

# CateCom: a practical data-centric approach to categorization of computational models.

Alexander Zech<sup>1,2</sup> and Timur Bazhirov<sup>1</sup>

<sup>1</sup>Exabyte Inc., San Francisco, CA, United States

<sup>2</sup>University of California, Berkeley, Berkeley, CA, United States

September 29, 2021

timur@exabyte.io

## Abstract

The advent of data-driven science in the 21st century brought about the need for well-organized structured data and associated infrastructure able to facilitate the applications of Artificial Intelligence and Machine Learning. We present an effort aimed at organizing the diverse landscape of physics-based and data-driven computational models in order to facilitate the storage of associated information as structured data. We apply object-oriented design concepts and outline the foundations of an open-source collaborative framework that is: (1) capable of uniquely describing the approaches in structured data, (2) flexible enough to cover the majority of widely used models, and (3) utilizes collective intelligence through community contributions. We present example database schemas and corresponding data structures and explain how these are deployed in software at the time of this writing.

**Keywords:** data structures, data standards, artificial intelligence, machine learning.

## 1 Introduction

The proliferation of data-driven science creates the need for systematically organized and machine-readable data formats.<sup>[1]</sup> While in the past considerable effort has been dedicated to structuring simulation results, the organization of simulation metadata has only recently gained attention. One major aspect of such an undertaking is the classification of numerous computational models. Early forms<sup>[2–4]</sup> of the categorization of quantum chemical models were based on only a few distinguishing descriptors, i.e. the treatment of electron correlation and one-particle basis set, as well as the type of Hamiltonian. More recently, projects emerged which not only collect results and metadata from output files of simulation packages but also define database schemas for their storage. For instance, the Novel Materials Discovery (NOMAD)<sup>[5]</sup> repository includes a structured collection of computational model metadata as part of its `metainfo` component. Another example is QCSchema<sup>[6]</sup> by the Molecular Science Software Institute (MolSSI), which provides software-independent data structures for quantum chemistry geared towards unified and consistent workflows. Organizing the computational models and their results can also be achieved through an ontology, often expressed in the Web Ontology Language (OWL)<sup>[7]</sup>. One such example is the OntoCompChem ontology<sup>[8]</sup>, which is applied to quantum chemistry calculations as part of the MolHub<sup>[9]</sup> web service. In the domain of materials science, ontologies are more prevalent but often focus on specific subdomains such as nanoparticles<sup>[10]</sup>. There are, however, examples of general ontologies, such as the Elementary Multiperspective Material Ontology (EMMO)<sup>[11]</sup> or the Materials Design Ontology (MDO)<sup>[12]</sup>. There exist a number of NIST/MGI-led prior efforts where large-scale high-throughput computational approaches have been used to screen thousands of compounds with subsequent web-based dissemination, databasing, and data-mining.<sup>[5,13–19]</sup>

Most current data structures for computational models include little information beyond just its name relying on the description in the scientific literature. Such an approach makes it difficult to construct data-driven predictions. We elaborate on the existing approaches by constructing a framework able to utilize previously obtained data (also allowing the generation of new data) to categorize the important descriptive features for a set of entities (materials, simulation workflows/models, computational methods) and target properties of interest (electronic, chemical, thermodynamic, structural properties) to construct associative maps, and organize "actionable" data in this extremely diverse and complex domain.

Our effort follows an object-oriented approach by building a basis of unit models, which are small inseparable sets of equations pertaining to a specific physical description of reality (e.g. Kohn-Sham Density Functional Theory). A given level of theory may further be expressed as a combination of unit models. Such modularity allows us to cover a diverse range of use cases. In addition to the categorization of the computational models, we discuss a semantic layer in the form of an ontology, which not only facilitates a more accurate description of the relationships between models, but also provides the foundation for further applications (e.g. building a knowledge graph, improved search, or AI/ML engines). The design of the proposed data structures is also coupled with the application thereof inside an online software platform<sup>[20]</sup> allowing to create a very short feedback loop and improve the resulting implementation based on feedback from thousands of users of the platform. Our standards facilitate the development of artificial intelligence tools that can reduce the dimensionality and complexity of the research work in materials science and chemistry, with the aim of eventually enabling inverse design. Our goal is to build collective intelligence utilizing contributions from a large audience of materials scientists (well beyond the select few experts in the computational field) and chemists in a controllable and high-level fashion.

The following section of this paper will introduce the simulation entities and the principles on which the model categorization scheme is based. The subsequent section illustrates the representation of the entities as data structures using a set of selected examples. In the fourth section, we discuss the relevancy and limitations of the categorization scheme and propose a community-driven approach to extend the categorization scheme. Finally, the principal conclusions are presented together with a perspective on future applications.

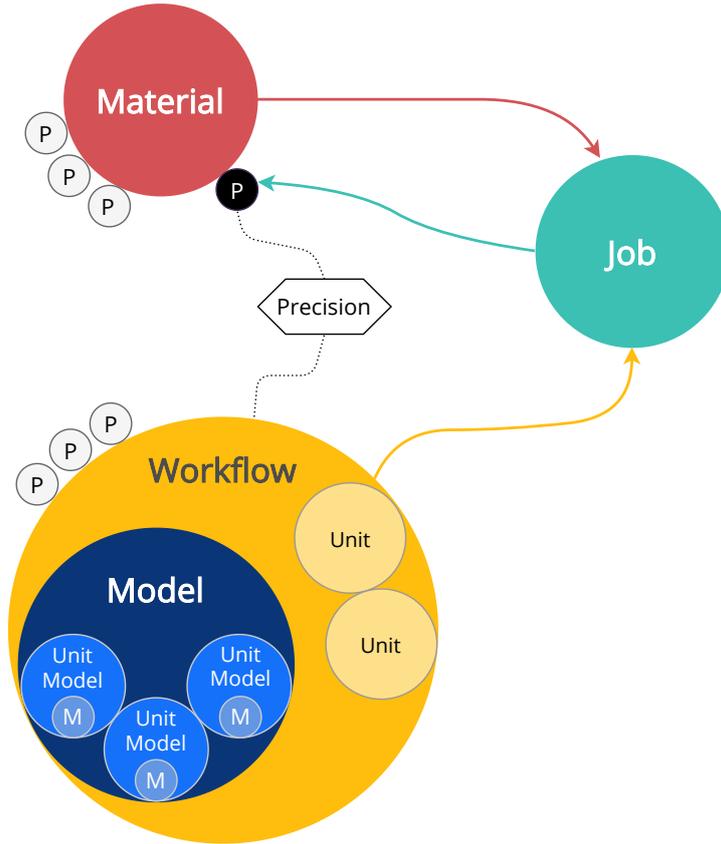
## 2 Methodology

### 2.1 General Approach

On an abstract level, we choose to represent a computational simulation in the form of several key entities (see Fig. 1). First of all, the *material* entity defines the chemical composition of the system under investigation. Although this entity is termed *material*, it may also constitute a non-periodic molecular species. The *workflow* entity organizes the sequence of tasks, for instance, execution of simulation software or input/output operations. The workflow also includes the specification of the theoretical model, which in turn is represented by the *model* entity. It should be noted that both the workflow and model entities are composed of reusable units, which may be combined in various ways. Applying the workflow on a material in order to produce one or more properties, i.e. connecting the workflow and materials entities, is achieved by the *job* entity. This entity does not only serve as a container for material and workflow entities but also stores more technical information related to high-performance computing and resource allocation. Properties may either be derived from existing entities (gray) or occur as a direct result of the job (black). In the latter case, one can associate a *precision* entity that is derived from the workflow. An important factor in determining the precision is the practical approach for solving the theoretical model, which is represented by the *method* entity (M). As part of a model (or unit model), it stores the selection of algorithms, thresholds, and other practical parameters. For the utilization on the Exabyte.io platform<sup>[20]</sup> and document-based, NoSQL (not only structured query language) databases in general, these entities are formulated as database schemas. Document-based, NoSQL databases are a convenient choice for such an application due to a few advantageous properties: (a) data that is accessed together is stored together (as opposed to joined from multiple data tables), (b) the organization of data can be as complex as one chooses thus supporting parent-child hierarchical structures, (c) data structures are not fixed and may be changed in response to new data models. The proposed database schemas have been implemented in the Exabyte Source of Schemas and Examples (ESSE)<sup>[21]</sup> using the JSON Schema notation (Draft-04)<sup>[22]</sup>.

### 2.2 Components and Entities

The ESSE module comprises several main schemas for simulation data, such as workflow, material, property, method, and model. The following sections focus on the latter two entities, while all other schemas are briefly summarized in Sec. 2.2.3.



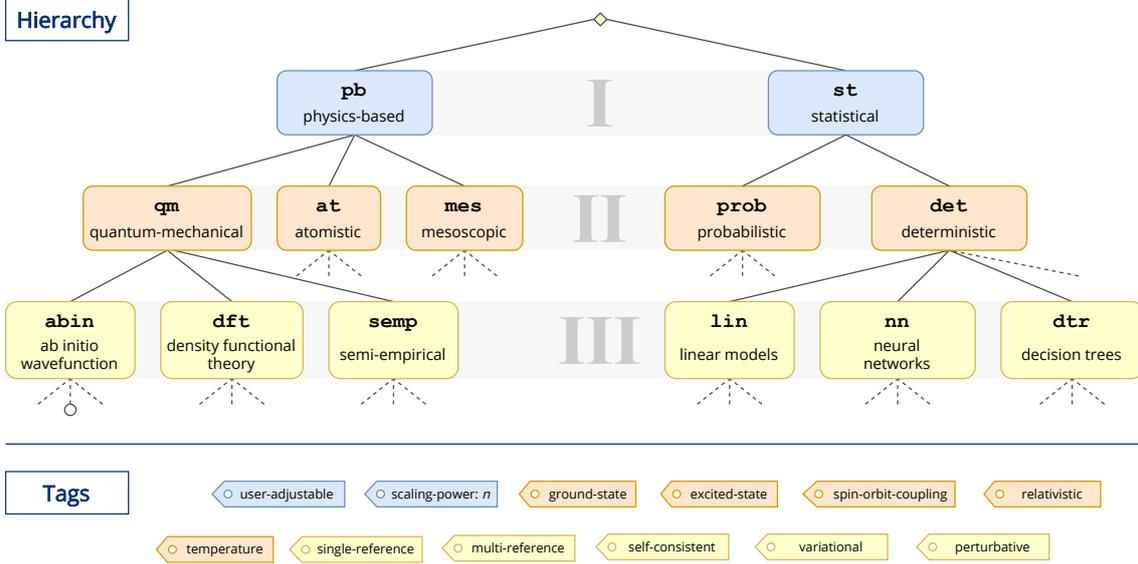
**Figure 1:** A visual representation of the key Entities considered in the current approach and relationships between them. A Job represents a research (computational) task involving a combination of Workflow and Material(s) Entities in order to produce initial *a priori* (gray) and target *a posteriori* properties (black) denoted by the letter "P". Resulting properties have a certain Precision. A Workflow is further made of individual workflow Units responsible for elementary computational operations (e.g. a "one-shot" run of simulation software). A Workflow has a (compound, in the general case) Model associated with it, which in turn is composed of Unit Models. The categorization of such Unit Models represents the focus of the current work.

### 2.2.1 Model

The function of the model schema is to define a given computational model as accurately as possible and to simultaneously store all of the necessary metadata. We chose to represent a given computational model in terms of one or more reusable, independent components, which we will refer to as *unit models* in the further course. The final model which is applied to a system is a combination of said unit models and termed *compound model*. We define a unit model as the smallest, logically consistent set of equations or operators associated with a central property (e.g. electronic energy). In practice, a unit model may not always be unambiguously defined, i.e. such a case may require further partitioning into a set of unit models that would go beyond the scope of this classification scheme. The classification, therefore, requires a subtle balance of exactness and pragmatism. The classification tree involves three main tiers (Figure 2) in order to presort the models into families of models, for instance *quantum mechanical* or *classical*. Following tier III the models are further divided into more specific categories, which are also organized hierarchically using **type**, **subtype** specifiers. The design of the schemas follows an object-oriented approach whereby schemas share fields through inheritance (by means of the `allOf` keyword). At each level of the classification tree, a given schema thus includes all categorization specifiers of the preceding levels and may serve as a prototype for the following level. One advantage of this concept is that changes to categories or implementation of new categories only occur locally and are propagated automatically to the lower levels.

In addition to the classification specifiers, each unit model comprises a **tag** field, which is also passed on through the categorization hierarchy (see Table 1 for examples). The tags describe attributes of the unit model not included in the categorization and indicate whether a modifier or augmentation has been applied to the unit model. We define a modifier as an addition to a model, which expands upon the underlying physical principle without fundamentally changing the working equations (e.g. in

a linear fashion). An example of a modifier is the inclusion of an additional external potential (e.g. due to point charges). An augmentation, on the other hand, is defined as an addition to a model, which does not change the underlying physical principles of the model. For instance augmentations include acceleration techniques such as resolution-of-the-identity or localization schemes (e.g. Edmiston-Ruedenberg localization<sup>[23]</sup>). Apart from the specifications above, the `tags` field may also hold user-defined labels.



**Figure 2:** Schematic representation of the tiers of the CateCom categorization. Tier I to tier III of the unit model classification hierarchy and descriptive tags colored according to the tier in which they first appear.

- Tier I At the primary level unit models are distinguished between *physics-based* (**pb**) and *statistical* (**st**). The latter category pertains to data-driven approaches which employ statistical relations in order to predict a result. Models based on fundamental laws of physics are assigned to the *physics-based* category even if the unit model heavily relies on statistical elements (see also Sec. 2.3).
- Tier II Within the *physics-based* group, the quantum mechanical (**qm**), atomistic (**at**) or mesoscopic (**mes**) category are used depending on the type of particle represented in the equations of the unit model. For instance, Kohn-Sham density functional theory (KS-DFT) naturally falls into the **qm** category due to its explicit dependence on electronic coordinates, while a force field such as CHARM22<sup>[24]</sup> only depends on atomic variables and is thus assigned to **at**. In the *statistical* group a unit model falls into the probabilistic (**prob**) category if the predicted result incorporates some aspect of random variation, whereas the deterministic category (**det**) is chosen if it does not.
- Tier III The quantum-mechanical models are further divided into three categories. The *ab initio* category (**abin**) comprises first-principle wavefunction models, such as Hartree-Fock theory, many-body perturbation theory, or coupled cluster theory, which do not require additional information about the system. With the electron density as the central quantum mechanical descriptor, the realizations of density functional theory are collected in the **dft** category. The *semi-empirical* category (**semp**) contains parametrized quantum-mechanical models, which usually only describe valence electrons for computational efficiency. As for the statistical model branch, Figure 2 shows how the deterministic models can be further subdivided using the example of three prominent machine learning model categories. Linear models (**lin**) span the space of models which assume a linear relationship between the input variables ( $\mathbf{X}$ ) and the dependent variable ( $\mathbf{y}$ ). The neural network category (**nn**) contains all models that are based on a neural network architecture, i.e. a network of interconnected processing nodes whereby connections between nodes are represented by weights. The decision tree category (**dtr**), on the other hand, comprises models which can generate a prediction based on recursively splitting the dataset into subsets. The outcome of such a procedure is then a linear acyclic graph of decision nodes and 'leaves' (endpoints, which do not split the data any further). The decision tree approach is usually applied in an ensemble (random forest model), which in the CateCom scheme is represented as a compound model. Other examples of deterministic

models not explicitly shown in Figure 2 support vector machines or clustering algorithms, such as k-means.

**Table 1:** Selected model attributes and the corresponding categorization levels.

Label	Explanation	Category
relativistic	Inclusion of relativistic effects.	pb
user-adjustable	The model contains additional parameters to fine-tune results.	pb
scaling-power: $n$	The model exhibits a formal scaling of $n$ -th power.	pb
self-consistent	Non-linearity in the model is solved through self-consistent optimization.	pb/qm
temperature	The model describes non-zero temperature effects.	pb/qm
excited-states	Access to electronically excited states.	pb/qm
spin-orbit coupling	The model accounts for spin-orbit coupling.	pb/qm
variational	The model follows the variational principle.	pb/qm
single-reference	The wavefunction is based on a single reference determinant.	pb/qm
multi-reference	The wavefunction is based on multiple reference determinants.	pb/qm
perturbative	The model contains elements of perturbation theory.	pb/qm/abin

### 2.2.2 Method

While the unit model pertains to the accuracy of a computational simulation, the *method* concerns its precision. The CateCom collection, therefore, includes a `method` schema for parameters concerning the computational methodology, such as convergence thresholds or hyperparameters (machine learning). As it is closely related to a model, method schemas are part of unit models and compound models. The method schema has three main attributes: associated with method parameters, method data, and precision. The `parameters` attribute holds a list of annotated control variables, which apart from the central key-value pair also encompasses a categorization keyword and, if applicable, a definition of the value’s unit. The method `data` attribute contains other input variables which may require additional files, such as user-generated pseudopotentials or basis sets.

In principle, the method schema contains all relevant information for the precision of a given choice of model and material. If one is able to formulate suitable scoring functions, the precision parameters can be turned into numeric features for a regression model. In conjunction with other factors, such as simulation time or memory usage, it would be highly desirable to predict an optimal model/method for a given material-property combination.

### 2.2.3 Other Entities

In the following, we briefly outline other notable ESSE data schemas. A more detailed definition can be found in Ref. 25. For materials data to be searchable, traceable, and reproducible, it is crucial to have a concise and informative way to describe materials and their properties. The *material* schema comprises descriptive properties that uniquely specify a material, such as Bravais lattice vectors and the unit cell basis. Note that the material schema is not limited to periodic systems and will also support molecular descriptors in an upcoming future release. A workflow defines the logical composition of simulation tasks that derive from one or several simulation engines or may take other forms such as Python scripts. The workflow as we define it is also hierarchically organized in three consecutive levels (from top to bottom): workflow, subworkflow, and workflow unit. In the *workflow* schema the logical composition is represented in terms of a directed acyclic graph (DAG), whereby each node is a workflow or subworkflow. A workflow may contain several subworkflows or other workflows. The organization of the simulation results is managed by the *properties* schema. Apart from the property data, this schema assigns a property group for easier access/findability and includes the unit of the property (if applicable).

## 2.3 Classification Rules

Since the classification and hierarchy of models involve some arbitrariness, we propose a set of classification rules to guide the categorization of hybrid models or edge cases. Although the rules listed below (Table 2) only pertain to a part of the categorization tree, they are intended to demonstrate how specific

cases can be distinguished and assigned to a category. Instead of aiming at a final set of categorization rules allowing to uniquely classify each computational model (and without discussing whether such an approach is even possible), we suggest the readers consider our approach below as a stepping stone toward a practically applicable implementation.

As an example, let us consider the Quantum Monte Carlo (QMC) model - an approach to find highly accurate solutions to the quantum many-body problem and is often used to study materials and molecular systems.<sup>[26,27]</sup> Due to its stochastic foundation, QMC results involve a quantifiable random error. Although stochastic sampling is an important component of the model, the aim of most QMC models is to solve for the ground state wave function (or density matrix).<sup>[26]</sup> As such the model was assigned to the *ab initio* wave function model category (**abin**). This example prompts another guideline to introduce for the categorization of computational models: in case of ambiguity, one should categorize a model based on its objective (e.g. solving the Schrödinger equation) rather than its components or derivation.

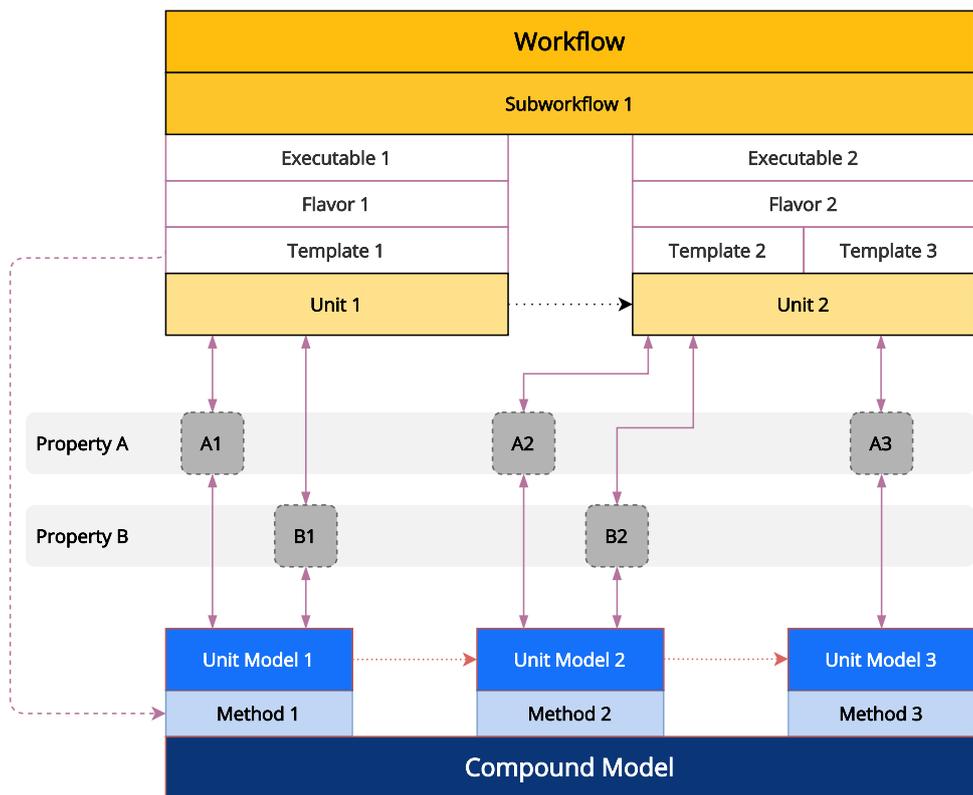
Another guideline concerns the relationship of categorization tiers and explicit realizations, i.e. instances of unit models. As outlined above, tier I to tier III serve as identifiers for groups of models. To guarantee a consistent usage of the unit model object, it should thus be avoided to equate a unit model instance with one of these three tiers.

**Table 2:** List of CateCom classification rules.

No.	Category	Categorization Rule
1	<b>pb</b>	The model is based on physical laws.
1.1	<b>pb/qm</b>	The model depends on electronic coordinates or involves an electronic or nuclear wavefunction.
1.1.1	<b>pb/qm/abin</b>	The model is based on first-principle wavefunction approximations.
1.1.2	<b>pb/qm/dft</b>	The model is based on density functional theory.
1.1.3	<b>pb/qm/semp</b>	The model only treats valence electrons explicitly and/or involves parametrization of two-electron integrals.
1.2	<b>pb/at</b>	The model depends on atom (nuclear) coordinates only (without using wavefunctions).
1.3	<b>pb/mes</b>	The model involve a conflated representation of particles.
2	<b>st</b>	The model predicts results based on data rather than physical laws.
2.1	<b>st/prob</b>	The model involves randomness and cannot predict a result with an exact formula. Often the result is characterized by a mean and a distribution.
2.2	<b>st/det</b>	The model does not include randomness and always gives the same prediction.
2.2.1	<b>st/det/lin</b>	The model comprises a linear combination of features (or kernels).
2.2.2	<b>st/det/nn</b>	The model employs a neural network architecture.
2.2.3	<b>st/det/dtr</b>	The model is based on decision trees.

## 2.4 Entity Interoperation

This section describes how the entities presented in Sec. 2.2 interact to facilitate the research process. In particular, we present the relation between workflow, compound model, and properties. Starting with the workflow, a hierarchical organization similar to the model entity is used. A subworkflow, which is associated with one application or software package, contains one or several workflow units. There are various types of workflow units each managing a different role, for instance, input/output operations or conditional operators (see Ref. 25 for a full list of types). The workflow unit type shown in Fig. 3 is of the *execution* type, which refers to an executable of the simulation software package. An executable may require one or several input files, which are generated by templates. The *flavor* entity, which is part of the workflow unit, matches templates to an executable for the purpose of obtaining a selected set of properties. In a broader sense, the flavor therefore represents a specific group of simulations, for instance, single-point calculations or geometry optimizations. The *compound model* defines the model and simulation parameters for a given subworkflow. As outlined in Sec. 2.2.1 each compound model is comprised of one or more unit models. Although a unit model is only associated with one workflow unit, the opposite does not apply. Allowing one workflow unit to map to several unit models makes this framework very flexible and consistent across different simulation packages. Each unit model as well as



**Figure 3:** Interoperation of a workflow (top) and compound model (bottom). The general example involves a workflow consisting of two workflow units and a compound model with three unit models. The modular approach of unit models allows a one-to-many mapping of workflow units to unit models and tracking of intermediate properties (e.g. A1 and B1). The entities presented here follow the color scheme of Figure 1.

the compound model itself contains a method object which stores information related to *how* a model is solved and is used to populate templates. Finally, associating properties with individual workflow units and unit models enables the user to monitor the progress of a given property across the compound model. Furthermore, since some unit models may not give rise to a certain property (cf. Unit model 3 and Property B in Fig. 3), the concept allows for convenient access to the last (or other criteria) occurrence of the property.

### 3 Examples

In the present section, we elaborate the above-introduced data structures for Unit Model, Compound Model, and Method by means of common-use examples. For the sake of brevity, only a few examples are shown and we refer the reader to the ESSE repository<sup>[21]</sup> for an extensive collection of examples. It should be noted that the reoccurring key "`_id`" is not a component of the CateCom data structure, but pertains to the document-based database software and will thus be omitted in the discussion below.

#### 3.1 Unit Models

Kohn-Sham Density Functional Theory (KS-DFT)<sup>[28,29]</sup> is a widely used quantum-mechanical model in material science as well as molecular science. As the solution of KS-DFT is often used as an input for a subsequent model, for instance perturbation theory<sup>[30]</sup> or as an alternative reference determinant<sup>[31]</sup> in *ab initio* wavefunction models, it is well suited to be represented as a unit model. Listing 1 shows the unit model data structure for a generalized gradient approximation (GGA) KS-DFT model using the Perdew-Burke-Ernzerhof (PBE)<sup>[32]</sup> exchange-correlation functional. An example of the range-separated hybrid functional HSE06<sup>[33]</sup> is presented in the appendix (see Listing A.1). According to the CateCom approach introduced in Sec. 2.2.1, the three tiers for KS-DFT are *physics-based*, *quantum-mechanical* and *density functional theory*, respectively. Since there exist multiple realizations of DFT (for instance orbital-free DFT<sup>[34,35]</sup>), the `type` field further specifies the variation of DFT. In the data structure each

tier (as well as type and subtype if applicable) is mapped to a simple object which contains a human-readable descriptor (`name`) and a machine-readable token (`slug`). The annotation fields introduced in Sec. 2.2.1 (`augmentation`, `modifier` and `tag`) give further descriptive information about the unit model and facilitate the search for unit models. In addition, references to the literature can be given using the `reference` field (left empty in Listing 1 for brevity). Each unit model includes a so-called `flowchartId` field, which is used to uniquely identify a unit model within a compound model and which serves as a reference for representing the compound model as a directed acyclic graph (DAG).

**Listing 1:** Example of a unit model data structure for Kohn-Sham Density Functional Theory. Apart from the categorization fields the data structure also holds a detailed representation of the exchange-correlation functional (here: PBE).

```

1 {
2   "tier1": {
3     "name": "physics-based",
4     "slug": "pb"
5   },
6   "tier2": {
7     "name": "quantum-mechanical",
8     "slug": "qm"
9   },
10  "tier3": {
11    "name": "density functional theory",
12    "slug": "dft"
13  },
14  "type": {
15    "name": "Kohn-Sham DFT",
16    "slug": "ksdft"
17  },
18  "subtype": {
19    "name": "Generalized Gradient Approximation",
20    "slug": "gga"
21  },
22  "functional": {
23    "name": "Perdew-Burke-Ernzerhof 1996",
24    "slug": "PBE",
25    "components": [
26      {
27        "name": "PBE",
28        "slug": "pbe-c",
29        "type": "correlation",
30        "fraction": 1.0
31      },
32      {
33        "name": "PBE",
34        "slug": "pbe-x",
35        "type": "exchange",
36        "fraction": 1.0
37      }
38    ]
39  },
40  "tags": [
41    "self-consistent"
42  ],
43  "method": {...},
44  "reference": {...},
45  "flowchartId": "u123fndff333",
46  "_id": "ffffff55555"
47 }

```

Apart from the categorization and annotation fields, the CateCom approach also supports additional fields that are exclusive to a certain unit model. For instance, the multitude of density functional approximations (DFA) warrants a separate key (termed `functional`), which captures the different categories of DFAs. The functional object contains identifier fields `name` and `slug`, which include the commonly used approximations and acronyms for DFAs. Many exchange-correlation functionals are comprised of several components (here referred to as *unit functionals*), for instance, separate approximations for exchange and correlation as well as a fraction of exact exchange (hybrid functionals). The `functional` data structure lists these contributions under the key `components`. Each component is characterized by nominal descriptors (`name` and `slug`), the type of functional component (*vide infra*) and the fraction with which the component enters the model. Besides the two unit functional types presented in Listing 1 a unit functional may adopt the type of a non-separable exchange-correlation functional (e.g.

GAM<sup>[36]</sup>), a kinetic energy functional (e.g. Thomas-Fermi<sup>[37,38]</sup>), or a non-local correlation functional such as VV10<sup>[39]</sup>. The `method` field in this example is left empty as it will be discussed separately in the next section.

## 3.2 Methods

As mentioned in Sec. 2.2.2, each unit model (and compound model) comprise method data that pertains to the precision of the computational simulation. For the above example of KS-DFT, this data holds, among others, information about the employed basis (e.g. plane wave energy cutoff) and integrals (e.g. k-point grid size). The `method` data structure is organized as follows. First, a simple categorization is given by the two required fields `type` and `subtype`. The type/subtype categorization for the methods is a preliminary solution, which -if the need arises - will be replaced by a more elaborate system comparable to the unit model categorization. In the specific example of Listing 2, the given type defines the use of plane waves in combination with pseudopotentials, while the subtype further specifies the use of ultra-soft pseudopotentials (`us`).

**Listing 2:** Data structure for the `method` entity. This example pertains to the use of plane waves and pseudopotentials for the description of the electronic structure. The data structure also includes non-default input parameters, two of which are deemed to have an effect on the overall precision.

```

1 {
2   "type": "pseudopotential",
3   "subtype": "us",
4   "parameters": [
5     {
6       "name": "ecutrho",
7       "value": 1e-8,
8       "categories": ["wavefunction", "precision"],
9       "units": "Ry"
10    },
11    {
12      "name": "ecutwfc",
13      "value": 1e-6,
14      "categories": ["wavefunction", "precision"],
15      "units": "Ry"
16    },
17    {
18      "name": "occupations",
19      "value": "smearing",
20      "categories": ["brillouin-zone"]
21    }
22  ],
23  "precision": [
24    {
25      "name": "ecutrho",
26    },
27    {
28      "name": "ecutwfc",
29    }
30  ],
31  "data": {
32    "searchText": "",
33    "pseudo": [...]
34  }
35 }

```

The method data structure also holds a list of non-default input variables in the `parameters` field. Each parameter is given in the form of an annotated key-value pair containing the name of the input variable (key) as it appears in the input file, its value, the corresponding categories (*vide infra*) and, if applicable, the unit of the value. As each *flavor* (cf. Sec. 2.4) is associated with a set of default input variables, the `parameters` field only needs to store input variables which deviate from the default value. The parameters in Listing 2 stem from a plane wave DFT calculation employing the Quantum ESPRESSO software package.<sup>[40]</sup> In particular, they define the kinetic energy cutoff for the charge density (`ecutrho`) and the wavefunction (`ecutwfc`) as well as the approach for sampling the Brillouin-zone (`occupations`). Parameters labeled with the *precision* category are considered to influence the precision of the corresponding unit model and thus fulfill a special role. For example the precision score (cf. Sec. 2.2.2) is calculated based on these parameters. For fast access, the `precision` field collects the names of these input parameters. The `data` field stores additional data specific to the method. For

example, in case of the plane-wave pseudopotential method, the data field contains the pseudopotentials themselves (`pseudo`). In addition, `data` contains a keyword (`searchText`) for filtering or searching the `data` attribute.

### 3.3 Compound Models

Following the definition of the CateCom schema, we illustrate its practical use by an example, which corresponds to established models for materials and molecules, respectively. In addition, the ESSE repository contains an extensive set of examples covering classical mechanics and machine learning models. Due to the variety of unit models also multi-level models, such as the combined quantum mechanical/molecular mechanical (QM/MM) approach<sup>[41]</sup>, can be realized as compound models.

#### 3.3.1 DFT+GW Model

Although DFT is arguably one of the most popular electronic structure models, its deficiencies inhibit an accurate simulation of some experiments, for instance, photoemission spectroscopy.<sup>[42]</sup> The *GW* approximation provides a way to improve upon the single-particle states obtained from DFT in a perturbative fashion.<sup>[30]</sup> While the GW approximation allows for an accurate description of "charged excitations", i.e. electronic excitations whereby an electron is added or removed from the *N*-electron system, neutral excitations, which preserve the number of electrons in the system, can be described using the Bethe-Salpeter equation (BSE).<sup>[43]</sup> The starting point for the GW approximation are the eigenfunctions  $\{\phi^{\text{KS}}\}$  and eigenvalues  $\{\epsilon^{\text{KS}}\}$  of the KS-DFT mean-field Hamiltonian (Hartree-Fock or DFT). The dynamically screened Coulomb potential and the single-particle energy levels obtained from GW, in turn, are input quantities for the calculation of optical excitations using the BSE. As such, the cascade of DFT, GW, and BSE are well suited to be expressed in terms of unit models.

**Listing 3:** Compound model data structure for the combination of Density Functional Theory and GW.

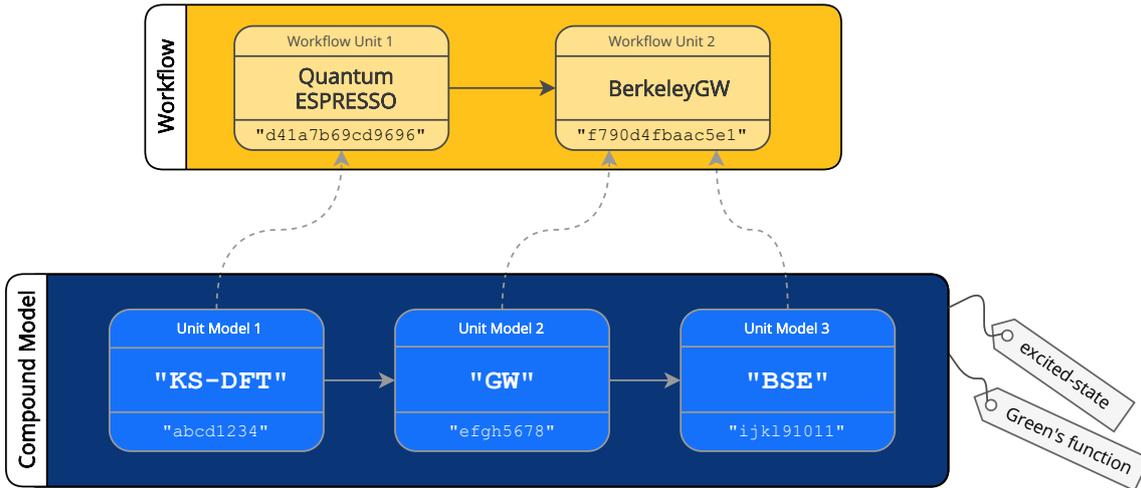
```

1 {
2   "_id": "com111",
3   "name": "KS-DFT + GW + BSE",
4   "modelGraph": [
5     {
6       "name": "KS-DFT",
7       "slug": "pb-qm-dft-ksdft",
8       "flowchartId": "abcd1234",
9       "next": "efgh5678",
10      "head": true,
11      "workflowUnitId": "d41a7b69cd9696"
12    },
13    {
14      "name": "GW",
15      "slug": "pb-qm-abin-gw",
16      "flowchartId": "efgh5678",
17      "next": "ijkl91011",
18      "head": false,
19      "workflowUnitId": "f790d4fbaac5e1"
20    }
21  ],
22  {
23    "name": "BSE",
24    "slug": "pb-qm-abin-bse",
25    "flowchartId": "ijkl91011",
26    "next": null,
27    "head": false,
28    "workflowUnitId": "f790d4fbaac5e1"
29  }
30 ],
31 "method": {...},
32 "tags": [
33   "excited-states",
34   "Green's function"
35 ]
36 }

```

Turning to the compound model data structure (Listing 3), a slightly different object composition can be seen. Since most of the information pertaining to the overall model is stored within the unit model data structures, the compound model does not need to repeat the information in its own data structure.

Consequently, the compound model holds the arrangement of unit models (`modelGraph`) and a global methods object (`method`). The `modelGraph` field contains a list of nodes, each representing a unit model. Each node contains the necessary fields to construct a directed acyclic graph, i.e. a unique identifier (`flowchartId`), a pointer to the next object (`next`) and an boolean indicator of the first node (`head`). In addition, a human-readable (`name`) and a machine-readable (`slug`) label are included. Finally, each unit model node also maps to a workflow unit by means of the `workflowUnitId` key, which does not necessarily have to be unique to each node. For instance, in the above example, three unit model nodes are mapped to two workflow units. A possible scenario for this is calculating the KS-DFT solution from one software package (e.g. Quantum ESPRESSO) and subsequently applying GW and BSE using another (e.g. BerkeleyGW<sup>[44]</sup>). Just like the unit model, the compound model data structure also includes a `method` key, which refers to a *global* method configuration.



**Figure 4:** Relation of compound model and workflow for the many-body theory example of GW and Bethe-Salpeter Equation (BSE) based on a reference from Kohn-Sham density functional theory (KS-DFT). In this example, the calculation of KS-DFT reference and application of many-body theories is performed in two steps and thus two independent workflow units.

## 4 Discussion

The CateCom approach laid out in previous sections introduces a model data structure combined with a systematic categorization of computational models. The model data structure is designed to be composed of one or more reusable components (unit models), each of which is assigned a model category. The model and other Entities representing the computational Workflow, Properties, and Materials altogether form the research-work-related data and metadata. Our approach is not meant as a final "polished" solution, but rather as a proof-of-concept but practically deployable implementation. Admittedly, it is largely limited to physics-based models in the current implementation and is, perhaps, heavily biased toward atomistic and nanoscale simulations. The uniqueness of categorization - whether to enforce it and how - is another topic that requires further clarification. Without uniqueness, the chosen categorization scheme can be seen as one linear realization of a non-linear "model graph". The rest of this section discusses several important aspects of CateCom.

### 4.1 The Material-Model-Property Categorization

#### 4.1.1 Material, Model, and Property relationship

For any practical applications, the fidelity of the modeling approaches can usually only be established within a certain class of materials and their associated properties. For example, it is well known that conventional Density Functional Theory significantly under-estimates the electronic band gap values in semiconductors, while providing adequate predictions for other properties such as lattice constants and/or vibrational spectra. Therefore a certain fidelity metric predicting the quality of a particular model would only stand with respect to certain material and property types.

### 4.1.2 Material and Property Categorization

To facilitate data-driven science, a coupled approach is needed where materials (and chemicals) have to be categorized as well as their derived properties. This way one can construct the associative relationships that can assist in identifying the most successful combinations of Materials, Workflow/Model, and Properties. Following the example in the previous section, such an approach should be able to identify, for example, that for III-V semiconductors a Model/Workflow containing both Density Functional Theory and GW Approximation provides higher fidelity than Density Functional Theory alone. Although the exact nature of such categorization is a topic of a separate discussion, and the categorization related to each of the entities can be interdependent, we still believe that starting with the model categorization provides a viable and practically useful first step.

## 4.2 FAIR Principles

There have been many efforts to collect and systematically organize research data to build large publicly available datasets.<sup>[5,13–19]</sup> Simultaneously, a new paradigm for scientific discovery, data-driven science, emerged, which aims to detect patterns or anomalies in these types of datasets.<sup>[5,45,46]</sup> In a cooperative effort including representatives from academia, industry, funding agencies, and scholarly publishers, the FAIR guidelines<sup>[47,48]</sup> (findable, accessible, interoperable, and reusable) were developed in order to enhance data reusability. The following subsections demonstrate how the CateCom approach ties in with the FAIR principles.

### 4.2.1 Findability

According to the FAIR guiding principles<sup>[47]</sup>, findability involves the use of globally unique and persistent identifiers. As described in Sec. 3.1, each unit model in the CateCom scheme encompasses such a unique identifier in the form of the `flowchartID`. Additionally, unit models and method parameters are enriched with `tags` metadata facilitating a search or filtering for certain properties. In this way, associations between unit models can be made, which are not represented in the CateCom tree. For instance, the `scaling-power-3` tag may be used to filter all unit models which formally exhibit a cubic scaling even if they are located in different categorization branches.

### 4.2.2 Accessibility

One aspect of the FAIR guidelines pertains to the accessibility of the data and the protocols applied in that process. The CateCom approach implements database schema corresponding to the Unit Model, Compound Model, or workflow using the JSON schema vocabulary. As such, the resulting Entity objects are represented in the widely adopted JSON format. The model Entities are stored in a document-based database and are thus easily retrievable using the unique database identifier or through a database search. The CateCom schemas themselves are publicly available since they are part of the open-source ESSE repository<sup>[21]</sup>. In addition, the JSON format is widely used on the world wide web in many web applications. The latter implementation provides packages for two of the currently most widely used languages - Python and JavaScript - with an intent to allow for easy adoption in software development efforts providing user interface components that in turn aid Accessibility.

### 4.2.3 Interoperability

Interoperability encompasses the integration with other data and cross-functional cooperation with applications. The CateCom Unit Models do not possess dependencies to specific software implementations of the models, such that the models can, in principle, be associated with any software package (given that the model is implemented therein). Furthermore, storing the entities as JSON objects has the advantage that there are plenty of resources available which directly accept this format or are able to convert it to a different format (see also Sec 4.4). The software implementation mentioned in the previous section provides additional opportunities for building interoperable systems.

### 4.2.4 Reusability

The CateCom scheme also addresses the reusability of data. In particular, the partitioning of models into unit models serves the purpose of reusing components that make up a model. A good example for this are many-body perturbation theory models which generally require the solution to an unperturbed

Hamiltonian  $\hat{H}_0$ . Consequently, the Unit Model corresponding to the unperturbed Hamiltonian can be combined with different perturbation theory models or, in the case of Hartree-Fock theory, with post-HF wavefunction models such as configuration interaction (CI). Furthermore, the storage of the method data plays a crucial part in recording the provenance of the final property data.

## 4.3 Predictive AI/ML

### 4.3.1 Avoidance of Duplicates

With the ability to quantify and store the metadata about the digital approaches comes the ability to avoid repetition. In case a particular workflow/model/property combination, for example - pseudopotential Density Functional Theory with a certain wavefunction and charge density cutoffs - has been applied to a specific material, our proposed data management model will be able to provide a way to generate a unique fingerprint. Based on such a unique fingerprint any further duplicate attempts can be avoided leading to improved efficiency of the research work.

### 4.3.2 The AI "Chemist-in-the-cloud"

The ultimate goal of the categorization described here is to enable the creation of an AI-powered digital computational chemist/materials scientist ("brain") able to suggest the best model/method combinations for characterizing materials. The complexity of materials science and chemistry is to a large degree defined by the diversity of the problem sets and the parametric conditions associated with them. Although computational techniques have been around for over half a century, the ability to apply them successfully with high fidelity still has a significant "art" component requiring very specialized knowledge limited to a select group of scientists only. And even this select group in practical applications often relies on their "intuition" derived through years of experience in dealing with specific problems rather than purely deterministic. With the help of the categorization scheme proposed and provided sufficient training data based on expert decisions such intuition can be instead represented as data-driven AI/ML approaches.

## 4.4 Community, Ecosystem, and Future Outlook

### 4.4.1 The Global Digital Ecosystem for Materials R&D

The categorization framework and associated ontologies both represent fundamental critical steps in the implementation of a global digital ecosystem for materials R&D. Having data standards is an important fundamental step in the design and implementation of the data- and software infrastructure for such an ecosystem. Object-oriented design for the entities and data structures naturally enables modularity when building the software components of such ecosystem, and greatly streamlines its implementation and long-term maintenance. A version of the present categorization framework have been previously deployed as part of an online software platform<sup>[20]</sup> with applications demonstrated for multiple use cases, including metallic alloys<sup>[49]</sup>, electronic properties of semiconductors<sup>[50,51]</sup>, vibrational properties of materials<sup>[52]</sup>, adsorption and catalysis in zeolites<sup>[53]</sup>, adhesive strength of composite materials<sup>[54]</sup>, and beyond. Materials R&D spans a complex and multi-dimensional landscape, and requires an extremely large variety of characterization data at multiple time- and length scales. Once obtained, the data must be stored and managed in an efficient way. As more and more of materials research is performed in a way that involves digital handling of data, ontologies and categorization becomes important. This will facilitate the availability of ever increasing amounts of materials data on the web with contributions from the global community.

### 4.4.2 Community Contributions

Simulation scientists are able to resort to a myriad of statistical and physics-based models, whereby the number of models is constantly growing. This circumstance makes a systematic mapping of the "model landscape" rather difficult for a small team since profound knowledge of a model is required in order to systematically arrange its properties and variants. Of course, expert knowledge is also indispensable for identifying reusable components of these models and examining edge cases that may fit more than one category. Thus, a more effective approach is to involve the community of experts directly in the maintenance and expansion of the categorization scheme. To this end, we propose to follow the collaborative strategy typical for code development platforms such as GitHub. These platforms also allow

interested contributors to discuss new features (e.g. a new category) and raise issues about existing ones. In practice, a contributor first obtains their own server-side clone ('fork') of the original repository (e.g. ESSE<sup>[21,25]</sup>). The implementation of new features, such as a new unit model, is then carried out in a feature branch located in the cloned repository. Once the new feature is ready, the contributor then issues a request for the integration of the new feature ('pull request') to the maintainer of the original repository. At this stage details of the new feature can be discussed and modified until the maintainer accepts the incoming changes.

#### 4.4.3 Interfacing with other approaches

Of course, the task of developing a global digital ecosystem for materials research and development cannot be accomplished without involving a global community and interfacing with other efforts. Despite recent comprehensive approaches<sup>[12,55]</sup>, it is still common for the materials science community to develop standards which are tailored to a specific sub-branch of research. Such "artisanal"<sup>[56]</sup> approach has led to several competing standards with a relatively small impact on the field. Our goal is to provide a common denominator allowing the key contributors to realize their ambitions while at the same time facilitating the level of quality required for practical real-world applications.

CateCom has a natural connection to the principle of an ontology and the similarities facilitate interoperation of CateCom with ontologies defining model classes, for the ontology-based data access (OBDA).<sup>[57-59]</sup> Since ontology vocabulary is based on different formats (RDF, OWL, etc.), OBDA usually requires an access interface translating queries and responses.<sup>[58,59]</sup> Ontologies may also be used for semantic annotation, i.e. metadata enrichment.

#### 4.4.4 Example Interfaces with other approaches

The Materials Design Ontology (MDO) defines a *Computational Method* class which as of this writing is limited to density functional theory (DFT) and Hartree-Fock (HF) theory. Nonetheless, the definition of these models shares commonalities with the CateCom scheme. For instance, the DFT class of MDO has a *Exchange Correlation Energy Functional* property implementing the most common groups of density functional approximations which are also supported by the KS-DFT unit model in CateCom. Part of the CateCom unit model (in JSON format) can thus be mapped to RDF format in order to be used with MDO. A similar mapping to the RDF format using a SPARQL-Generate script<sup>[60]</sup> has been described in Ref. 12. Furthermore, a model class is also defined as part of the Elementary Multiperspective Material Ontology (EMMO).<sup>[11]</sup> Although specific models (e.g. DFT) are not explicitly represented, its subclasses are similarly organized as tier I and tier II of the CateCom scheme, for instance *DataBasedModel* and *PhysicsBasedModel* exactly correspond to the *st* and *pb* categories (cf. Figure 2). Interoperation of EMMO and CateCom could, for instance, be achieved by annotation or translation as described above.

In practical terms, any interfacing most likely involves a conversion between two database schemas. The Novel Materials Discovery (NOMAD<sup>[5]</sup>) metainfo schema defines the structure of material-science-related data. The schema contains a very extensive list of entities and properties, such that mapping is not limited to CateCom unit models but can in principle extend to other ESSE<sup>[21,25]</sup>, Entities (e.g. Method, Property). Apart from a one-to-one mapping, one could also populate the CateCom data structures based on a descriptive string. For instance, JARVIS-DFT database<sup>[18]</sup> contains a *functional* property, which contains enough information in order to be converted to a *ksdft* CateCom model object. Such object generation might not be fully complete, neglecting the method object entirely. Other approaches such as OPTIMADE<sup>[55]</sup> provides a universal application programming interface (API) to access material data across several databases. The OPTIMADE specification always includes a structure attribute, whereas properties other than structural or chemical information are provider-specific. As a consequence, model-related metadata (and therefore the mapping to CateCom) may or may not be available.

#### 4.4.5 Future Outlook

In our vision, CateCom presents a fundamental building block in facilitating mainstream data-driven research in materials science and chemicals. Our goal is to engage a large community of people possessing specialized knowledge about materials and chemicals in digital work resulting in the creation of novel AI/ML techniques. We see that community effort is critical in obtaining the "critical mass" of data and creating network effects allowing to sustain the effort for the long term. To understand the future outlook, we draw analogies with the Computer-Aided Design and Electronic Design Automation

industries. In both, as the transition from exploratory science-centric to practically applied engineering-focused research work was progressing, the number of data representation standards was consolidated to 3-5. These consolidated standards emerged behind the software development efforts amassing the largest user community - such as AutoDesk, Synopsys, Dassault, etc. We expect a similar progression of events to happen for data-driven digital materials R&D in the near future.

Apart from the general goal, there are also more technical questions that could be addressed in future work. More specifically, the current approach allows the combination of unit models to compound models without any restrictions. In order to prevent improper combinations, a strategy for interfacing unit models is needed. A potential solution would be to track input and output quantities of each model, such that combinations can be evaluated in terms of the intersection of input and output quantities. Another point concerns the representation of time-dependent simulations. Future work should thus examine whether unit models may include a "time-propagation" operator or whether a compound model analog (e.g. "dynamical compound model") would be a viable option.

Regarding the categorization, it would be desirable to extend an existing ontology or create a new ontology for the multitude of computational models spanning both physics-based and data-driven models. Such an ontology would be helpful for building a knowledge graph of materials science research containing the semantic relationships between material, model, and property entities.

## 5 Conclusion

We introduced an approach for the categorization of computational models in conjunction with database schemas representing the Models and Methods. The proposed data-centric categorization scheme follows an object-oriented design concept, whereby a given model is expressed in terms of reusable, indivisible components (Unit Models). This modular unit model approach allows for a consistent description of model properties across software packages and is able to describe multi-level models, such as QM/MM. The data structures derived from the proposed schemas have been elucidated based on the examples, such as Density Functional Theory (DFT), GW Approximation, and similar. It has been demonstrated how the CateCom scheme complies with the FAIR guiding principles. In particular, possible mechanisms for the interoperation of CateCom with other approaches of the digital materials science ecosystem have been presented. In order to manage limit cases and guide new additions of categories, a set of categorization rules has been presented.

In its current state of development, CateCom represents a proof-of-concept with an emphasis on physics-based models. With the aim of leveraging expert knowledge, we discussed a community-driven approach for the extension of the CateCom scheme. Just as many other categorization efforts the CateCom scheme does not claim uniqueness with regards to the chosen categories. The organization of model data as presented herein allows for several convenient benefits such as transferability of a given model from one problem to another. The model categorization also allows for the generation of unique fingerprints, which facilitate the research process by avoiding duplicates. We share the current implementation of the categorization as part of an open-source online codebase<sup>[21]</sup> and demonstrate some of the applications of the underlying data infrastructure in the online platform<sup>[25]</sup>.

The ideas expressed in the present manuscript build upon the Materials Genome Initiative<sup>[61]</sup>, and are designed to facilitate collaboration between materials scientists, chemists, computer/data scientists to create, deploy and analyze a set of curated methodologies to rapidly study materials at multiple time- and length scales. In our view, the present data convention is aimed to facilitate the next generation of computer-aided design tools and enables advanced R&D capabilities that facilitate the development of new kinds of products in critical industries including semiconductor, photovoltaics, energy storage, oil & gas, specialty chemicals, aerospace and automotive and others and has the potential to transform the materials sector at large.

## References

- [1] Monya Baker. Reproducibility: Seek out stronger science. *Nature*, 537(7622):703–704, September 2016. ISSN 0028-0836. doi: 10.1038/nj7622-703a.
- [2] J A Pople. Two-Dimensional chart of quantum chemistry. *J. Chem. Phys.*, 43(10):S229–S230, November 1965. ISSN 0021-9606. doi: 10.1063/1.1701495.

- [3] Martin Karplus. Three-dimensional “pople diagram”. *J. Phys. Chem.*, 94(14):5435–5436, July 1990. ISSN 0022-3654, 1541-5740. doi: 10.1021/j100377a002.
- [4] György Tarczay, Attila G Császár, Wim Klopper, and Harry M Quiney. Anatomy of relativistic energy corrections in light molecular systems. *Mol. Phys.*, 99(21):1769–1794, November 2001. ISSN 0026-8976. doi: 10.1080/00268970110073907.
- [5] Claudia Draxl and Matthias Scheffler. NOMAD: The FAIR concept for big data-driven materials science. *MRS Bull.*, 43(9):676–682, September 2018. ISSN 0883-7694, 1938-1425. doi: 10.1557/mrs.2018.208.
- [6] MolSSI. Molssi/qcschema. URL <https://github.com/MolSSI/QCSchema>. Date accessed: 2021-9-16.
- [7] OWL 2 web ontology language document overview (second edition). <https://www.w3.org/TR/owl2-overview/>. Accessed: 2021-9-16.
- [8] Nenad Krdzavac, Sebastian Mosbach, Daniel Nurkowski, Philipp Buerger, Jethro Akroyd, Jacob Martin, Angiras Menon, and Markus Kraft. An ontology and semantic web service for quantum chemistry calculations. *J. Chem. Inf. Model.*, 59(7):3154–3165, July 2019. ISSN 1549-9596, 1549-960X. doi: 10.1021/acs.jcim.9b00227.
- [9] Weerapong Phadungsukanan, Markus Kraft, Joe A Townsend, and Peter Murray-Rust. The semantics of chemical markup language (CML) for computational chemistry : CompChem. *J. Cheminform.*, 4(1):15, August 2012. ISSN 1758-2946. doi: 10.1186/1758-2946-4-15.
- [10] Dennis G Thomas, Rohit V Pappu, and Nathan A Baker. NanoParticle ontology for cancer nanotechnology research. *J. Biomed. Inform.*, 44(1):59–74, February 2011. ISSN 1532-0464, 1532-0480. doi: 10.1016/j.jbi.2010.03.001.
- [11] European Materials Modelling Council (EMMC). Elementary multiperspective material ontology, 2021. URL <https://github.com/emmo-repo/EMMO/>. Date accessed: 2021-9-16.
- [12] Huanyu Li, Rickard Armiento, and Patrick Lambrix. An ontology for the materials design domain. In *The Semantic Web – ISWC 2020*, Lecture Notes in Computer Science, pages 212–227, Cham, Switzerland, November 2020. Springer. ISBN 9783030624651, 9783030624668. doi: 10.1007/978-3-030-62466-8\_14.
- [13] Stefano Curtarolo, Wahyu Setyawan, Gus L W Hart, Michal Jahnatek, Roman V Chepulskii, Richard H Taylor, Shidong Wang, Junkai Xue, Kesong Yang, Ohad Levy, Michael J Mehl, Harold T Stokes, Denis O Demchenko, and Dane Morgan. AFLOW: An automatic framework for high-throughput materials discovery. *Comput. Mater. Sci.*, 58:218–226, June 2012. ISSN 0927-0256. doi: 10.1016/j.commatsci.2012.02.005.
- [14] Anubhav Jain, Shyue Ping Ong, Geoffroy Hautier, Wei Chen, William Davidson Richards, Stephen Dacek, Shreyas Cholia, Dan Gunter, David Skinner, Gerbrand Ceder, and Kristin A Persson. Commentary: The materials project: A materials genome approach to accelerating materials innovation. *APL Materials*, 1(1):011002, July 2013. doi: 10.1063/1.4812323.
- [15] James E Saal, Scott Kirklin, Muratahan Aykol, Bryce Meredig, and C Wolverton. Materials design and discovery with High-Throughput density functional theory: The open quantum materials database (OQMD). *JOM*, 65(11):1501–1509, November 2013. ISSN 1543-1851. doi: 10.1007/s11837-013-0755-4.
- [16] Camilo E Calderon, Jose J Plata, Cormac Toher, Corey Oses, Ohad Levy, Marco Fornari, Amir Natan, Michael J Mehl, Gus Hart, Marco Buongiorno Nardelli, and Stefano Curtarolo. The AFLOW standard for high-throughput materials science calculations. *Comput. Mater. Sci.*, 108:233–238, October 2015. ISSN 0927-0256. doi: 10.1016/j.commatsci.2015.07.019.
- [17] Scott Kirklin, James E Saal, Bryce Meredig, Alex Thompson, Jeff W Doak, Muratahan Aykol, Stephan Rühl, and Chris Wolverton. The open quantum materials database (OQMD): assessing the accuracy of DFT formation energies. *npj Computational Materials*, 1(1):1–15, December 2015. ISSN 2057-3960, 2057-3960. doi: 10.1038/npjcompumats.2015.10.

- [18] Kamal Choudhary, Kevin F Garrity, Andrew C E Reid, Brian DeCost, Adam J Biacchi, Angela R Hight Walker, Zachary Trautt, Jason Hattnick-Simpers, A Gilad Kusne, Andrea Centrone, Albert Davydov, Jie Jiang, Ruth Pachter, Gowoon Cheon, Evan Reed, Ankit Agrawal, Xiaofeng Qian, Vinit Sharma, Houlong Zhuang, Sergei V Kalinin, Bobby G Sumpter, Ghanshyam Pilania, Pinar Acar, Subhasish Mandal, Kristjan Haule, David Vanderbilt, Karin Rabe, and Francesca Tavazza. The joint automated repository for various integrated simulations (JARVIS) for data-driven materials design. *npj Computational Materials*, 6(1):1–13, November 2020. ISSN 2057-3960, 2057-3960. doi: 10.1038/s41524-020-00440-1.
- [19] Sebastiaan P Huber, Spyros Zoupanos, Martin Uhrin, Leopold Talirz, Leonid Kahle, Rico Häuselmann, Dominik Gresch, Tiziano Müller, Aliaksandr V Yakutovich, Casper W Andersen, Francisco F Ramirez, Carl S Adorf, Fernando Gargiulo, Snehal Kumbhar, Elsa Passaro, Conrad Johnston, Andrius Merkys, Andrea Cepellotti, Nicolas Mounet, Nicola Marzari, Boris Kozinsky, and Giovanni Pizzi. AiiDA 1.0, a scalable computational infrastructure for automated reproducible workflows and data provenance. *Sci Data*, 7(1):300, September 2020. ISSN 2052-4463. doi: 10.1038/s41597-020-00638-4.
- [20] Exabyte.io, 2015. URL <https://exabyte.io/>. Date accessed: 2021-9-24.
- [21] Exabyte source of schemas and examples, 2021. URL <https://github.com/Exabyte-io/esse>. Date accessed: 2021-9-16.
- [22] JSON schema, 2017. URL <https://json-schema.org/>. Date accessed: 2021-9-16.
- [23] Clyde Edmiston and Klaus Ruedenberg. Localized atomic and molecular orbitals. *Rev. Mod. Phys.*, 35(3):457–464, July 1963. ISSN 0034-6861. doi: 10.1103/RevModPhys.35.457.
- [24] A D MacKerell, D Bashford, M Bellott, R L Dunbrack, J D Evanseck, M J Field, S Fischer, J Gao, H Guo, S Ha, D Joseph-McCarthy, L Kuchnir, K Kuczera, F T Lau, C Mattos, S Michnick, T Ngo, D T Nguyen, B Prodhom, W E Reiher, B Roux, M Schlenkrich, J C Smith, R Stote, J Straub, M Watanabe, J Wiórkiewicz-Kuczera, D Yin, and M Karplus. All-atom empirical potential for molecular modeling and dynamics studies of proteins. *J. Phys. Chem. B*, 102(18):3586–3616, April 1998. ISSN 1520-6106, 1520-5207. doi: 10.1021/jp973084f.
- [25] Timur Bazhurov. Data-centric online ecosystem for digital materials science. February 2019. URL <https://arxiv.org/abs/1902.10838>.
- [26] W M C Foulkes, L Mitas, R J Needs, and G Rajagopal. Quantum monte carlo simulations of solids. *Rev. Mod. Phys.*, 73(1):33–83, January 2001. ISSN 0034-6861. doi: 10.1103/RevModPhys.73.33.
- [27] Brian M Austin, Dmitry Yu Zubarev, and William A Lester, Jr. Quantum monte carlo and related approaches. *Chem. Rev.*, 112(1):263–288, January 2012. ISSN 0009-2665, 1520-6890. doi: 10.1021/cr2001564.
- [28] P Hohenberg and W Kohn. Inhomogeneous electron gas. *Phys. Rev.*, 136(3B):B864–B871, November 1964. ISSN 0959-8472. doi: 10.1103/PhysRev.136.B864.
- [29] W Kohn and L J Sham. Self-Consistent equations including exchange and correlation effects. *Phys. Rev.*, 140(4A):A1133–A1138, November 1965. ISSN 0959-8472. doi: 10.1103/PhysRev.140.A1133.
- [30] Dorothea Golze, Marc Dvorak, and Patrick Rinke. The GW compendium: A practical guide to theoretical photoemission spectroscopy. *Front Chem*, 7:377, July 2019. ISSN 2296-2646. doi: 10.3389/fchem.2019.00377.
- [31] Adam Rettig, Diptarka Hait, Luke W Bertels, and Martin Head-Gordon. Third-Order Møller-Plesset theory made more useful? the role of density functional theory orbitals. *J. Chem. Theory Comput.*, 16(12):7473–7489, December 2020. ISSN 1549-9618, 1549-9626. doi: 10.1021/acs.jctc.0c00986.
- [32] J P Perdew, K Burke, and M Ernzerhof. Generalized gradient approximation made simple. *Phys. Rev. Lett.*, 77(18):3865–3868, October 1996. ISSN 0031-9007, 1079-7114. doi: 10.1103/PhysRevLett.77.3865.

- [33] Aliaksandr V Krukau, Oleg A Vydrov, Artur F Izmaylov, and Gustavo E Scuseria. Influence of the exchange screening parameter on the performance of screened hybrid functionals. *J. Chem. Phys.*, 125(22):224106, December 2006. ISSN 0021-9606. doi: 10.1063/1.2404663.
- [34] Tomasz A Wesolowski and Yan Alexander Wang. *Recent Progress in Orbital-free Density Functional Theory*. World Scientific, 2013. ISBN 9789814436731.
- [35] William C Witt, Beatriz G del Rio, Johannes M Dieterich, and Emily A Carter. Orbital-free density functional theory for materials research. *J. Mater. Res.*, 33(7):777–795, April 2018. ISSN 0884-2914, 2044-5326. doi: 10.1557/jmr.2017.462.
- [36] Haoyu S Yu, Wenjing Zhang, Pragya Verma, Xiao He, and Donald G Truhlar. Nonseparable exchange-correlation functional for molecules, including homogeneous catalysis involving transition metals. *Phys. Chem. Chem. Phys.*, 17(18):12146–12160, May 2015. ISSN 1463-9076, 1463-9084. doi: 10.1039/c5cp01425e.
- [37] L H Thomas. The calculation of atomic fields. *Math. Proc. Cambridge Philos. Soc.*, 23(5):542–548, January 1927. ISSN 1469-8064, 0305-0041. doi: 10.1017/S0305004100011683.
- [38] E Fermi. Eine statistische methode zur bestimmung einiger eigenschaften des atoms und ihre anwendung auf die theorie des periodischen systems der elemente. *Zeitschrift für Physik*, 48(1):73–79, January 1928. ISSN 0044-3328. doi: 10.1007/BF01351576.
- [39] Oleg A Vydrov and Troy Van Voorhis. Nonlocal van der waals density functional: the simpler the better. *J. Chem. Phys.*, 133(24):244103, December 2010. ISSN 0021-9606, 1089-7690. doi: 10.1063/1.3521275.
- [40] Paolo Giannozzi, Oscar Basergio, Pietro Bonfà, Davide Brunato, Roberto Car, Ivan Carnimeo, Carlo Cavazzoni, Stefano de Gironcoli, Pietro Delugas, Fabrizio Ferrari Ruffino, Andrea Ferretti, Nicola Marzari, Iurii Timrov, Andrea Urru, and Stefano Baroni. Quantum ESPRESSO toward the exascale. *J. Chem. Phys.*, 152(15):154105, April 2020. ISSN 0021-9606, 1089-7690. doi: 10.1063/5.0005082.
- [41] Hans Martin Senn and Walter Thiel. QM/MM methods for biomolecular systems. *Angew. Chem. Int. Ed Engl.*, 48(7):1198–1229, 2009. ISSN 1433-7851, 1521-3773. doi: 10.1002/anie.200802019.
- [42] M van Schilfgaarde, Takao Kotani, and S Faleev. Quasiparticle self-consistent GW theory. *Phys. Rev. Lett.*, 96(22):226402, June 2006. ISSN 0031-9007. doi: 10.1103/PhysRevLett.96.226402.
- [43] Xavier Blase, Ivan Duchemin, Denis Jacquemin, and Pierre-François Loos. The Bethe-Salpeter equation formalism: From physics to chemistry. *J. Phys. Chem. Lett.*, 11(17):7371–7382, September 2020. ISSN 1948-7185. doi: 10.1021/acs.jpcclett.0c01875.
- [44] Jack Deslippe, Georgy Samsonidze, David A Strubbe, Manish Jain, Marvin L Cohen, and Steven G Louie. BerkeleyGW: A massively parallel computer package for the calculation of the quasiparticle and optical properties of materials and nanostructures. *Comput. Phys. Commun.*, 183(6):1269–1289, June 2012. ISSN 0010-4655. doi: 10.1016/j.cpc.2011.12.006.
- [45] Anthony J G Hey, Stewart Tansley, Kristin Michele Tolle, and Others. *The fourth paradigm: data-intensive scientific discovery*, volume 1. Microsoft research Redmond, WA, 2009.
- [46] Claudia Draxl and Matthias Scheffler. Big Data-Driven materials science and its FAIR data infrastructure. *Handbook of Materials Modeling*, pages 1–25, 2019. doi: 10.1007/978-3-319-42913-7\_104-1.
- [47] Mark D Wilkinson, Michel Dumontier, IJsbrand Jan Aalbersberg, Gabrielle Appleton, Myles Axton, Arie Baak, Niklas Blomberg, Jan-Willem Boiten, Luiz Bonino da Silva Santos, Philip E Bourne, Jildau Bouwman, Anthony J Brookes, Tim Clark, Mercè Crosas, Ingrid Dillo, Olivier Dumon, Scott Edmunds, Chris T Evelo, Richard Finkers, Alejandra Gonzalez-Beltran, Alasdair J G Gray, Paul Groth, Carole Goble, Jeffrey S Grethe, Jaap Heringa, Peter A C ’t Hoen, Rob Hooft, Tobias Kuhn, Ruben Kok, Joost Kok, Scott J Lusher, Maryann E Martone, Albert Mons, Abel L Packer, Bengt Persson, Philippe Rocca-Serra, Marco Roos, Rene van Schaik, Susanna-Assunta Sansone, Erik Schultes, Thierry Sengstag, Ted Slater, George Strawn, Morris A Swertz, Mark Thompson, Johan van der Lei, Erik van Mulligen, Jan Velterop, Andra Waagmeester, Peter Wittenburg, Katherine Wolstencroft, Jun Zhao, and Barend Mons. The FAIR guiding principles for scientific data management and stewardship. *Sci Data*, 3:160018, March 2016. ISSN 2052-4463. doi: 10.1038/sdata.2016.18.



# Appendix

## A Schema Examples

**Listing A.1:** Unit model example for Kohn-Sham density functional theory using the HSE06 exchange-correlation functional.

```
1 {
2   "tier1": {
3     "name": "physics-based",
4     "slug": "pb"
5   },
6   "tier2": {
7     "name": "quantum-mechanical",
8     "slug": "qm"
9   },
10  "tier3": {
11    "name": "density functional theory",
12    "slug": "dft"
13  },
14  "type": {
15    "name": "Kohn-Sham DFT",
16    "slug": "ksdft"
17  },
18  "subtype": {
19    "name": "Hybrid functional",
20    "slug": "hybrid"
21  },
22  "functional": {
23    "name": "HSE06",
24    "slug": "hse06",
25    "components": [
26      {
27        "name": "PBE",
28        "slug": "pbe-x",
29        "type": "exchange",
30        "fraction": 0.75,
31        "range": "short-range"
32      },
33      {
34        "name": "Exact exchange",
35        "slug": "hf",
36        "type": "exchange",
37        "fraction": 0.25,
38        "range": "short-range"
39      },
40      {
41        "name": "PBE",
42        "slug": "pbe-x",
43        "type": "exchange",
44        "fraction": 1.0,
45        "range": "long-range"
46      },
47      {
48        "name": "PBE",
49        "slug": "pbe-c",
50        "type": "correlation",
51        "fraction": 1.0
52      }
53    ]
54  },
55  "tags": [
56    "user-adjustable",
57    "range-separated"
58  ],
59  "method": { ... },
60  "reference": { ... },
61  "flowchartId": "u123fndff444",
62  "_id": "ffffff666666"
63 }
```

**Listing A.2:** Unit model example regarding coupled cluster theory including singly and doubly excited determinants (CCSD).

```
1 {
2   "tier1": {
3     "name": "physics-based",
4     "slug": "pb"
5   },
6   "tier2": {
7     "name": "quantum-mechanical",
8     "slug": "qm"
9   },
10  "tier3": {
11    "name": "ab-initio",
12    "slug": "abin"
13  },
14  "type": {
15    "name": "coupled cluster",
16    "slug": "cc"
17  },
18  "subtype": {
19    "name": "CC Singles Doubles",
20    "slug": "ccsd"
21  },
22  "tags": [
23    "single-reference",
24    "ground-state"
25  ],
26  "method": { ... },
27  "reference": { ... },
28  "flowchartId": "u123fndff444",
29  "_id": "ffffff666666"
30 }
```

**Listing A.3:** Unit model example for the ordinary least-squares regression model from the statistical-based model branch.

```
1 {
2   "tier1": {
3     "name": "statistical",
4     "slug": "st"
5   },
6   "tier2": {
7     "name": "deterministic",
8     "slug": "det"
9   },
10  "tier3": {
11    "name": "linear models",
12    "slug": "lin"
13  },
14  "type": {
15    "name": "ordinary least-squares",
16    "slug": "ols"
17  },
18  "tags": [
19    "supervised",
20    "regression"
21  ],
22  "method": { ... },
23  "reference": { ... },
24  "flowchartId": "u123fndff444",
25  "_id": "ffffff666666"
26 }
```