 concepts, 32 object properties, and 32 data properties. Figure 6 shows an overview of the ontology. The ontology specification is also publicly accessible at w3id.org at <https://w3id.org/mdo/full/1.1/>. The competency questions can be answered using the concepts and relations in the different modules (CQ1 and CQ2 by *Core*, CQ3 to CQ8 by *Structure*, CQ9 and CQ10 by *Calculation*, and CQ11 to CQ13 by *Provenance*).

The **Core** module, shown in Figure 7, consists of the top-level concepts and relations of MDO, which are also reused in other modules. Figure 8 shows the description logic axioms for the *Core* module. The module represents general information of materials calculations. The concepts *Calculation* and *Structure* represent materials calculations and materials’ structures, respectively, while *Property* represents a quantifiable aspect of one material or materials system and is defined as a sub-concept of *Quantity* from QUDT (Core4). *Property* is specialized into the disjoint concepts *CalculatedProperty* and *PhysicalProperty* (Core1, Core2, Core3).² In our view, a *CalculatedProperty* is the result of a calculation, i.e., the outcome of executing a specific computational method. In contrast, a *PhysicalProperty* refers to a specific objective property over a physical system, i.e., something that at least in principle can be measured in an experiment. A *PhysicalProperty* does not represent the outcome of such a measurement, only the abstract notion of the property. For example, *OpticalBandGap* is an instance of *PhysicalProperty*. An approximation can be computed by using a *DensityFunctionalTheoryMethod* in the **Calculation** module using Heyd-Scuseria-Ernzerhof (HSE), an instance of *HybridGeneralizedGradientApproximation*. This approximation can be represented by *HSEKohn-ShamBandGap*, an instance of *CalculatedProperty*.

²It would be possible to extend MDO with a *MeasuredProperty* to represent the outcome of a specific experiment as a parallel concept to *CalculatedProperty*. However, since the current version of MDO does not aim to describe experiments, this was determined to be outside of its scope at the moment.

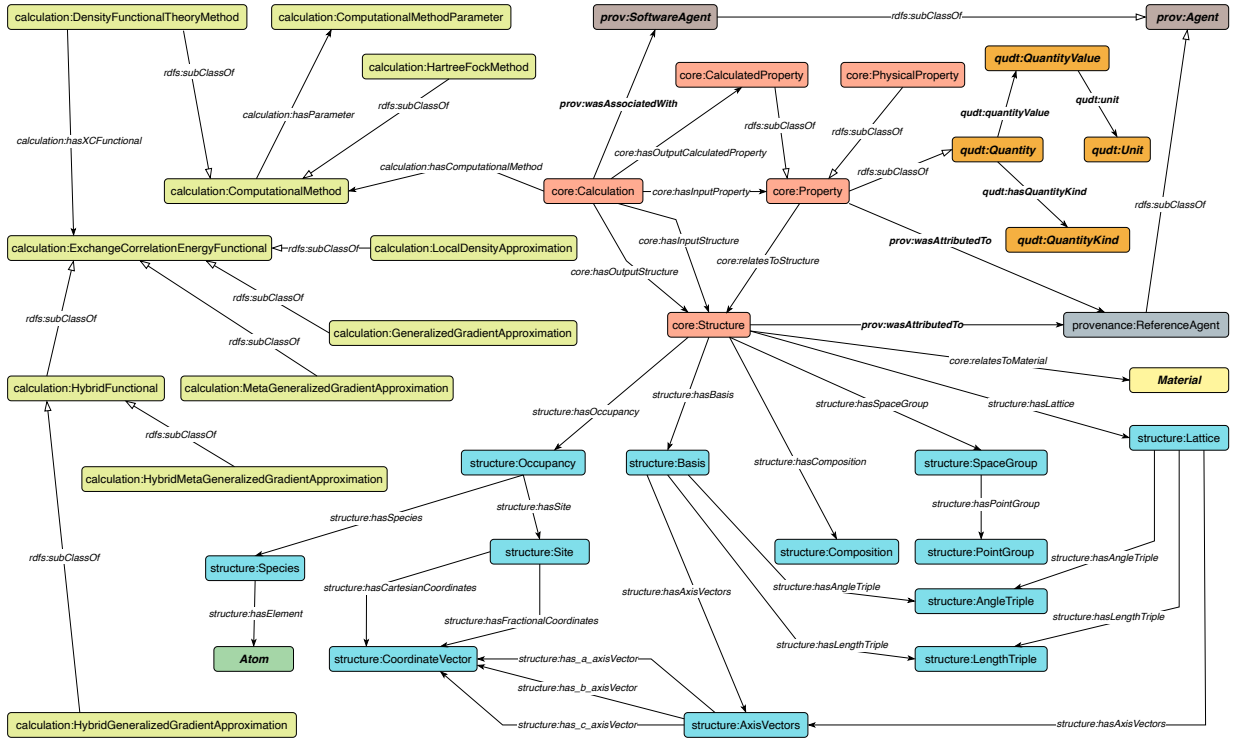


Figure 6. An overview of MDO.

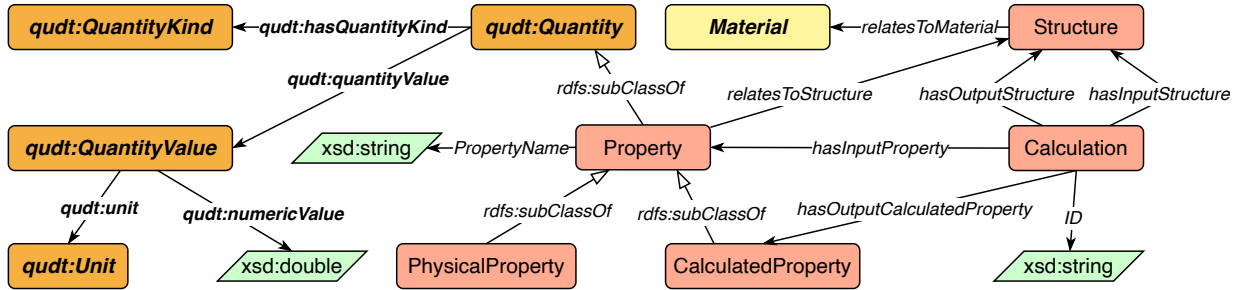


Figure 7. Concepts and relations in the Core module.

Properties are also related to *structures* (Core5). When a calculation is applied on materials structures, each calculation takes some structures and properties as input, and may output structures and calculated properties (Core6, Core7). Further, we define a concept *Material* and state that each structure is related to some material (Core8).

The **Structure** module, shown in Figure 9, represents the structural information of materials. Figure 10 shows the description logic axioms for the *Structure* module. Each structure has exactly one composition which represents what chemical elements compose the structure and the ratio of elements in the structure (Struc1). The composition has different representations of chemical formulas. The occupancy of a structure relates the sites with the species, i.e., the specific chemical elements, that occupy the site (Struc2 - Struc5). Each site has at most one representation of coordinates in Cartesian format and at most one in fractional format (Struc6, Struc7). The spatial information regarding structures is essential to reflect physical characteristics such as melting point and strength of materials. To represent this spatial information, we state that each structure is represented by some bases and a (periodic) structure can also be represented by one or more lattices (Struc8). Each basis and each lattice can be identified by

(Core1) $\text{CalculatedProperty} \sqsubseteq \text{Property}$
 (Core2) $\text{PhysicalProperty} \sqsubseteq \text{Property}$
 (Core3) $\text{CalculatedProperty} \sqcap \text{PhysicalProperty} \sqsubseteq \perp$
 (Core4) $\text{Property} \sqsubseteq \text{Quantity}$
 (Core5) $\text{Property} \sqsubseteq \forall \text{ relatesToStructure.Structure}$
 (Core6) $\text{Calculation} \sqsubseteq \exists \text{ hasInputStructure.Structure} \sqcap \forall \text{ hasInputStructure.Structure} \sqcap \forall \text{ hasOutputStructure.Structure}$
 (Core7) $\text{Calculation} \sqsubseteq \exists \text{ hasInputProperty.Property} \sqcap \forall \text{ hasInputProperty.Property} \sqcap \forall \text{ hasOutputCalculatedProperty.CalculatedProperty}$
 (Core8) $\text{Structure} \sqsubseteq \exists \text{ relatesToMaterial.Material} \sqcap \forall \text{ relatesToMaterial.Material}$

Figure 8. Description logic axioms for the Core module.

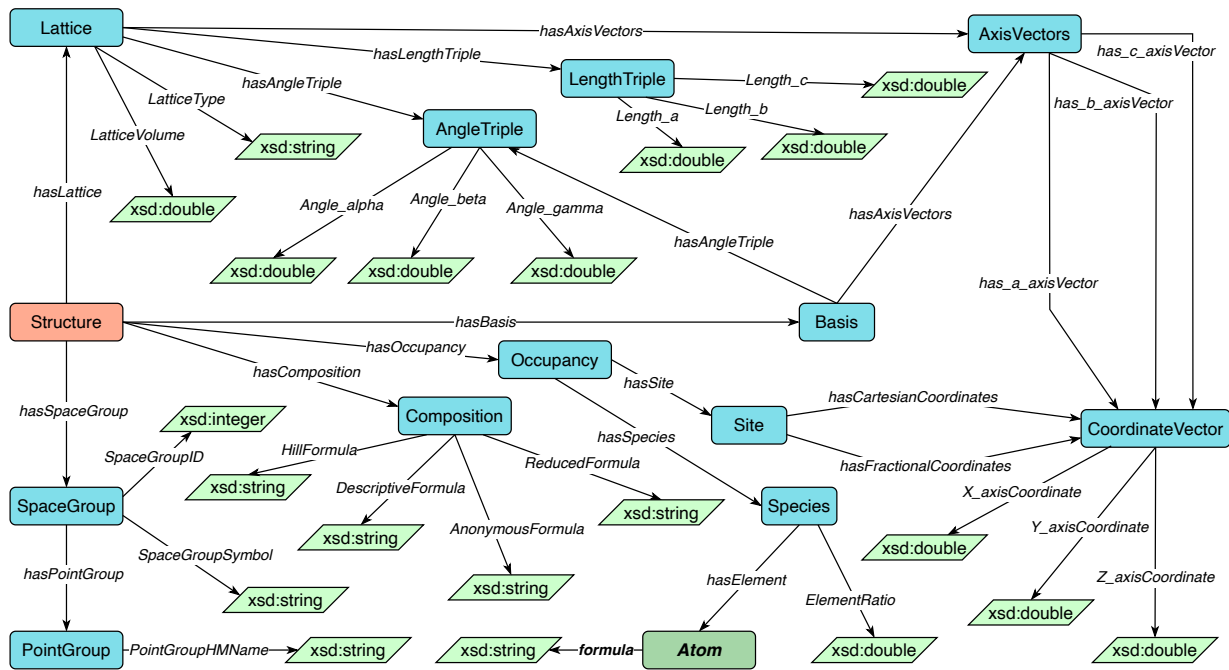


Figure 9. Concepts and relations in the Structure module.

one *axis-vectors* set or one *length triple* together with one *angle triple* (Struc9, Struc10). An *axis-vectors* set has three connections to *coordinate vector* representing the coordinates of three translation vectors respectively, which are used to represent a (minimal) repeating unit (Struc11). These three translation vectors are often called a, b, and c. Point groups and space groups are used to represent information of the symmetry of a structure. The *space group* represents a symmetry group of patterns in three dimensions of a *structure* and the *point group* represents a group of linear mappings which correspond to the group of motions in space to determine the symmetry of a *structure*. Each *structure* has one corresponding *space group* (Struc12). Based on the definition from International Tables for Crystallography, each *space group* also has some corresponding *point groups* (Struc13).

The **Calculation** module, shown in Figure 11, represents the classification of different computational methods. Figure 12 shows the description logic axioms for the *Calculation* module. Each *calculation*³ is achieved by a specific *computational method* (Cal1). Each *computational method* has some *parameters* (Cal2). In the latest version (v1.1)

³At the moment we have not considered calculations consisting of different steps, but each step would be a calculation on its own. Dealing with more complex representations is left for future work.

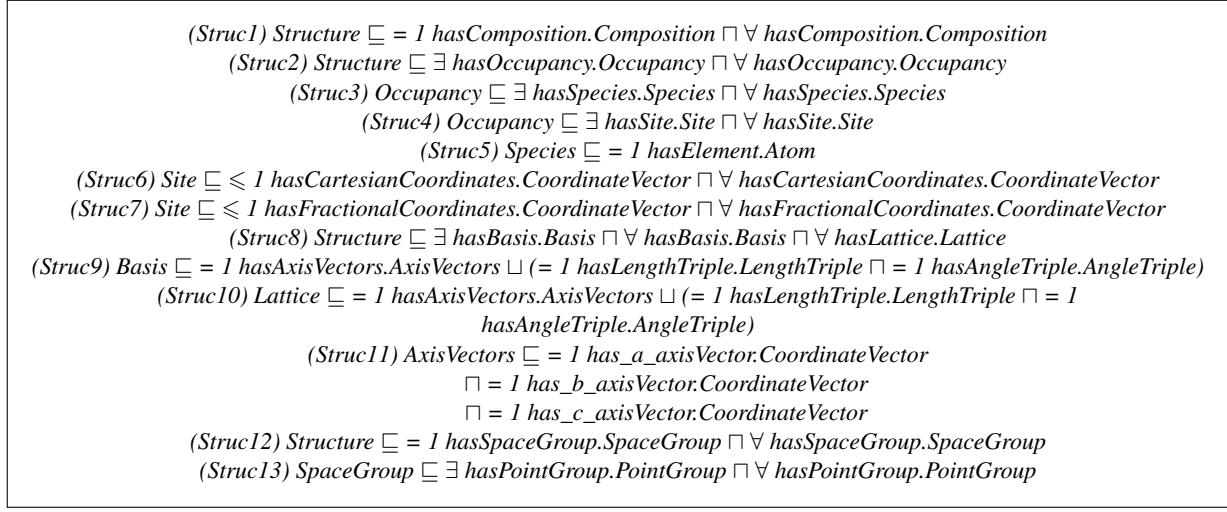


Figure 10. Description logic axioms for the Structure module.

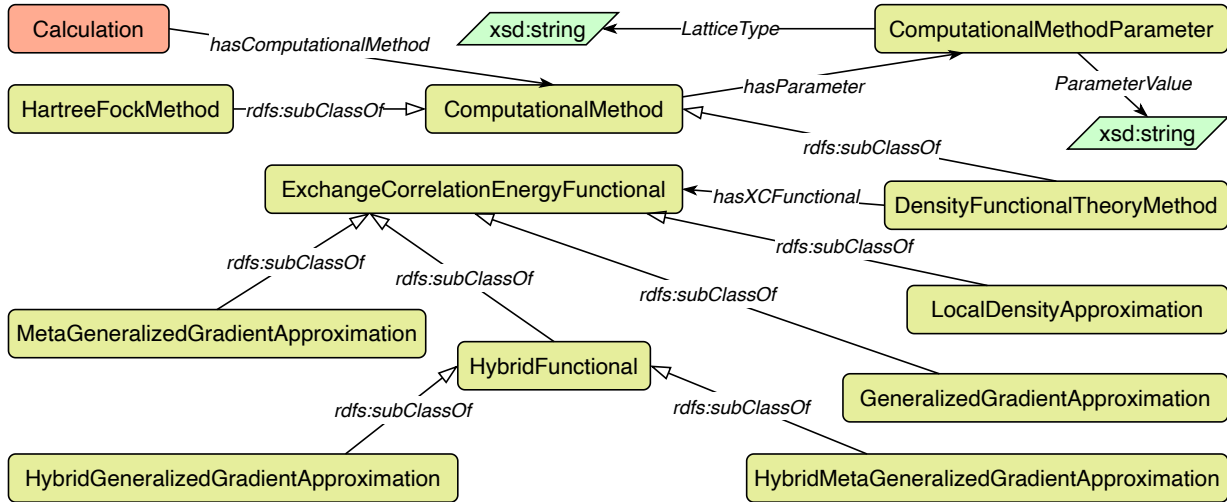


Figure 11. Concepts and relations in the Calculation module.

of this module, we represent two different methods, the *density functional theory method* and the *HartreeFock method* (Cal3, Cal4). In particular, the density functional theory method is frequently used in materials design to investigate the electronic structure. Such method has at least one corresponding *exchange correlation energy functional* (Cal5) which is used to calculate the exchange–correlation energy of a system. There are different kinds of functionals to calculate exchange–correlation energy (Cal6 - Cal11).

The **Provenance** module, shown in Figure 13, represents the provenance information of materials data and calculation. Figure 14 shows the description logic axioms for the *Provenance* module. We reuse part of PROV-O and define a new concept *ReferenceAgent* as a sub-concept of PROV-O's agent (Prov1). We state that each *structure* and *property* can be published by *reference agents* which could be databases or publications (Prov2, Prov3). Each *calculation* is produced by a specific *software* (Prov4).

In Figure 15 we exemplify the use of MDO to represent a specific materials calculation and related data in an instantiation. The example is from one of the 85 stable materials published in Materials Project in [46]. The calculation is about one kind of elpasolites, with the composition $Rb_2Li_1Ti_1Cl_6$. To not overcrowd the figure, we

(Cal1) $\text{Calculation} \sqsubseteq = 1 \text{ hasComputationalMethod.ComputationalMethod}$
 (Cal2) $\text{ComputationalMethod} \sqsubseteq \exists \text{ hasParameter.ComputationalMethodParameter} \sqcap \forall$
 $\text{hasParameter.ComputationalMethodParameter}$
 (Cal3) $\text{DensityFunctionalTheoryMethod} \sqsubseteq \text{ComputationalMethod}$
 (Cal4) $\text{HartreeFockMethod} \sqsubseteq \text{ComputationalMethod}$
 (Cal5) $\text{DensityFunctionalTheoryMethod} \sqsubseteq \exists \text{ hasXCFunctional.ExchangeCorrelationEnergyFunctional}$
 $\sqcap \forall \text{ hasXCFunctional.ExchangeCorrelationEnergyFunctional}$
 (Cal6) $\text{GeneralizedGradientApproximation} \sqsubseteq \text{ExchangeCorrelationEnergyFunctional}$
 (Cal7) $\text{LocalDensityApproximation} \sqsubseteq \text{ExchangeCorrelationEnergyFunctional}$
 (Cal8) $\text{metaGeneralizedGradientApproximation} \sqsubseteq \text{ExchangeCorrelationEnergyFunctional}$
 (Cal9) $\text{HybridFunctional} \sqsubseteq \text{ExchangeCorrelationEnergyFunctional}$
 (Cal10) $\text{HybridGeneralizedGradientApproximation} \sqsubseteq \text{HybridFunctional}$
 (Cal11) $\text{HybridmetaGeneralizedGradientApproximation} \sqsubseteq \text{HybridFunctional}$

Figure 12. Description logic axioms for the Calculation module.

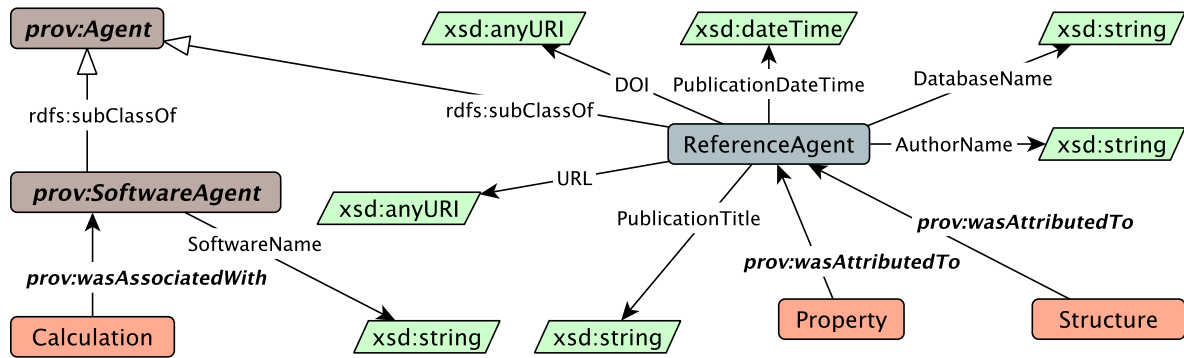


Figure 13. Concepts and relations in the Provenance module.

(Prov1) $\text{ReferenceAgent} \sqsubseteq \text{Agent}$
 (Prov2) $\text{Structure} \sqsubseteq \forall \text{ wasAttributedTo.ReferenceAgent}$
 (Prov3) $\text{Property} \sqsubseteq \forall \text{ wasAttributedTo.ReferenceAgent}$
 (Prov4) $\text{Calculation} \sqsubseteq \exists \text{ wasAssociatedwith.SoftwareAgent}$

Figure 14. Description logic axioms for the Provenance module.

only show the instances corresponding to the calculation's output structure, and for multiple calculated properties, species and sites, we only show one instance respectively. Connected to the instances of the Core module's concepts, are instances representing the structural information of the output structure, the provenance information of the output structure and calculated property, and the information about the computational method used for the calculation.

5. Extending MDO

In this section we use the approach in [29] to propose new concepts for MDO. The result of this work is a list of proposed concepts that are validated by a domain expert to be relevant to the domain. However, at this point the

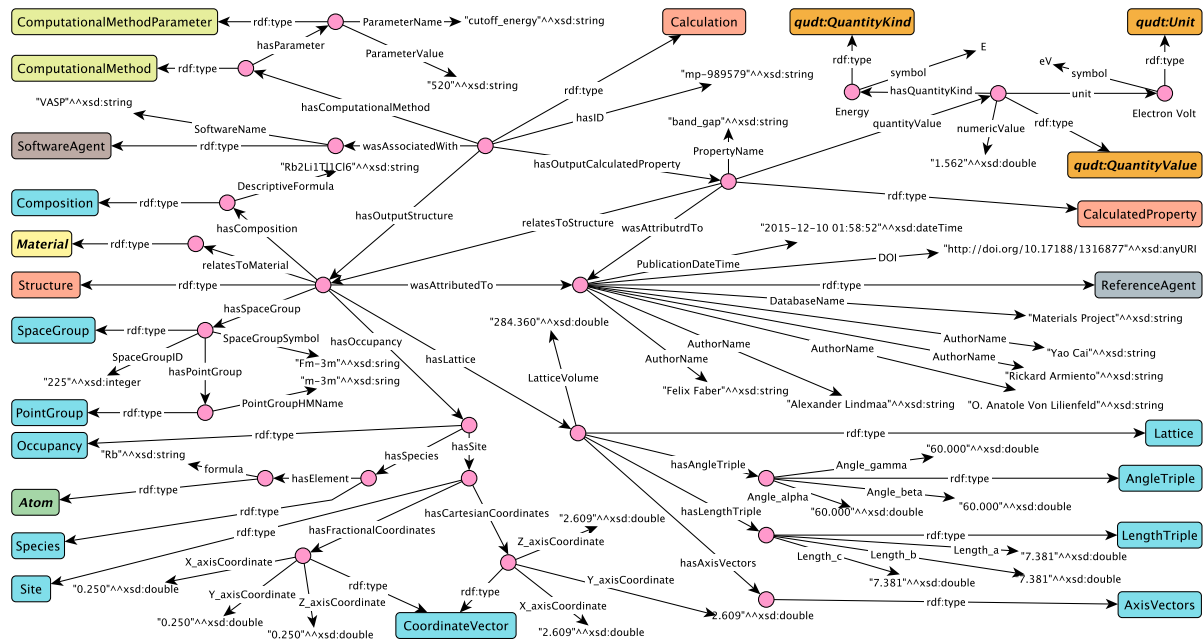


Figure 15. An instantiated materials calculation.

concepts are not yet included in the public version of MDO. Discussions are ongoing regarding the scope of the extension of MDO with respect to the domain and intended use of MDO.

A first step in the approach in [29] is to collect the corpus that is used as input. To be able to find as relevant information for MDO as possible, we used MDO as a seed for querying journal databases. The 37 concepts of MDO were used as search phrases for the titles and abstracts of two journals in the field of materials design, NPJ Computational Materials (<https://www.sciencedirect.com/journal/computational-materials-science>) and Computational Materials Science (<https://www.nature.com/npjcompumats/>), to find relevant articles. For these articles we retrieved the titles and abstracts. The final corpus contained titles and abstracts from 403 articles of NPJ Computational Materials and 8,193 from Computational Materials Science.

In the preprocessing step characters were set to lower case and punctuations were removed. Further, we removed words of length one or two. One consequence is that often materials symbols are removed. An advantage is that the phrases and words are usually not material dependent, but we miss cases where this is interesting.

After preprocessing there were 21,548 distinct words which together occur 808,862 times. An overview of the frequency of the words is presented in Table 2. Most of the words (72.27%) occur less than 10 times, while there are 17 words that occur more than 3000 times. These are 'based', 'properties', 'method', 'calculations', 'phase', 'materials', 'study', 'structure', 'temperature', 'density', 'results', 'energy', 'electronic', 'model', 'molecular', 'simulations', and 'surface'.

5.1. Frequent phrases

As explained in Section 2.3, the ToPMine system [40] identifies high-quality frequent phrases of a length up to a maximum length that is given as an input parameter. In our experiments this was set to 10. The second column of Table 3 shows the number of frequent phrases that ToPMine⁴ generates for different values of *min_support*. The higher the *min_support*, the fewer frequent phrases are generated.

⁴Note that ToPMine, as described in [40], does not use the preprocessing step.

Table 2

The distribution of word frequency after preprocessing.

Frequency	Percentage of words
less than 10	72.27
10-30	13.25
31-100	7.76
101-500	5.25
501-1000	0.83
1001-2000	0.44
2001-3000	0.12
More than 3000	0.08

Table 3

Number of frequent phrases for *min_support* 10, 15, 20, 25 and 30, respectively, and three different versions of the ToPMine algorithm. For the ToPMine_max versions, *max_support_word* is set to 8000.

<i>max_support_word</i>	<i>min_support</i>	ToPMine	ToPMine_max without stemming	ToPMine_max with stemming
8000	10	6,901	6,478	5,452
	15	3,826	3,578	3,022
	20	2,542	2,402	2,046
	25	1,816	1,722	1,477
	30	1,375	1,298	1,119

Table 4

Number of frequent phrases for *min_support* 10 and for *max_support_word* 500, 1000, 3000, 5000, and 8000, respectively, for two different versions of the ToPMine_max algorithm.

<i>min_support</i>	<i>max_support_word</i>	ToPMine_max without stemming	ToPMine_max with stemming
10	8,000	6,478	5,452
	5,000	5,947	5,023
	3,000	4,692	4,090
	1,000	1,878	1,692
	500	932	866

We also defined a maximum support threshold *max_support_word*. We call the system that uses this additional threshold and the preprocessing step, TopMine_max. Words that occur more than *max_support_word* times were removed. These words were usually very general terms that are not interesting for an ontology or that would not be interesting for a domain ontology, but possibly for a top-level ontology. We do note, however, that some of these words could be useful such as ‘method’, ‘electronic’, ‘model’, and ‘molecular’. The second column in Table 4 shows how *max_support_word* influences the number of generated frequent phrases with a constant *min_support* of 10. The higher *max_support_word*, the more frequent phrases are generated. Note that no word occurs more than 8000 times in our corpus, so setting *max_support_word* to 8000 allows all words (or, in other words, *max_support_word* is not used).⁵

Another way to look at the influence of *min_support* and *max_support_word* is to compare how many of the frequent phrases are the same and different for different settings. In Figure 16 we show this comparison of different settings to the base setting (ToPMine_max with stemming) where *min_support* is 10 and *max_support_word* is 8000 (i.e., *max_support_word* is not used) which is shown in the middle of the figure. The ‘Same’ bars show how many generated phrases occur both in the base setting and the compared setting. The ‘Removed’ bars show how many frequent phrases occur in the base setting, but not in the compared setting. For the cases where we

⁵Therefore, the difference in numbers for ToPMine and ToPMine_max without stemming in Table 3 shows the influence of the preprocessing step.

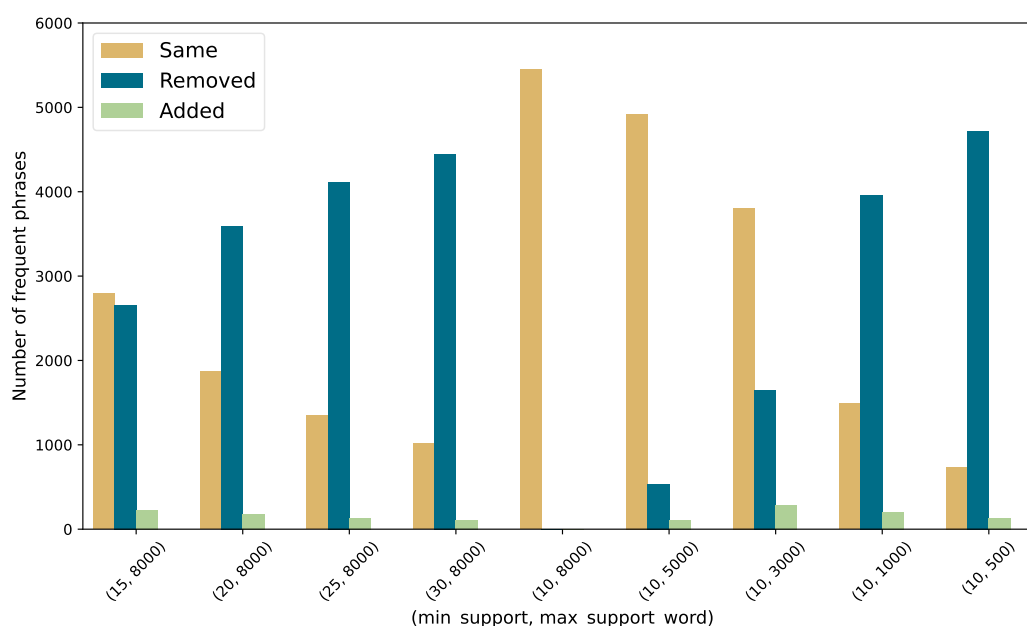


Figure 16. Comparison of the frequent phrases of ToPMine_max with stemming and *min_support* 10 (and *max_support_word* 8000) to settings with *min_support* 15, 20, 25 and 30, (and *max_support_word* 8000), respectively, and settings with *min_support* 10 and *max_support_word* 500, 1000, 3000, 5000, respectively.

change *min_support*, these would be phrases that are frequent phrases for *min_support* 10, but not for the higher *min_support* in the compared setting. For example ‘computational screening’ is removed for *min_support* set to 15. For the cases where we change the *max_support_word*, these would be phrases with words that occur more often than the *max_support_word* in the compared setting. For instance, ‘sheet metal forming’ contains the word ‘metal’ with frequency 3457 and would be removed for *max_support_word* set to 1000. The ‘Added’ bars show which frequent phrases occur newly in the compared settings. This happens, as stated before, because ToPMine does not generate all frequent phrases, but focuses on high-quality frequent phrases. As an example, ‘exchange correlation potential’ appears at least 10 times and less than 30 times and ‘exchange correlation’ appears at least 30 times. Both are frequent phrases for *min_support* set to 10. However, ToPMine does not generate ‘exchange correlation’ for *min_support* set to 10, but it does generate ‘exchange correlation potential’. For *min_support* set to 30 ‘exchange correlation potential’ is not a frequent phrase, while ‘exchange correlation’ is, and ToPMine does generate ‘exchange correlation’ as a frequent phrase.

Further, we investigated the influence of using stemming on the frequent phrases. For instance, the phrases ‘molecular dynamics simulations’, ‘molecular dynamics simulation’, ‘molecular dynamic simulations’ and ‘molecular dynamic simulation’ have the same stem ‘molecular dynam simul’. Stemming allows for removing redundant phrases and thus reduces the work of the domain expert. The influence on the number of generated phrases can be seen by comparing the last two columns in Tables 3 and 4. A disadvantage is that in some cases possible concept candidates may be removed. To alleviate this problem we show the domain expert for each of the stemmed frequent phrases the list of corresponding original phrases. This also helps the domain expert to choose terms to be added to the ontology.

In Table 5, we show the candidate concepts based on the validation of a domain expert on the frequent phrases from the experiment with *min_support* 30 and *max_support_word* 500. In total, 88 candidate concepts are suggested based on 81 out of 131 frequent phrases generated by the experiment. Some candidate concepts can be added into MDO as sub-concepts of existing concepts. For instance, ‘Linearized Augmented Plane Wave Method’ is a sub-concept of ‘Density Functional Theory Method’. Some candidate concepts are relevant to the materials design domain but may be not interesting for data access or data integration over materials design databases. For instance, ‘Covalent Bond’ is a bonding type that can be used to describe materials structures.

5.2. Topics

After the phrase mining we generated topics represented as sets of phrases. The number of topics (*num_topic*) is an input parameter to ToPMine. Each topic contains a set of phrases and these sets do not have to be disjoint. For instance, Figure 17 shows the overlap of phrases between topics for different settings of input parameters. In general, when we increase the number of topics, the number of frequent phrases in each topic decreases and the overlap between topics decreases as well.

The domain expert validates these topics and if possible, labels them to generate concepts for the ontology. In Table 6, we show the domain expert validation on 10 topics generated by ToPMine_max with stemming, *min_support* set to 30 and *max_support_word* set to 500. Among these topics, there are two topics (topics 0 and 9) that are interpreted with multiples labels, i.e., the domain expert divided the topic in different parts. The other topics received one label. Further, representative phrases are given for each topic. The labels and the representative phrases can all lead to new concepts.

Table 5

Candidate concepts based on domain expert validation on the experiment with *min_support* 30 and *max_support_word* 500.

Stacking Fault	Stone-wales Defect	Cement Paste
Van der Waals Force	Covalent Bond	Edurance Limit
Functionally Graded Material	Symmetric Tilt Grain Boundary Structure	Fatigue Limit
Linearized Augmented Plane Wave Method	Asymmetric Tilt Grain Boundary Structure	Perdew-Burke-Ernzerhof (PBE) Exchange-Correlation Functional
Face Centered Cubic	Rock Salt Structure	Porous Media
Boron Nitride	Rock Salt	Microstructural Features
Nearest Neighbor	Projector Augmented Wave Method	Hall-Petch Relation
Body Centered Cubic	Iron	Conduction Band
Coarse Grained Model	Cahn-Hilliard Equation	Slip Plane
Fiber Reinforced	Cauchy-Born Rule	Vapor Deposition
Zinc Blende	Domain Wall	Spinodal Decomposition
Core Shell	Armchair	Spontaneous Polarization
Rare Earth	Zigzag	Absorption Spectrum
Refractive Index	Double Walled Nanotube	Charpy Impact Test
Half metallicity	Power Factor	Alkaline Earth Metal
X-ray diffraction	Carbon Nanotube (cnt)	Contact Angle
Modified Embedded Atom Method	Mixed Mode Fracture	Vickers Hardness
Unit Cell	Homo-lumo Energy Gap	Rutile Titanium Dioxide (TiO ₂)
Absorption Spectra	Stainless Steels	Kinematic Hardening
Glass Formation	Vibrational Modes	Hexagonal Close Packed (hcp)
Brillouin Zone	Domain Switching	Anomalous Hall Effect
Lennard Jones	Sound Velocity	Valence Band
Dispersion Curves	Anatase (TiO ₂)	Voight Model
Cohesive Zone Model	Austenitic Stainless Steel	Reuss Model
Quasi-harmonic Debye Model	Crystallographic Orientation	Solute Segregation
Additive Manufacturing	Brittle Transition	Directional Solidification
Real Space Methods	Ductile Transition	Muffin-tin Orbital method
Quasi-harmonic Model	Brittle-Ductile Transition	Muffin-tin Orbital Approximation
Quantum Dot	Modified Becke-Johnson Exchange-Correlation Functional	
Hexagonal Boron Nitride	Kohn-Sham	

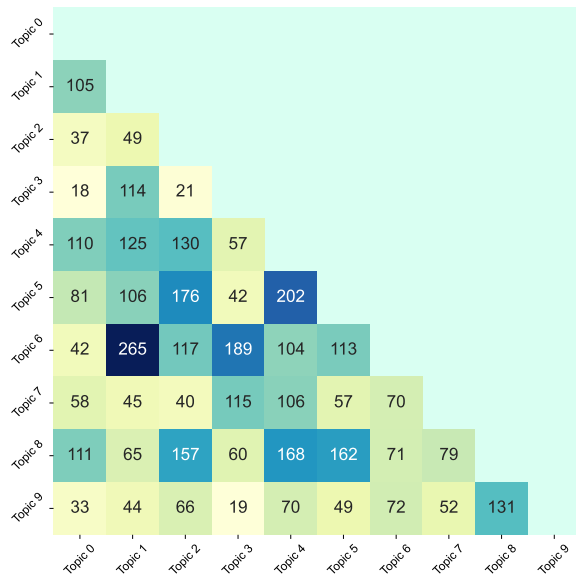
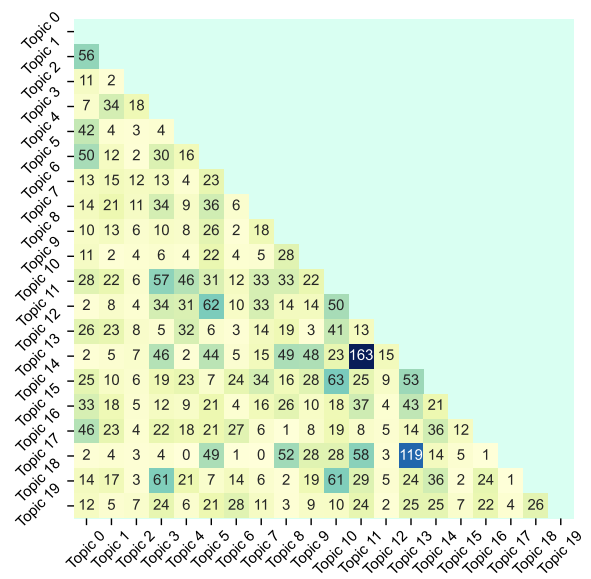
(a) $\min_support$ 10, num_topic 10(b) $\min_support$ 10, num_topic 20

Figure 17. Number of common phrases between pairs of topics.

Table 6

Topic labelling based on domain expert validation on the experiment with $\min_support$ 30 and $\max_support_word$ 500 (Up to five representative phrases are selected for each label).

No.	Topic labels	Representative Phrases
0	Computational Method Categories	Linearized Augmented Plane Wave Method, Hartree-Fock Method, Kohn-Sham, Perdew-Burke-Ernzerhof (PBE) Exchange-Correlation Functional, Modified Becke-Johnson Exchange-Correlation Functional,
	Materials Properties and Features	Absorption Spectrum, Refractive Index, Homo-lumo Energy Gap, Alkaline Earth Metal, Dispersion Curves
	Electronic Structure Features	Conduction Band, Valence Band
	Materials Categorizations	Half Metallicity, Rare Earth
	Experimental Method Categories	X-ray Diffraction
	Specific Materials	Zinc Blende
	Applications	Optoelectronic Devices
1	Hardness-related Materials Concepts	Quasi-harmonic Debye Model, Quasi-harmonic Model, Rock Salt, Sound Velocity, Zinc Blende
2	Materials Strength-related Concepts	Stacking Fault, Van der Waals Force, Tension Compression, Uniaxial Tension, Symmetric Tilt Grain Boundary Structure
3	Materials Fatigue/Fracture-related Concepts	Functionally Graded Material, Fiber Reinforced, Cohesive Zone Model, Unit Cell, Cement Paste
4	Materials Synthesis Concepts	Additive Manufacturing, Vapor Deposition, Directional Solidification, Microstructural Features, Crystallographic Orientations
5	Battery-related Materials Concepts	Ion Batteries, Anatase (TiO ₂), Lithium Ion Batteries, Rutile Titanium Dioxide (TiO ₂), Boron Nitride
6	Materials Structural Categorizations	Face Centered Cubic, Body Centered Cubic, Coarse Grained Model, Hexagonal Close Packed (hcp), Iron
7	Nanotube-related Concepts	Armchair, Boron Nitride, Hexagonal Boron Nitride, Carbon Nanotube (cnt), Cross Section
8	Artificial Intelligence-Methods	Artificial Neural, Neural Networks, Open Source, Degrees Freedom, Artificial Neural Networks
9	Materials Concepts for Solar-cells	Solar Cells, Quantum Dots, Domain Wall, Power Factor, Electric Fields
	Materials Magnetism Concepts	Domain Switching, Anomalous Hall Effect
	Materials Polarization Concepts	Spontaneous Polarization

6. Using MDO in ontology-based access to materials databases

In this section we show how MDO can be used for providing semantic and integrated access to materials databases. As a proof of concept we implemented data integration over two data sources, Materials Project [5] and OQMD [32] using a new GraphQL-based framework for data access and integration. This framework is introduced in [47, 48] and illustrated in Figure 18. The framework generates a GraphQL server that provides integrated access to data from heterogeneous data sources. These data sources may be based on different schemas and formats and may be accessed in different ways (e.g., tabular data accessed via SQL queries or JSON-formatted data accessed via a REST API). To address the heterogeneity, the framework relies on an ontology that provides an integrated view of the data from the different sources, and corresponding semantic mappings that define how the data from the underlying data sources is represented as instances of the ontology (arrows (a) and (b)). Furthermore, two processes are defined. The first process generates the GraphQL server. This includes generating both a GraphQL schema for the API provided by the server (arrow (i)) and a generic resolver function (arrow (ii)). This process does not need to be repeated unless the ontology or the mappings change. After this generation process, the GraphQL server can be set up. The second process deals with query answering and is performed after the GraphQL server is set up. During this process the query is validated against the GraphQL schema (arrow (1)); the underlying data sources are accessed via resolver functions, the retrieved data is combined, and the data is structured according to the schema (arrows (2) and (3)), and finally the query result is returned (arrow (4)). Details are available in [47, 48].

6.1. An MDO-based data access and integration system

In our proof of concept implementation we use MDO as the ontology to generate the GraphQL server. The GraphQL server contains a GraphQL schema generated based on MDO, and a generic resolver function that allows for accessing underlying data sources and restructuring the obtained data according to the GraphQL schema. This generic resolver function is implemented based on RML [49, 50] semantic mappings defined using MDO terminology.

In a GraphQL API, the GraphQL schema defines types, their fields, and the value types of the fields. An object type represents a list of fields and each field has a value of a specific type such as object type or scalar type. A scalar is used to represent a value such as a string. An input object type can be used to define an input object with a set of input fields; the input fields are either scalars, or other input objects. A GraphQL schema also supports defining types that represent operations such as query and mutation. The schema presumes *Query* type as the query root operation type. The part of the final GraphQL schema shown in Listing 1 contains two basic object type definitions which are *Calculation* and *Structure*. Both have field definitions which represent the relationships to scalar types or to other object types. For instance, the *Calculation* type has a field definition *ID* of which the value type is *String*, and a field definition *hasOutputStructure* of which the value type is *Structure*. The *Query* type has a field

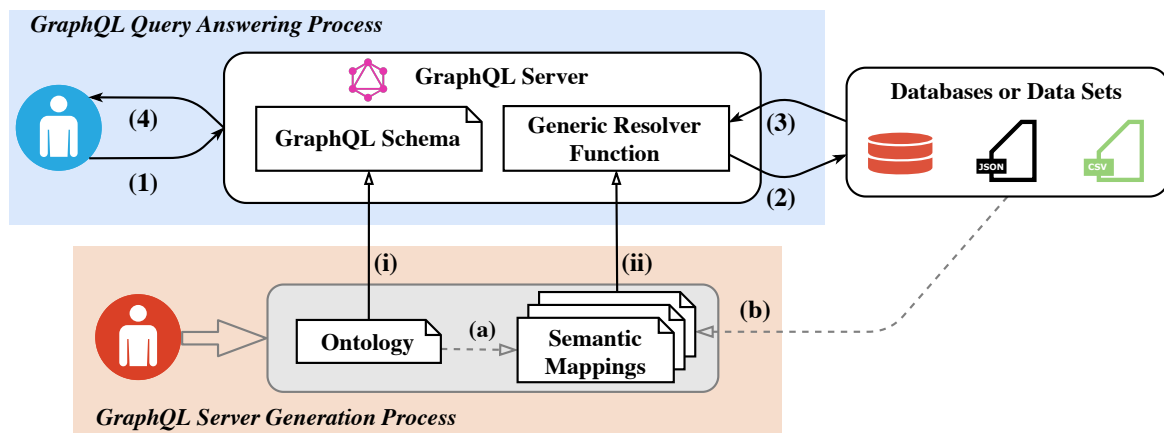


Figure 18. Framework of ontology-based GraphQL server generation (OBG-gen).

Listing 1: An excerpt of the GraphQL schema generated based on MDO.

```

1  type Query {
2    CalculationList(filter: CalculationFilter): [Calculation]
3  }
4  type Calculation {
5    hasOutputStructure: [Structure]
6    hasOutputCalculatedProperty: [CalculatedProperty]
7    ID: String
8  }
9  type Structure {
10   StructureID: String
11   hasComposition: Composition
12 }
13 input CalculationFilter {
14   hasOutputCalculatedProperty: CalculatedPropertyFilter
15   ID: StringFilter
16   _and: [CalculationFilter]
17   _or: [CalculationFilter]
18   _not: CalculationFilter
19 }
20 input CalculatedPropertyFilter {
21   PropertyName: StringFilter
22   numericalValue: FloatFilter
23   _and: [CalculatedPropertyFilter]
24   _or: [CalculatedPropertyFilter]
25   _not: CalculatedPropertyFilter
26 }
27 input StringFilter {
28   _eq: String
29   _gt: String
30   _in: [String]
31 }
32 input FloatFilter {
33   _eq: Float
34   _gt: Float
35   _in: [Float]
36 }

```

definition which is *CalculationList*. This definition allows users to write a GraphQL query that accesses all the entities of *Calculation* type, as shown in the example in Listing 12. Further, the schema contains four input object type definitions which are *CalculationFilter*, *CalculatedPropertyFilter*, *StringFilter* and *FloatFilter*, for capturing notions of filtering conditions that should be taken into account in the evaluation of the GraphQL query. As an example, the *CalculationFilter* can be used as an argument of the *CalculationList* query. For instance, in the query example of Listing 12, the argument of the query from line 3 to line 10 is used to represent a conjunctive filter expression with the meaning ‘property name is band gap’ and ‘the value of the property is greater than 5’. In our generic implementation of GraphQL resolver functions, a filter expression represented by an input object type will be parsed and evaluated against the underlying data.

Listing 2 shows an example of mappings in RML related to ‘band gap’ which is a *CalculatedProperty*. In general, an RML document has one or more *Triples Maps* which declare how input data is mapped into triples of the form (subject, predicate, object). A *Triples Map* contains the following three components; *Logical Source*, *Subject Map* and a set of *Predicate-Object Maps*. A logical source declares the source of input data to be mapped (e.g., line 2 to line 6). It contains definitions of *source* locating the input data source, *reference formulation* declaring how to refer to the input data, and *logical iterator* declaring the iteration loop used to map the input data. A subject map declares the rule for generating subjects when mapping input data into triples (e.g, line 7 to line 10). A predicate-object map

consists of one or more predicate maps declaring how to generate predicates of triples (e.g., line 12), one or more object maps or referencing object maps defining how to generate objects of triples (e.g., line 19 to line 25). An object map can be a reference-valued term map (e.g., line 46 to line 48) or a constant-valued term map (e.g., line 13 to line 15). In Listing 3, we show an excerpt of the JSON response based on Materials Project for the query to retrieve the data in which the task_id of the calculation is mp-989579. All material needed to generate the server is available online at <https://github.com/LiUSemWeb/OBG-gen> (including, e.g., files with the source code, mappings, the queries and documentation).

Listing 2: An excerpt of the RML mappings defined based on MDO.

```

1 <BandGapPropertyMapping>
2   rr:logicalSource [
3     rml:source "http://example.com/mp-989579_Rb2LiTlCl6.json";
4     rml:referenceFormulation ql:JSONPath;
5     rml:iterator "$.data[*]";
6   ];
7   rr:subjectMap [
8     rr:template "http://example.com/mdo/bandgapproperty/{task_id}";
9     rr:class core:CalculatedProperty;
10  ];
11  rr:predicateObjectMap [
12    rr:predicate core:hasPropertyName;
13    rr:objectMap [
14      rr:constant "band_gap";
15    ];
16  ];
17  rr:predicateObjectMap [
18    rr:predicate qudt:quantityValue;
19    rr:objectMap [
20      rr:parentTriplesMap <BandGapQuantityValueMapping>
21      rr:joinCondition [
22        rr:child "task_id";
23        rr:parent "task_id";
24      ];
25    ];
26  ];
27 <BandGapQuantityValueMapping>
28   rr:logicalSource [
29     rml:source "http://example.com/mp-989579_Rb2LiTlCl6.json";
30     rml:referenceFormulation ql:JSONPath;
31     rml:iterator "$.data[*]";
32   ];
33   rr:subjectMap [
34     rr:template "http://example.com/mdo/bandgapquantityvalue/{task_id}";
35     rr:class qudt:QuantityValue;
36   ];
37   rr:predicateObjectMap [
38     rr:predicate qudt:unit;
39     rr:objectMap [
40       rr:constant qudt_unit:EV;
41     ];
42   ];
43   rr:predicateObjectMap [
44     rr:predicate qudt:numericalValue;
45     rr:objectMap [
46       rr:reference "BandGap";
47     ];
48   ];
49 ]

```

Listing 3: An excerpt of the JSON response based on Materials Project API.

```

1  {
2    "data": [
3      {
4        "band_gap": 1.5623,
5        "density": 3.474406325286245,
6        "elements": ["Rb", "Li", "Tl", "Cl"],
7        "final_energy": -31.6397249,
8        "full_formula": "Rb2Li1Tl1Cl6",
9        "spacegroup": {
10         "symprec": 0.1,
11         "source": "spglib",
12         "symbol": "Fm-3m",
13         "number": 225,
14         "point_group": "m-3m",
15         "crystal_system": "cubic"
16       },
17       "task_id": "mp-989579",
18       "volume": 284.3605319552886
19     ]
20   }
21 }

```

6.2. Comparison

We compare our tool, OBG-gen (Ontology-Based GraphQL Server Generation) with three systems: morph-rdb [51], HyperGraphQL [52], and UltraGraphQL [53]. Morph-rdb is a tool that can access a relational database by translating SPARQL queries into SQL queries based on R2RML mappings. HyperGraphQL and its extension UltraGraphQL are GraphQL interfaces to query Linked Data that may be provided by local RDF files and remote SPARQL endpoints.

The semantic mappings (for all the systems) are based on the MDO. OBG-gen generates the GraphQL schema based on MDO. For UltraGraphQL and HyperGraphQL we use a modified version of the generated schema since they require directive definitions, as additional configurations for object type or field definitions, to specify the context information when translating a GraphQL query to SPARQL query (e.g., for an object type in the GraphQL schema, what is the URL of the object type's corresponding class in the RDF data.).

6.2.1. Data

The data from Materials Project and OQMD represents five different types of entities (Calculation, Structure, Composition, Band Gap and Formation Energy). We collected data in the sizes of 1K (i.e., 1000 entries), 2K, 4K, 8K, 16K and 32K from each database for populating the five entity types. We represented this data in different formats, i.e., tabular data for relational databases and for CSV files, and JSON-formatted data for JSON files. Additionally, for the RDF-based systems in our evaluation, we created an RDF file based on RML mappings and MDO for each dataset setting. We used six dataset settings for the experiments, which are 1K-1K, 2K-2K, 4K-4K, 8K-8K, 16K-16K and 32K-32K. Taking 32K-32K as an example, for each entity type, the test data contains the 32K data from Materials Project and the 32K data from OQMD.

6.2.2. Systems

Morph-rdb is served with data stored in a single database instance containing data from Materials Project and OQMD in separate tables. HyperGraphQL and UltraGraphQL are served with the same RDF data for each dataset setting. We use OBG-gen with two input settings. OBG-gen-rdb is served with two MySQL database instances hosting data from Materials Project and OQMD respectively. Conceptually, OBG-gen-mix is also served with two database instances. However, each instance contains different formats of data such as data in MySQL database, CSV or JSON files.

Table 7
Query Characteristics.

Query	CQ	DI	Filter	Query	CQ	DI	Filter
Q1	CQ5			Q7	CQ5		✓
Q2	CQ2, CQ5	✓		Q8	CQ5	✓	✓
Q3	CQ1	✓		Q9	CQ6, CQ7	✓	✓
Q4	CQ1, CQ2, CQ5	✓		Q10	CQ6, CQ7	✓	✓
Q5	CQ1, CQ2, CQ5			Q11	CQ1, CQ2, CQ5		✓
Q6			✓	Q12	CQ5	✓	✓

6.2.3. Queries

The queries that are used in our experiments are listed in Appendix A. We describe their characteristics in Table 7. The ‘CQ’ column describes which competency questions from Section 3.1 are covered by the queries. As the selected data covers competency questions CQ1-2 and CQ5-7, these are the ones that are covered. However, the other competency questions would in principle be easily covered with other or extended datasets. The ‘DI’ column shows which queries are of particular interest in the domain, i.e., these are often used queries to the materials databases. The other queries are mainly used to evaluate system performance on technically difficult queries. The ‘Filter’ column indicates whether the query contains filters.

As example, query Q9 in Listing 12 requests all the entities of Calculation type of which the value of the band gap property is larger than 5 electron volt. For such calculation entities, the query requests the corresponding values of ID, and reduced chemical formula of the composition of the output structure. Query Q12 in Listing 15 requests all the entities of type Structure which contain the silicon element.

6.2.4. Experiments and measurements

We evaluate the query execution time (QET) of the different systems over the six dataset settings. For each query separately, we run the query four times and always consider the first run as a warm-up, then take the average of the values of the remaining three runs. Figure 19 and Figure 20 illustrate the measurements for all data sizes and all queries. The measures for all data sizes and all queries are available online at <https://github.com/LiUSemWeb/OBG-gen/evaluation/README.md>. For UltraGraphQL, we have measurements only for queries Q1–Q4 because UltraGraphQL does not support queries with filtering conditions. For HyperGraphQL, regarding queries with filter expressions, we only have the measurement for Q7 because the system can only deal with simple filter expressions.

6.2.5. Results and discussion

We observe that both GraphQL servers generated by OBG-gen-rdb and OBG-gen-mix can answer all the 12 queries and thus the covered competency questions of MDO.

We also observe that increasing dataset sizes lead to increasing QETs (Figure 20). For queries without filtering conditions (Q1–Q5) (Figures 19 and 20), all of the systems have increases of QETs as the size of the dataset increases. However, morph-rdb is less sensitive to the data size increase compared with other systems. UltraGraphQL and HyperGraphQL outperform other systems for some smaller datasets (e.g., UltraGraphQL’s QETs of Q1 and Q2, HyperGraphQL’s QETs for Q1 from 1K-1K to 4K-4K). We explain this by the fact that these two systems have additional context information declaring URIs of classes to which instances in the RDF data belong. This is in contrast with the other systems which have to make use of semantic mappings to output queries to be evaluated against the underlying data sources. OBG-gen-rdb outperforms morph-rdb for some queries in smaller datasets (e.g., Q1 in 1K-1K, Q5 in 1K-1K and 2K-2K). For some queries, OBG-gen-rdb and morph-rdb have close QETs (e.g., Q2 in 1K-1K).

Another observation is regarding how OBG-gen-rdb and morph-rdb perform for queries with filter conditions (Q6–Q12) (Figures 19 and 20). The two systems behave similarly for Q6 with stable QETs and Q12 with slight increases, as the data size increases. The result size of Q6 is a constant over all the datasets in different sizes. Additionally, the filter expressions for Q6 and Q12 are simpler compared with those of Q7–Q11. Therefore, the QETs consumed for evaluating filtering expressions for Q6 and Q12 are less than those of Q7–Q11. For other queries (Q7–Q11), morph-rdb outperforms OBG-gen-rdb, however the differences between the two systems are less than those for queries without filtering conditions (e.g., Q1–Q4). The filtering conditions in GraphQL queries for

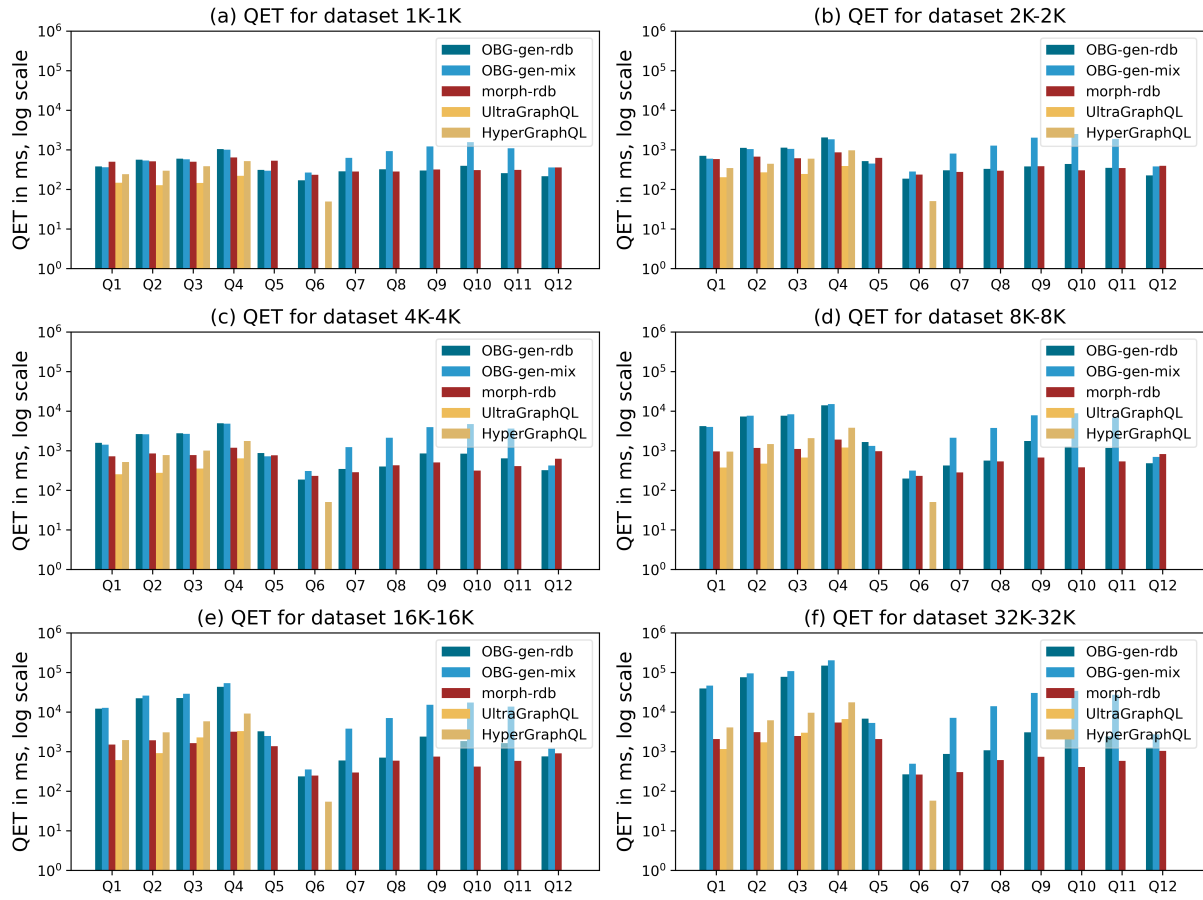


Figure 19. Query Execution Time (QET) per data size on materials dataset.

OBG-gen-rdb and in SPARQL queries for morph-rdb are written within *WHERE* clauses in SQL queries, thus will be evaluated against the back-end databases. A similar observation is also found in [54] where the experiments show that morph-rdb outperforms other systems (e.g., morph-graphql) as the size of dataset increases due to the SPARQL to SQL optimizations.

7. Conclusion

In this work we addressed the data access and interoperability issue for computational materials databases by developing MDO and providing a proof-of-concept implementation of an MDO-based data access and integration system for computational materials databases with a focus on solid-state physics and condensed matter theory. We have described MDO and a possible extension and showed that the proof-of-concept implementation can answer all competency questions for MDO, while not all of these could be answered by using the underlying databases' APIs.

One direction of future work is to extend the current proof-of-concept implementation in different ways. We want to integrate more databases as well as the OPTIMADE API. Further, as many end users in this domain may be more comfortable with form-based user interfaces, we will look into providing a form-based user interface or one that aids users to pose queries.

After discussion with domain experts we will extend the public version of MDO with the concepts and relations they deem appropriate. This includes discussing the concepts proposed in Section 5, but also looking into recent ongoing work in other projects such as EMMO-CIF (<https://github.com/emmo-repo/CIF-ontology>) and the VIMMP

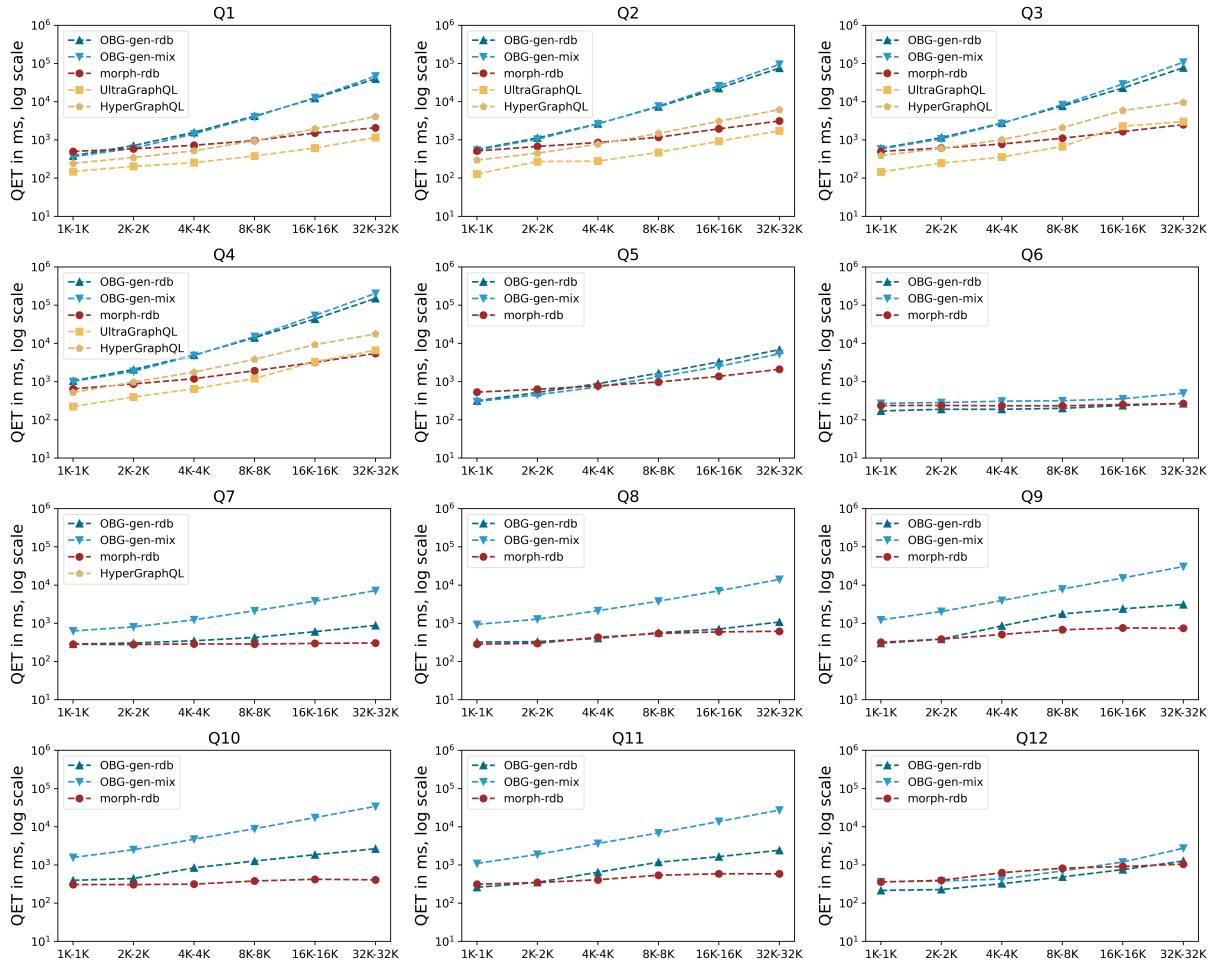


Figure 20. Query Execution Time (QET) per query on materials dataset.

ontologies [55, 56] (e.g., the VIMMP Ontology of Software) that may contain relevant concepts for extending MDO. In the latter case alignments may be provided to these ontologies. We will also use MDO in related domains as in a newly started project on interoperability of simulation systems.

We will also look into top-level ontologies and investigate which ontological commitments would be fitting MDO. This is, for instance, one of the topics of a recently accepted OntoCommons (<https://ontocommons.eu/>) demonstrator that we will lead.

Acknowledgements. This work has been financially supported by the Swedish e-Science Research Centre (SeRC), the Swedish National Graduate School in Computer Science (CUGS), the Swedish Research Council (Vetenskapsrådet, dnr 2018-04147), and the Swedish Agency for Economic and Regional and Growth (Tillväxtverket).

References

- [1] H. Li, R. Armiento and P. Lambrix, An Ontology for the Materials Design Domain, in: *The Semantic Web - ISWC 2020 - 19th International Semantic Web Conference, Athens, Greece, November 2-6, 2020, Proceedings, Part II*, J.Z. Pan, V.A.M. Tamma, C. d'Amato, K. Janowicz, B. Fu, A. Polleres, O. Seneviratne and L. Kagal, eds, Lecture Notes in Computer Science, Vol. 12507, Springer, 2020, pp. 212–227. doi:10.1007/978-3-030-62466-8_14.

- [2] M. Abd Nikooie Pour, H. Li, R. Armiento and P. Lambrix, A First Step towards Extending the Materials Design Ontology, in: *ESWC Workshop on Domain Ontologies for Research Data Management in Industry Commons of Materials and Manufacturing*, S. Chiacchiera, M.T. Horsch, J. Francisco Morgado and G. Goldbeck, eds, 2021, pp. 1–11.
- [3] European Commission and Directorate-General for Research and Innovation, *Materials research and innovation in the creative industries: report on the round table discussion, Brussels, 5 October 2012*, Publications Office, 2013. doi:doi/10.2777/30054.
- [4] G. Ceder and K. Persson, How Supercomputers Will Yield a Golden Age of Materials Science, *Scientific American* **309** (2013).
- [5] A. Jain, S.P. Ong, G. Hautier, W. Chen, W.D. Richards, S. Dacek, S. Cholia, D. Gunter, D. Skinner, G. Ceder and K.A. Persson, Commentary: The Materials Project: A materials genome approach to accelerating materials innovation, *APL Materials* **1**(1) (2013), 011002. doi:10.1063/1.4812323.
- [6] S. Curtarolo, G. Hart, M. Buongiorno-Nardelli, N. Mingo, S. Sanvito and O. Levy, The high-throughput highway to computational materials design, *Nature Materials* **12** (2013), 191–201. doi:10.1038/nmat3568.
- [7] K. Lejaeghere, G. Bihlmayer, T. Björkman, P. Blaha, S. Blügel, V. Blum, D. Caliste, I.E. Castelli, S.J. Clark, A.D. Corso, S.d. Gironcoli, T. Deutsch, J.K. Dewhurst, I.D. Marco, C. Draxl, M. Dulak, O. Eriksson, J.A. Flores-Livas, K.F. Garrity, L. Genovese, P. Giannozzi, M. Giantomassi, S. Goedecker, X. Gonze, O. Grånäs, E.K.U. Gross, A. Gulans, F. Gygi, D.R. Hamann, P.J. Hasnip, N.a.W. Holzwarth, D. Iuşan, D.B. Jochym, F. Jollet, D. Jones, G. Kresse, K. Koepnick, E. Küçükbenli, Y.O. Kvashnin, I.L.M. Locht, S. Lubeck, M. Marsman, N. Marzari, U. Nitzsche, L. Nordström, T. Ozaki, L. Paulatto, C.J. Pickard, W. Poelmans, M.I.J. Probert, K. Refson, M. Richter, G.-M. Rignanese, S. Saha, M. Scheffler, M. Schlipf, K. Schwarz, S. Sharma, F. Tavazza, P. Thunström, A. Tkatchenko, M. Torrent, D. Vanderbilt, M.J. van Setten, V.V. Speybroeck, J.M. Wills, J.R. Yates, G.-X. Zhang and S. Cottenier, Reproducibility in density functional theory calculations of solids, *Science* **351**(6280) (2016), aad3000. doi:10.1126/science.aad3000.
- [8] M.W. Gaultois, A.O. Oliynyk, A. Mar, T.D. Sparks, G.J. Mulholland and B. Meredig, Perspective: Web-based machine learning models for real-time screening of thermoelectric materials properties, *APL Materials* **4**(5) (2016), 053213. doi:10.1063/1.4952607.
- [9] R. Armiento, Database-Driven High-Throughput Calculations and Machine Learning Models for Materials Design, in: *Machine Learning Meets Quantum Physics*, K.T. Schütt, S. Chmiela, O.A. von Lilienfeld, A. Tkatchenko, K. Tsuda and K.-R. Müller, eds, Springer International Publishing, 2020, pp. 377–395. doi:10.1007/978-3-030-40245-7_17.
- [10] P. Lambrix, R. Armiento, A. Delin and H. Li, FAIR Big Data in the Materials Design Domain, in: *Encyclopedia of Big Data Technologies*, A.Y. Zomaya, J. Taheri and S. Sakr, eds, Springer, Cham, 2022. doi:10.1007/978-3-319-63962-8_293-2.
- [11] A. Agrawal and A. Choudhary, Perspective: materials informatics and big data: realization of the Fourth paradigm of science in materials science, *APL Materials* **4** (2016), 053208:1–10. doi:10.1063/1.4946894.
- [12] G.J. Mulholland and S.P. Paradiso, Perspective: Materials informatics across the product lifecycle: Selection, manufacturing, and certification, *APL Materials* **4**(5) (2016), 053207. doi:10.1063/1.4945422.
- [13] C.W. Andersen, R. Armiento, E. Blokhin, G.J. Conduit, S. Dwaraknath, M.L. Evans, Á. Fekete, A. Gopakumar, S. Gražulis, A. Merkys, F. Mohamed, C. Oses, G. Pizzi, G.-M. Rignanese, M. Scheidgen, L. Talirz, C. Toher, D. Winston, R. Aversa, K. Choudhary, P. Colinet, S. Curtarolo, D. Di Stefano, C. Draxl, S. Er, M. Esters, M. Fornari, M. Giantomassi, M. Govoni, G. Hautier, V. Hegde, M.K. Horton, P. Huck, G. Huhs, J. Hummelshøj, A. Karirya, B. Kozinsky, S. Kumbhar, M. Liu, N. Marzari, A.J. Morris, A.A. Mostofi, K.A. Persson, G. Petretto, T. Purcell, F. Ricci, F. Rose, M. Scheffler, D. Speckhard, M. Uhrin, A. Vaitkus, P. Villars, D. Waroquiers, C. Wolverton, M. Wu and X. Yang, OPTIMADE, an API for exchanging materials data, *Scientific Data* **8** (2021), 217. doi:10.1038/s41597-021-00974-z.
- [14] K. Cheung, J. Drennan and J. Hunter, Towards an Ontology for Data-driven Discovery of New Materials, in: *Semantic Scientific Knowledge Integration AAAI/SSS Workshop*, 2008, pp. 9–14.
- [15] X. Zhang, C. Hu and H. Li, Semantic query on materials data based on mapping MATML to an OWL ontology, *Data Science Journal* **8** (2009), 1–17. doi:10.2481/dsj.8.1.
- [16] T. Ashino, Materials Ontology: An Infrastructure for Exchanging Materials Information and Knowledge, *Data Science Journal* **9** (2010), 54–61, doi: 10.2481/dsj.008-041.
- [17] European Committee for Standardization, A Guide to the Development and Use of Standards Compliant Data Formats for Engineering Materials Test Data, 2010.
- [18] D.G. Thomas, R.V. Pappu and N.A. Baker, NanoParticle Ontology for cancer nanotechnology research, *Journal of Biomedical Informatics* **44**(1) (2011), 59–74, doi: 10.1016/j.jbi.2010.03.001.
- [19] J. Hastings, N. Jeliakova, G. Owen, G. Tsiliki, C.R. Munteanu, C. Steinbeck and E. Willighagen, eNanoMapper: harnessing ontologies to enable data integration for nanomaterial risk assessment, *Journal of Biomedical Semantics* **6** (2015), 10:1–15, doi: 10.1186/s13326-015-0005-5.
- [20] X. Zhang, D. Pan, C. Zhao and K. Li, MMOY: Towards deriving a metallic materials ontology from Yago, *Advanced Engineering Informatics* **30** (2016), 687–702. doi:10.1016/j.aei.2016.09.002.
- [21] F. Le Piane, M. Baldoni, M. Gaspari and F. Mercuri, Introducing MAMBO: Materials And Molecules Basic Ontology, in: *Proceedings of the Workshop on Domain Ontologies for Research Data Management in Industry Commons of Materials and Manufacturing (DORIC-MM 2021) co-located with the 18th European Semantic Web Conference (ESWC 2021)*, 2021, pp. 28–39. <http://purl.org/net/epubs/work/50300311>.
- [22] A. Zainul Ihsan, D. Dessi, M. Alam, H. Sack and S. Sandfeld, Steps towards a Dislocation Ontology for Crystalline Materials, in: *Proceedings of the Second International Workshop on Semantic Digital Twins co-located with the 18th Extended Semantic Web Conference (ESWC 2021)*, CEUR Workshop Proceedings, Vol. 2887, CEUR-WS.org, 2021. <http://ceur-ws.org/Vol-2887/paper4.pdf>.

- [23] M. Alam, H. Birkholz, D. Dessì, C. Eberl, H. Fliegl, P. Gumbsch, P. von Hartrott, L. Mädler, M. Niebel, H. Sack and A. Thomas, Ontology Modelling for Materials Science Experiments, in: *Joint Proceedings of the Semantics co-located events: Poster & Demo track and Workshop on Ontology-Driven Conceptual Modelling of Digital Twins co-located with Semantics 2021*, CEUR Workshop Proceedings, Vol. 2941, CEUR-WS.org, 2021. <http://ceur-ws.org/Vol-2941/paper11.pdf>.
- [24] S. Borgo, A. Galton and O. Kutz, Foundational ontologies in action, *Applied Ontology* **17**(1) (2022), 1–16. doi:10.3233/AO-220265.
- [25] B. Smith, On classifying material entities in Basic Formal Ontology, in: *Proceedings of the Third Interdisciplinary Ontology Meeting*, 2012, pp. 1–13. <https://philpapers.org/rec/smiocm>.
- [26] A. Gangemi, N. Guarino, C. Masolo, A. Oltramari and L. Schneider, Sweetening Ontologies with DOLCE, in: *Knowledge Engineering and Knowledge Management. Ontologies and the Semantic Web, 13th International Conference, EKAW 2002, Sigüenza, Spain, October 1-4, 2002, Proceedings*, A. Gómez-Pérez and V.R. Benjamins, eds, Lecture Notes in Computer Science, Vol. 2473, Springer, 2002, pp. 166–181. doi:10.1007/3-540-45810-7_18.
- [27] H. Herre, General Formal Ontology (GFO): A Foundational Ontology for Conceptual Modelling, in: *Theory and Applications of Ontology: Computer Applications*, Springer, 2010, pp. 297–345. doi:10.1007/978-90-481-8847-5_14.
- [28] J. Kaufman and E. Begley, MatML: A Data Interchange Markup Language, *Advanced Materials & Processes* **161** (2003).
- [29] H. Li, R. Armiento and P. Lambrix, A Method for Extending Ontologies with Application to the Materials Science Domain, *Data Science Journal* **18**(1) (2019). doi:10.5334/dsj-2019-050.
- [30] T. Pellissier Tanon, G. Weikum and F. Suchanek, YAGO 4: A Reason-able Knowledge Base, in: *The Semantic Web*, A. Harth, S. Kiriene, A.-C. Ngonga Ngomo, H. Paulheim, A. Rula, A.L. Gentile, P. Haase and M. Cochez, eds, Springer International Publishing, Cham, 2020, pp. 583–596.
- [31] P. de Matos, A. Dekker, M. Ennis, J. Hastings, K. Haug, S. Turner and C. Steinbeck, ChEBI: a chemistry ontology and database, *Journal of cheminformatics* **2**(S1) (2010), P6:1–, doi: 10.1186/1758-2946-2-S1-P6.
- [32] J.E. Saal, S. Kirklin, M. Aykol, B. Meredig and C. Wolverton, Materials Design and Discovery with High-Throughput Density Functional Theory: The Open Quantum Materials Database (OQMD), *JOM* **65** (2013), 1501–1509. doi:10.1007/s11837-013-0755-4.
- [33] C. Draxl and M. Scheffler, NOMAD: The FAIR concept for big data-driven materials science, *MRS Bulletin* **43**(9) (2018), 676–682, doi: 10.1557/mrs.2018.208.
- [34] S. Curtarolo, W. Setyawan, S. Wang, J. Xue, K. Yang, R. Taylor, L. Nelson, G. Hart, S. Sanvito, M. Buongiorno-Nardelli, N. Mingo and O. Levy, AFLOWLIB.ORG: A distributed materials properties repository from high-throughput ab initio calculations, *Computational Materials Science* **58**(Supplement C) (2012), 227–235. doi:10.1016/j.commatsci.2012.02.002.
- [35] V. Presutti, E. Blomqvist, E. Daga and A. Gangemi, Pattern-Based Ontology Design, in: *Ontology Engineering in a Networked World*, M.C. Suárez-Figueroa, A. Gómez-Pérez, E. Motta and A. Gangemi, eds, Springer, 2012, pp. 35–64. doi:10.1007/978-3-642-24794-1_3.
- [36] Z. Dragisic, P. Lambrix and E. Blomqvist, Integrating Ontology Debugging and Matching into the eXtreme Design Methodology, in: *Proceedings of the 6th Workshop on Ontology and Semantic Web Patterns*, E. Blomqvist, P. Hitzler, A. Krisnadhi, T. Narock and M. Solanki, eds, CEUR Workshop Proceedings, Vol. 1461, 2015. http://ceur-ws.org/Vol-1461/WOP2015_paper_1.pdf.
- [37] M.C. Suárez-Figueroa, A. Gómez-Pérez and M. Fernández-López, The NeOn methodology for ontology engineering, in: *Ontology engineering in a networked world*, M.C. Suárez-Figueroa, A. Gómez-Pérez, E. Motta and A. Gangemi, eds, Springer, 2012, pp. 9–34. doi:10.1007/978-3-642-24794-1_2.
- [38] M. Poveda-Villalón, A. Gómez-Pérez and M.C. Suárez-Figueroa, OOPS! (Ontology Pitfall Scanner!): An On-line Tool for Ontology Evaluation, *International Journal on Semantic Web and Information Systems (IJSWIS)* **10**(2) (2014), 7–34. doi:10.4018/ijswis.2014040102.
- [39] P. Lambrix and V. Ivanova, A unified approach for debugging is-a structure and mappings in networked taxonomies, *Journal of Biomedical Semantics* **4** (2013), 10. doi:10.1186/2041-1480-4-10.
- [40] A. El-Kishky, Y. Song, C. Wang, C.R. Voss and J. Han, Scalable Topical Phrase Mining from Text Corpora, *Proceedings of the VLDB Endowment* **8**(3) (2014), 305–316. doi:10.14778/2735508.2735519.
- [41] D.M. Blei, A.Y. Ng and M.I. Jordan, Latent Dirichlet Allocation, *Journal of Machine Learning Research* **3** (2003), 993–1022.
- [42] M. Abd Nikooie Pour, H. Li, R. Armiento and P. Lambrix, A First Step towards a Tool for Extending Ontologies, in: *Proceedings of the Sixth International Workshop on the Visualization and Interaction for Ontologies and Linked Data co-located with the 20th International Semantic Web Conference (ISWC 2021), Virtual Conference, 2021*, P. Lambrix, C. Pesquita and V. Wiens, eds, CEUR Workshop Proceedings, Vol. 3023, CEUR-WS.org, 2021, pp. 1–12. <http://ceur-ws.org/Vol-3023/paper2.pdf>.
- [43] C. Vardeman II, A. Krisnadhi, M. Cheatham, K. Janowicz, H. Ferguson, P. Hitzler and A. Buccellato, An Ontology Design Pattern and its use case for modeling material transformation, *Semantic Web* **8**(5) (2017), 719–731. doi:10.3233/SW-160231.
- [44] T. Lebo, S. Sahoo, D. McGuinness, K. Belhajjame, J. Cheney, D. Corsar, D. Garijo, S. Soiland-Reyes, S. Zednik and J. Zhao, PROV-O: The PROV Ontology, *W3C recommendation* (2013). <https://www.w3.org/TR/prov-o/>.
- [45] R. Haas, P.J. Keller, J. Hodges and J. Spivak, Quantities, units, dimensions and data types ontologies (QUDT), Accessed: 2022-07-29.
- [46] F.A. Faber, A. Lindmaa, O.A. Von Lilienfeld and R. Armiento, Machine Learning Energies of 2 Million Elpasolite (ABC_2D_6) Crystals, *Physical review letters* **117**(13) (2016), 135502. doi:10.1103/PhysRevLett.117.135502.
- [47] H. Li, O. Hartig, R. Armiento and P. Lambrix, Ontology-Based GraphQL Server Generation for Data Access and Integration, submitted.
- [48] H. Li, Ontology-Driven Data Access and Data Integration with an Application in the Materials Design Domain, PhD thesis, Linköping University, Sweden, 2022. doi:10.3384/9789179292683.
- [49] A. Dimou, M. Vander Sande, P. Colpaert, R. Verborgh, E. Mannens and R. Van de Walle, RML: A Generic Language for Integrated RDF Mappings of Heterogeneous Data, in: *Proceedings of the Workshop on Linked Data on the Web co-located with the 23rd International World Wide Web Conference (WWW 2014)*, Vol. 1184, 2014. http://ceur-ws.org/Vol-1184/ldow2014_paper_01.pdf.

- [50] A. Dimou, M. Vander Sande, J. Slepicka, P. Szekely, E. Mannens, C. Knoblock and R. Van de Walle, Mapping Hierarchical Sources into RDF Using the RML Mapping Language, in: *2014 IEEE International Conference on Semantic Computing*, 2014, pp. 151–158. doi:10.1109/ICSC.2014.25.
- [51] Morph-rdb, version 3.12.5, Accessed: 2022-07-29. <https://github.com/oeg-upm/morph-rdb/releases/tag/v3.12.5>.
- [52] Semantic Integration Ltd., HyperGraphQL, version 2.0.0, Accessed: 2022-07-29. <https://github.com/hypergraphql/hypergraphql/releases/tag/2.0.0>.
- [53] Semantic Integration Ltd., UltraGraphQL, version 1.0.0, Accessed: 2022-07-29. <https://git.rwth-aachen.de/i5/ultragraphql>.
- [54] D. Chaves-Fraga, F. Priyatna, A. Alobaid and O. Corcho, Exploiting Declarative Mapping Rules for Generating GraphQL Servers with Morph-GraphQL, *International Journal of Software Engineering and Knowledge Engineering* **30**(06) (2020), 785–803. doi:10.1142/S0218194020400070.
- [55] M.T. Horsch, S. Chiacchiera, W.L. Cavalcanti and B. Schembera, *Data Technology in Materials Modelling*, Zenodo, 2021. doi:10.1007/978-3-030-68597-3.
- [56] M.T. Horsch, S. Chiacchiera, M.A. Seaton, I.T. Todorov, D. Toti and G. Goldbeck, Introduction to the VIMMP ontologies, Zenodo, 2021, The VIMMP project is funded from the European Union’s Horizon 2020 research and innovation programme under grant agreement no. 760907. doi:10.5281/zenodo.4411422.

Appendix A. GraphQL queries

A.1. Queries without filter expressions

Query 1: List all the structures with the reduced chemical formula of each structure's composition.

Listing 4: Q1.

```
1 {
2   StructureList{
3     hasComposition{
4       ReducedFormula
5     }
6   }
7 }
```

Query 2: List all the calculations with the reduced chemical formula of each output structure's composition.

Listing 5: Q2.

```
1 {
2   CalculationList{
3     hasOutputStructure{
4       hasComposition{
5         ReducedFormula
6       }
7     }
8   }
9 }
```

Query 3: List all the calculations with the name and value of each output calculated property.

Listing 6: Q3.

```
1 {
2   CalculationList{
3     hasOutputCalculatedProperty{
4       PropertyName
5       numericalValue
6     }
7   }
8 }
```

Query 4: List all the calculations with the reduced chemical formula of each output structure's composition, and with the name and value of each output calculated property.

Listing 7: Q4.

```
1 {
2   CalculationList{
3     hasOutputStructure{
4       hasComposition{
5         ReducedFormula
6       }
7     }
8   }
9 }
```

```

6      }
7      }
8      hasOutputCalculatedProperty{
9          PropertyName
10         numericalValue
11     }
12 }
13 }

```

Query 5: List all the calculations and structures with their identifier information.

Listing 8: Q5.

```

1  {
2    ThingList{
3      ... on Calculation{iri}
4      ... on Structure{iri}
5    }
6  }

```

A.2. Queries with filter expressions

Query 6: List all the calculations with the ID of each calculation, and with the name and value of each output calculated property, where the ID of each calculation is in a given list of values.

Listing 9: Q6.

```

1  {
2    CalculationList (
3      filter: {
4        ID: {
5          _in: [ "6332", "8088", "21331", "mp-561628", "mp-614918" ]
6        }
7      }
8    )
9    {
10     ID
11     hasOutputCalculatedProperty {
12       PropertyName
13       numericalValue
14     }
15   }
16 }

```

*Query 7: List all the calculations with the ID of each calculation, and with the name and value of each output calculated property, where the ID of each calculation is in a given list of values **and** the reduced chemical formula of each output structure's composition is in a given list of values.*

Listing 10: Q7.

```

1  {
2    CalculationList (
3      filter: {

```

```

4      _and: [
5        {
6          ID: {
7            _in: ["6332", "8088", "21331", "mp-561628", "mp-614918"]
8          }
9        }
10       {
11         hasOutputStructure: {
12           hasComposition: {
13             ReducedFormula: {
14               _in: ["MnCl2", "YClO"]
15             }
16           }
17         }
18       }
19     ]
20   }
21 )
22 {
23   ID
24   hasOutputCalculatedProperty {
25     PropertyName
26     numericalValue
27   }
28 }
29 }

```

*Query 8: List all the calculations with the ID of each calculation, and with the name and value of each output calculated property, where the ID of each calculation is in a given list of values **and** the reduced chemical formula of each output structure's composition is in a given list A **or** a given list B.*

Listing 11: Q8.

```

1  {
2    CalculationList(
3      filter: {
4        _and: [
5          {
6            ID: {
7              _in: ["6332", "8088", "21331", "mp-561628", "mp-614918"]
8            }
9          }
10         {
11           _or: [
12             {
13               hasOutputStructure: {
14                 hasComposition: {
15                   ReducedFormula: { _in: ["MnCl2", "YClO"] }
16                 }
17             }
18           ]
19         }
20         {
21           hasOutputStructure: {
22             hasComposition: {
23               ReducedFormula: { _in: ["CeCrS2O", "SiO2", "O"] }
24             }
25         }
26       ]

```

```

27     }
28   ]
29 }
30 )
31 {
32   ID
33   hasOutputCalculatedProperty {
34     propertyName
35     numericalValue
36   }
37 }
38 }

```

Query 9: List all the calculations with the ID of each calculation, and with the reduced chemical formula of each output structure's composition, where the value of output band gap property is greater than 5.

Listing 12: Q9.

```

1 {
2   CalculationList(
3     filter: {
4       hasOutputCalculatedProperty: {
5         _and: [
6           { propertyName: { _eq: "Band Gap" } }
7           { numericalValue: { _gt: 5 } }
8         ]
9       }
10    }
11  )
12  {
13    ID
14    hasOutputStructure {
15      hasComposition {
16        ReducedFormula
17      }
18    }
19  }
20 }

```

*Query 10: List all the calculations with the ID of each calculation, and with the reduced chemical formula of each output structure's composition, and with the name and value of each output calculated property, where the value of output band gap property is greater than 5 **and** the reduced chemical formula of each output structure's composition is in a given list of values.*

Listing 13: Q10.

```

1 {
2   CalculationList(
3     filter: {
4       _and: [
5         {
6           hasOutputStructure: {
7             hasComposition: {
8               ReducedFormula: { _in: ["MnCl2", "YClO"] }
9             }
10          }

```



```

11      }
12      {
13          hasOutputCalculatedProperty: {
14              _and: [
15                  { PropertyName: { _eq: "Band Gap" } }
16                  { numericalValue: { _gt: 5 } }
17              ]
18          }
19      }
20  ]
21  }
22  )
23  {
24      ID
25      hasOutputStructure {
26          hasComposition {
27              ReducedFormula
28          }
29      }
30      hasOutputCalculatedProperty {
31          PropertyName
32          numericalValue
33      }
34  }
35  }

```

*Query 11: List all the calculations with the ID of each calculation, where, as first condition the property name is "Band Gap" **and** the value is greater than 4, **and** as second condition the property name is "Band Gap" **and** the value is greater than 4, **or** the reduced chemical formula is in a given list. We note that the filter condition can be simplified to only contain the first condition, i.e., the property name is "Band Gap" **and** the value is greater than 4.*

Listing 14: Q11.

```

1  {
2      CalculationList(
3          filter: {
4              _and: [
5                  {
6                      hasOutputCalculatedProperty: {
7                          _and: [
8                              { PropertyName: { _eq: "Band Gap" } }
9                              { numericalValue: { _gt: 4 } }
10                         ]
11                     }
12                 }
13             ]
14             _or: [
15                 {
16                     hasOutputCalculatedProperty: {
17                         _and: [
18                             { PropertyName: { _eq: "Band Gap" } }
19                             { numericalValue: { _gt: 4 } }
20                         ]
21                     }
22                 }
23             ]
24             hasOutputStructure: {
25                 hasComposition: {

```

```
26         ReducedFormula: { _in: ["YClO", "CsCl"] }
27     }
28 }
29 }
30 ]
31 }
32 ]
33 }
34 )
35 {
36     ID
37 }
38 }
```

Query 12: List all the structures that contain the Silicon element, with the reduced chemical formula of each structure’s composition.

Listing 15: Q12.

```
1  {
2      StructureList(
3          filter: {
4              hasComposition: {
5                  ReducedFormula: { _like: "%Si%" }
6              }
7          }
8      )
9      {
10         hasComposition {
11             ReducedFormula
12         }
13     }
14 }
```