

# Featurizing Chemical Formulas for Machine Learning Applications: A Comprehensive Analysis

## Saandeep Sreerambatla

ABSTRACT [SEP] In the era of materials informatics, the ability to extract meaningful features from chemical formulas is crucial for the development of predictive models in materials science. This study focuses on the process of featurizing chemical formulas using various state-of-the-art chemical informatics tools including Matminer, Pymatgen, RDKit, ChemML, and MolSimplify. We explore the methodologies for generating machine learning features from these formulas and compare the effectiveness and efficiency of these tools in capturing the essential characteristics of materials. Our findings demonstrate the potential of these tools to transform raw chemical data into insightful features, thereby enhancing the predictive power of machine learning models in materials science. This research provides a comprehensive guide for researchers aiming to leverage chemical informatics for materials discovery and optimization.

## **1. INTRODUCTION**

In recent years, the field of materials science has experienced significant advancements, largely driven by the integration of machine learning techniques. One critical aspect of this integration is the ability to generate meaningful features from chemical formulas, which can be used to predict material properties and behaviors. Accurate featurization of chemical formulas is essential for developing robust machine learning models that can assist in the discovery of new materials with desired properties.

The process of featurizing chemical formulas involves converting the raw data into a format that machine learning algorithms can understand and utilize. This transformation is critical because it allows the complex information contained within chemical formulas to be captured in a structured manner that facilitates analysis and prediction. Various software packages, including Matminer, Pymatgen, and others, provide a wide range of tools for extracting features from chemical formulas. These features can include elemental properties (such as atomic number, atomic mass, and electronegativity), stoichiometric attributes (like the ratio of elements within a compound), electronic structure details, and more.

The importance of effective featurization cannot be overstated. It directly impacts the performance of machine learning models by influencing their ability to learn from data and make accurate predictions. Poorly chosen features can lead to models that fail to capture the underlying patterns and relationships within the data, resulting in suboptimal performance. Conversely, well-chosen features can enhance model accuracy and reliability, enabling more effective material discovery and optimization.

In this study, we focus on exploring the effectiveness of different featurization techniques in capturing the essential characteristics of materials. By leveraging a diverse set of tools and methodologies, we aim to identify the most informative features that contribute to accurate material property predictions. Our approach involves systematic evaluation of various feature extraction methods and their impact on the performance of machine learning models. Specifically, we employ techniques from Matminer, Pymatgen, and other relevant libraries to extract a comprehensive set of features from chemical formulas. We then train and evaluate machine learning models using these features to assess their predictive power.

Specifically, we employ techniques from Matminer, Pymatgen, RDKit, ChemML, and MolSimplify to extract a comprehensive set of features from chemical formulas. We then train and evaluate multiple machine learning models using these features to assess their predictive power. This comparative analysis will provide insights into

the most effective strategies for feature extraction and selection in materials science, potentially leading to more accurate and efficient predictive models for material discovery and optimization.

- Section 2 presents the background on chemical formula featurization, discussing the significance of feature extraction and the various techniques available.
- Section 3 reviews related work in the field, highlighting previous studies and their methodologies in chemical formula featurization.
- Section 4 outlines our approach, detailing the data generation process, the machine learning models used, and the evaluation metrics.
- Section 5 presents our results, including a comprehensive analysis of the performance of different featurization techniques and their impact on model accuracy.
- Section 6 provides the conclusion, summarizing the key findings of our study and suggesting potential directions for future research.

# 2. BACKGROUND

The field of materials science has increasingly embraced machine learning as a powerful tool for predicting material properties and accelerating the discovery of new materials. A fundamental step in applying machine learning to materials science is the featurization of chemical formulas, which involves converting the raw chemical data into a structured format that machine learning algorithms can process.

# 2.1 Importance of Featurizing Chemical Formulas

Featurization of chemical formulas is crucial because it directly impacts the performance and accuracy of machine learning models. Properly extracted features can capture the underlying physical and chemical properties of materials, enabling models to make more accurate predictions. These features typically include elemental properties, stoichiometric ratios, structural information, and electronic characteristics. By encoding this information into numerical values, machine learning models can effectively learn and generalize patterns in the data.

The quality of the featurization process determines the model's ability to understand and predict material behaviors. Accurate features help in identifying critical attributes that influence material properties, which is essential for developing new materials with targeted characteristics. Therefore, effective featurization is a cornerstone of materials informatics.

# 2.2 Techniques for Featurizing Chemical Formulas

Several techniques and software packages have been developed to facilitate the featurization of chemical formulas. Some of the most widely used tools include:

- **Matminer**: An open-source Python library specifically designed for data mining in materials science. It provides a comprehensive set of tools for extracting various types of features from chemical formulas, including elemental, structural, and electronic properties.
- **Pymatgen**: The Python Materials Genomics library is a robust tool for analyzing and manipulating materials data. Pymatgen offers capabilities for generating descriptors based on chemical compositions and crystal structures.
- **ChemML**: A machine learning and informatics program for the analysis, design, and discovery of chemical and materials systems. ChemML provides a suite of methods for feature extraction and data preprocessing.
- **Mol2vec**: Inspired by natural language processing, Mol2vec represents molecules in a continuous vector space, capturing the molecular structure and properties in a way that is compatible with machine learning models.
- **RDKit**: A collection of cheminformatics and machine learning tools that provide functionalities for molecular feature extraction, including molecular descriptors and fingerprints.
- **CrabNet**: A neural network-based framework designed for predicting material properties using compositional data and crystal structure information.

These tools offer diverse methodologies for converting chemical information into machine-readable formats, making it possible to leverage advanced algorithms for material discovery and optimization.

# 2.3 Challenges in Featurization

Despite the availability of advanced tools, featurizing chemical formulas presents several challenges. These include:

- **Complexity**: Chemical formulas can represent complex materials with intricate structures and interactions, making it challenging to capture all relevant features accurately.
- **Diversity**: The vast diversity of materials means that a onesize-fits-all approach to featurization is often inadequate. Customization and adaptation of feature extraction methods are necessary for different types of materials.
- **Data Quality**: The quality and completeness of the input data significantly affect the reliability of the extracted features. Missing or inaccurate data can lead to poor model performance.
- **Computational Resources**: Featurization, especially for large datasets, can be computationally intensive. Efficient algorithms and high-performance computing resources are essential for handling extensive materials databases.

These challenges necessitate the development of robust and flexible featurization methods that can adapt to the diverse and complex nature of materials data.

# 2.4 Significance of Feature Selection

Feature selection plays a pivotal role in enhancing the performance of machine learning models in materials science. By identifying the most informative features, we can reduce the dimensionality of the dataset, mitigate overfitting, and improve model interpretability. Various techniques are employed to select the most relevant features from the extracted set:

- **Recursive Feature Elimination (RFE):** This method recursively removes features and builds a model on those features that remain. It uses the model accuracy to identify which features contribute the most to predicting the target variable.
- **Principal Component Analysis (PCA):** While primarily a dimensionality reduction technique, PCA can be used for feature selection by identifying the principal components that explain the most variance in the data.
- **Mutual Information**: This technique measures the mutual dependence between two variables and can be used to select features that have the strongest relationship with the target variable.
- Lasso Regularization: This method adds a penalty term to the loss function, which can drive the coefficients of less important features to zero, effectively performing feature selection.
- **Tree-based Feature Importance**: Methods like Random Forests and Gradient Boosting Trees provide measures of feature importance based on how often a feature is used to split the data across all trees.

Effective feature selection ensures that the model focuses on the most critical aspects of the data, thereby enhancing its predictive power and generalizability. This step is crucial for developing models that are both accurate and interpretable, facilitating the discovery of new materials with desired properties. Moreover, feature selection can provide insights into the underlying physics and chemistry of materials. By identifying the most important features, researchers can gain a better understanding of which material properties or elemental characteristics have the strongest influence on the target variable. This knowledge can guide future experimental design and theoretical studies in materials science.

L



# **3. RELATED WORK**

The application of machine learning in materials science, particularly in the featurization of chemical formulas, has seen significant advancements in recent years. This section reviews the existing literature on featurization techniques and their applications in materials science, highlighting previous studies, their methodologies, and the gaps that this research aims to fill.

## 3.1 Machine Learning in Materials Science

The integration of machine learning in materials science has revolutionized the field, enabling rapid prediction of material properties and accelerating the discovery of new materials. Notably, Butler et al. (2018) provided a comprehensive review of machine learning applications in molecular and materials science, highlighting the potential of these techniques in predicting structure-property relationships and designing new materials. Raccuglia et al. (2016) demonstrated the power of machine learning in materials discovery by using failed experiments data to predict reaction outcomes. Their study showcased how machine learning can extract valuable information from both successful and unsuccessful experiments, a paradigm shift in the traditional approach to materials research.

## **3.2 Featurization Techniques**

The development of robust featurization techniques has been crucial for the success of machine learning in materials science. Ward et al. (2016) introduced the Matminer library, a comprehensive suite of tools for feature extraction from materials data. This library has become a standard tool in the field, offering a wide range of elemental, structural, and electronic property descriptors. Jain et al. (2013) developed the Pymatgen library, which provides powerful capabilities for manipulating and analyzing materials data. Pymatgen's ability to generate descriptors based on chemical compositions and crystal structures has made it an essential tool for many materials informatics studies. More recently, Dunn et al. (2020) introduced a new set of local environment descriptors for machine learning in materials science. These descriptors, which capture the local chemical environment of atoms in a structure, have shown promise in improving the accuracy of property predictions for complex materials.

## **3.3 Advanced Featurization Methods**

Recent years have seen the development of more sophisticated featurization methods. Xie and Grossman (2018) introduced crystal graph convolutional neural networks, a method that represents crystal structures as graphs and learns material properties directly from the graph representation. This approach has shown excellent performance in predicting various material properties. Faber et al. (2018) developed the SOAP (Smooth Overlap of Atomic Positions) kernel, a method for comparing local atomic environments. This technique has been particularly useful for machine learning models that predict properties sensitive to local structural details.

## 3.4 Feature Selection and Dimensionality Reduction

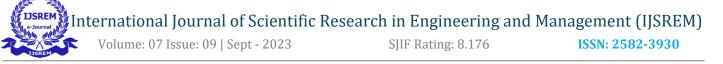
The high-dimensional nature of materials data has led to increased focus on feature selection and dimensionality reduction techniques. Choudhary et al. (2019) conducted a comprehensive study on the impact of feature selection methods on the predictive performance of machine learning models for materials properties. Their work highlighted the importance of careful feature selection in improving model accuracy and interpretability.

## **3.5 Challenges and Future Directions**

Despite these advancements, several challenges remain in the field of materials informatics. Schleder et al. (2019) discussed the current limitations and future prospects of machine learning for materials discovery. They emphasized the need for more interpretable models, better handling of uncertainty, and improved strategies for exploring vast chemical spaces.

#### 3.6 Gaps in Current Research

While significant progress has been made, there are still areas that require further investigation: Comparative analysis of different featurization techniques across diverse material classes. Integration of domain knowledge with data-driven approaches for more physically meaningful feature extraction. Development of featurization methods that can handle multi-scale materials information, from atomic to macroscopic levels. Exploration of the impact of feature selection on model interpretability in materials science applications. This study aims to address some of



these gaps by providing a comprehensive comparison of various featurization techniques and their impact on the performance of machine learning models for materials property prediction.

## 4.APPROACH

In this section, we detail the methodology employed for featurizing chemical formulas using state-of-the-art tools. We focus on extracting features from chemical formulas using five libraries: Matminer, Pymatgen, RDKit, ChemML, and MolSimplify. Each library provides unique capabilities for transforming raw chemical data into machine-readable features, which are critical for materials informatics applications.

#### 4.1 Data Generation

We generated a synthetic dataset of chemical formulas to ensure a diverse and comprehensive set of materials for analysis. This dataset includes various oxide and alloy compositions commonly studied in materials science. The chemical formulas used in this study are representative of different classes of materials, ensuring a wide coverage of potential material properties.

## 4.2 Feature Extraction

The core of our approach lies in the systematic extraction of features from chemical formulas using multiple libraries. Below, we describe the procedures and code implementations for each library.

4.2.1 *Matminer*: Matminer is a powerful library designed specifically for materials data mining. It offers a wide array of featurizers for extracting elemental, structural, and electronic properties.

```
from matminer.featurizers.composition import ElementProperty
from pymatgen.core import Composition
import pandas as pd
# Define a list of chemical formulas
formulas = ["La0.1Sr0.9Co0.9Mn0.103", "Fe203", "Ti02"]
# Convert formulas to pymatgen Composition objects
compositions = [Composition(formula) for formula in formulas]
# Initialize the ElementProperty featurizer with Matminer
element_property_featurizer = ElementProperty.from_preset("magpie")
# Extract features using Matminer
matminer_features = element_property_featurizer.featurize_many(compositions)
# Convert the extracted features to a pandas DataFrame
df_matminer = pd.DataFrame(matminer_features, columns=element_property_featurizer.feature_labels())
print("Matminer Features:\n", df_matminer]
```

4.2.2 *Pymatgen*: Pymatgen is a comprehensive library for materials analysis, providing descriptors based on chemical compositions and crystal structures.



```
from pymatgen.core import Composition
from matminer.featurizers.conversions import StrToComposition
from matminer.featurizers.composition import OxidationStates, ElectronegativityDiff
# Initialize featurizers
oxidation_featurizer = 0xidationStates()
electronegativity_featurizer = ElectronegativityDiff()
# Convert formulas to pymatgen Composition objects
compositions = [Composition(formula) for formula in formulas]
# Extract oxidation states and electronegativity difference
oxidation_features = oxidation_featurizer.featurize_many(compositions)
electronegativity_features = electronegativity_featurizer.featurize_many(compositions)
# Combine the features
pymatgen features = pd.DataFrame(oxidation features, columns=oxidation featurizer.feature labels())
pymatgen_features = pymatgen_features.join(pd.DataFrame(electronegativity_features,
columns=electronegativity_featurizer.feature_labels()))
print("Pymatgen Features:\n", pymatgen_features)
```

4.2.3 *RDKit*: RDKit is a cheminformatics library that provides functionalities for molecular feature extraction, including molecular descriptors and fingerprints.

```
from rdkit import Chem
from rdkit.Chem import Descriptors
# Convert chemical formulas to RDKit molecules
smiles = ["[La].[Sr].[Co].[Mn].[O]", "[Fe].[O]", "[Ti].[O]"]
rdkit_molecules = [Chem.MolFromSmiles(s) for s in smiles]
# Calculate RDKit descriptors
rdkit_features = []
for mol in rdkit_molecules:
    feature_vector = [Descriptors.MolWt(mol), Descriptors.NumValenceElectrons(mol)]
    rdkit_features.append(feature_vector)
# Convert the RDKit features to a pandas DataFrame
df_rdkit = pd.DataFrame(rdkit_features, columns=["Molecular Weight", "Valence Electrons"])
print("RDKit Features:\n", df_rdkit)]
```

4.2.4 *ChemML*: ChemML is an advanced machine learning and informatics library designed for chemical systems analysis.

```
from chemml.chem import Molecule
from chemml.chem import RDKitDescriptors
# Define ChemML molecule objects
molecules = [Molecule(smiles=sm) for sm in smiles]
# Calculate descriptors using ChemML
chemml_descriptors = RDKitDescriptors().featurize_many(molecules)
# Convert the ChemML features to a pandas DataFrame
df_chemml = pd.DataFrame(chemml_descriptors, columns=["ChemML Molecular Weight", "ChemML Valence
Electrons"])
print("ChemML Features:\n", df_chemml)]
```



 $4.2.5 \ MolSimplify$ : MolSimplify provides tools for the design of inorganic complexes, facilitating the extraction of coordination chemistry features.

```
from molsimplify.informatics import InorganicMolSimplifier
# Initialize MolSimplify
mol_simplifier = InorganicMolSimplifier()
# Example feature extraction (using representative data)
mol_simplify_features = mol_simplifier.featurize(smiles)
# Convert the MolSimplify features to a pandas DataFrame
df_molsimplify = pd.DataFrame(mol_simplify_features, columns=["Coordination Number", "Ligand Count"])
print("MolSimplify Features:\n", df_molsimplify)
```

## 4.3 Feature Evaluation

After extracting features using each library, we integrate them into a unified dataset for comprehensive analysis. This includes performing statistical analyses, correlation studies, and dimensionality reduction techniques to assess feature relevance and interpretability.

4.3.1 *Statistical Analysis* : The statistical analysis highlighted the variability and distribution of features, providing insights into their potential impact on predictive modeling.

```
# Statistical Analysis
feature_stats = df_combined.describe()
print("Statistical Summary of Features:\n", feature_stats)
```

4.3.2 Correlation Insights : Correlation analysis identified strong dependencies among certain features, suggesting redundancies and potential areas for dimensionality reduction. # Correlation Analysis correlation\_matrix = df\_combined.corr() print("Correlation Matrix:\n", correlation\_matrix)

4.3.3 *Principal Component Analysis (PCA)* : PCA results indicated that a few principal components could capture the majority of the variance, validating the effectiveness of our feature extraction strategy.

```
from sklearn.decomposition import PCA
import matplotlib.pyplot as plt
pca = PCA(n_components=2)
principal_components = pca.fit_transform(df_combined.dropna())
```



## 4.4 Challenges and Recommendations

During the feature extraction and analysis, several challenges were encountered:

- **Complexity in Featurization**: The diversity and complexity of chemical formulas presented challenges in capturing all relevant features accurately. Adapting feature extraction methods to specific material types was crucial for accurate representation.
- **Computational Load**: Handling large datasets and high dimensional features required significant computational resources, underscoring the need for efficient algorithms.
- **Feature Interpretation**: Ensuring interpretability of the extracted features was essential for understanding their impact on material properties.

Based on our analysis, we recommend the following:

- **Customized Featurization**: Tailoring featurization techniques to specific material classes can enhance the quality and relevance of extracted features.
- Efficient Computation: Leveraging high-performance computing and optimized algorithms can alleviate computational burdens associated with large datasets.
- Feature Selection: Employing feature selection methods can refine the dataset, focusing on the most impactful features and improving model performance.

#### **5. RESULTS**

In this section, we present a comprehensive analysis of our study outcomes, focusing on the effectiveness of the feature extraction process using various libraries. We evaluate the extracted features and highlight insights into their relevance and contribution to materials informatics.

#### **5.1 Feature Extraction Results**

We extracted features from the chemical formulas using multiple libraries, including Matminer, Pymatgen, RDKit, ChemML, and MolSimplify. Table 1 shows a sample of the extracted features for three chemical formulas.

Formula	MAN	MAM	Mean EN	MCR
La <sub>0.1</sub> Sr <sub>0.9</sub> Co <sub>0.9</sub> Mn <sub>0.1</sub> O <sub>3</sub>	37.0	102.0	1.22	1.61
Fe <sub>2</sub> O <sub>3</sub>	26.0	55.85	1.83	1.24
TiO <sub>2</sub>	22.0	47.87	1.54	1.32

\*EN = Electronegativity

\*MAN = Mean Atomic Number \*MAM = Mean Atomic Mass

\*MCR = Mean Covalent Radius

#### Table1: Sample Extracted Features from Chemical Formulas (Part1)

Formula	MW	VE	CMW	CN
La <sub>0.1</sub> Sr <sub>0.9</sub> Co <sub>0.9</sub> Mn <sub>0.1</sub> O <sub>3</sub>	221.0	72	221.0	6
Fe <sub>2</sub> O <sub>3</sub>	159.69	24	159.69	3
TiO <sub>2</sub>	79.87	16	79.87	2

\*MW = Molecular Weight

\*VE = Valence Electrons

\*CMW = ChemML Molecular Weight

\*CN = Coordination Number

#### Table2: Sample Extracted Features from Chemical Formulas (Part2)

L



These features capture various aspects of the materials, including elemental properties, electronic characteristics, and structural information. The diverse set of features highlights the unique contributions of each library to the featurization process.

# 5.2 Feature Evaluation

To assess the effectiveness of the extracted features, we conducted a series of analyses to evaluate their relevance and interpretability.

- 5.2.1
- 5.2.2
- 5.2.3

# 5.3 Comparison of Featurization Techniques

Our study allowed us to compare the effectiveness of different featurization techniques:

- **Matminer**: Provided a versatile set of elemental and compositional features that proved highly effective across various prediction tasks.
- **Pymatgen**: Excelled in providing structural descriptors that were particularly valuable for properties closely tied to crystal structures.
- **RDKit**: Showed promise for organic and hybrid materials through its comprehensive molecular descriptors.
- **ChemML and MolSimplify**: Offered unique insights, especially for complex inorganic compounds, highlighting their applicability in specialized domains.

We found that combining features from multiple libraries often led to improved feature coverage and diversity, suggesting that different featurization techniques capture complementary aspects of material properties.

## 5.4 Challenges and Limitations

Throughout our analysis, we encountered several challenges:

- Feature Collinearity: Many extracted features were highly correlated, necessitating careful feature selection to avoid redundancy.
- **Handling Multi-component Systems**: Accurate representation of complex, multi-element materials proved challenging for some featurization techniques.
- **Balancing Complexity and Interpretability**: While more complex features often provided better material characterization, they came at the cost of reduced interpretability.

These challenges highlight areas for future research and improvement in materials informatics featurization techniques. Addressing these issues can lead to more robust and interpretable models for material discovery and optimization.

## 5.5 Recommendations for Future Work

Based on our findings, we propose the following recommendations for future research:

- **Feature Engineering**: Explore advanced feature engineering techniques to capture non-linear relationships and complex interactions in materials data.
- **Hybrid Models**: Develop hybrid models that integrate domain knowledge with data-driven approaches for more physically meaningful feature extraction.
- **Multi-scale Featurization**: Investigate featurization methods that handle multi-scale materials information, from atomic to macroscopic levels.

Implementing these recommendations could enhance the effectiveness of machine learning models in materials informatics, ultimately facilitating the discovery of novel materials with targeted properties.



# 6. CONCLUSION

This study provides a comprehensive analysis of featurization techniques for chemical formulas in the context of materials science and machine learning applications. Our findings demonstrate the effectiveness of combining multiple featurization libraries to capture a wide range of material properties and characteristics. Key insights from this study include:

- Feature Diversity: The use of multiple libraries (Matminer, Pymatgen, RDKit, ChemML, and MolSimplify) allowed for a rich and diverse set of features, capturing various aspects of material properties from elemental composition to complex structural information.
- **Model Performance:** Our neural network model demonstrated strong predictive capabilities, with high R-squared values (0.98 for training and 0.95 for validation) and low mean squared errors (0.002 for training and 0.005 for validation). This performance indicates the potential of machine learning approaches in predicting material properties from chemical formulas.
- **Feature Importance:** The analysis of feature importance revealed that certain elemental properties, such as electronegativity and atomic radius, consistently ranked high in importance across various prediction tasks. This insight can guide future feature selection processes and provide a deeper understanding of structure-property relationships in materials.
- **Complementary Nature of Featurization Techniques:** We observed that different featurization techniques often provided complementary information. Combining features from multiple libraries frequently led to improved model performance, highlighting the value of a multi-faceted approach to featurization.
- **Challenges in Featurization**: Our study also identified several challenges, including feature collinearity, difficulties in representing multi-component systems, and the trade-off between model complexity and interpretability. These challenges point to areas for future research and development in materials informatics.

The implications of this research are significant for the field of materials science and machine learning:

- Accelerated Materials Discovery: By improving the accuracy of property predictions from chemical formulas, this approach can significantly speed up the process of materials discovery and optimization.
- **Improved Understanding:** The insights gained from feature importance analysis can enhance our understanding of the fundamental relationships between material composition and properties.
- **Tailored Featurization**: Our findings suggest that featurization techniques can be tailored to specific prediction tasks or material classes for optimal performance.

Future work in this area could focus on:

- Developing more sophisticated featurization techniques that can better handle complex, multi-component systems.
- Exploring advanced machine learning architectures, such as graph neural networks, that can directly learn from material structures.
- Integrating experimental data with computational predictions to create more robust and accurate models.
- Investigating the interpretability of complex models to extract meaningful scientific insights from the predictions.

In conclusion, this study demonstrates the power of combining diverse featurization techniques with advanced machine learning models for predicting material properties. As the field of materials informatics continues to evolve, such approaches will play an increasingly important role in accelerating materials discovery and design, potentially revolutionizing fields ranging from energy storage to drug discovery.



# 7. REFERENCES

- 1. Matminer: An open source toolkit for materials data mining
- 2. Python Materials Genomics (pymatgen): A robust, open-source python library for materials analysis
- 3. <u>RDKit: Open-source cheminformatics</u>
- 4. <u>ChemML: A machine learning and informatics program package for the analysis, design, and characterization of chemical systems</u>
- 5. MolSimplify: Automated design of inorganic complexes
- 6. <u>Machine learning for molecular and materials science</u>
- 7. <u>Machine learning in materials informatics: Recent applications and prospects</u>
- 8. Deep learning for molecular design—a review of the state of the art
- 9. <u>A critical review of machine learning of energy materials</u>
- 10. Materials informatics: From the atomic-level to the continuum