

This is a self-archived version of an original article. This version may differ from the original in pagination and typographic details.

Author(s): Chang, Jin Hyun; Kleiven, David; Melander, Marko; Akola, Jaakko; Garcia-Lastra, Juan Maria; Vegge, Tejs

Title: CLEASE : A versatile and user-friendly implementation of Cluster Expansion method

Year: 2019

Version: Published version

Copyright: © 2019 IOP Publishing Ltd.

Rights: CC BY 4.0

Rights url: <https://creativecommons.org/licenses/by/4.0/>

Please cite the original version:

Chang, J. H., Kleiven, D., Melander, M., Akola, J., Garcia-Lastra, J. M., & Vegge, T. (2019). CLEASE : A versatile and user-friendly implementation of Cluster Expansion method. *Journal of Physics: Condensed Matter*, 31(32), Article 325901. <https://doi.org/10.1088/1361-648X/ab1bbc>

PAPER • OPEN ACCESS

CLEAVE: a versatile and user-friendly implementation of cluster expansion method

To cite this article: Jin Hyun Chang *et al* 2019 *J. Phys.: Condens. Matter* **31** 325901

View the [article online](#) for updates and enhancements.



IOP | ebooks™

Bringing you innovative digital publishing with leading voices to create your essential collection of books in STEM research.

Start exploring the collection - download the first chapter of every title for free.

CLEAVE: a versatile and user-friendly implementation of cluster expansion method

Jin Hyun Chang¹, David Kleiven², Marko Melander³, Jaakko Akola^{2,4}, Juan Maria Garcia-Lastra¹ and Tejs Vegge¹

¹ Department of Energy Conversion and Storage, Technical University of Denmark, DK-2800 Kgs. Lyngby, Denmark

² Department of Physics, Norwegian University of Science and Technology, NO-7491 Trondheim, Norway

³ Department of Chemistry, University of Jyväskylä, Jyväskylä, Finland

⁴ Computational Physics Laboratory, Tampere University, PO Box 692, FI-33014 Tampere, Finland

E-mail: jchang@dtu.dk

Received 12 February 2019, revised 8 April 2019

Accepted for publication 23 April 2019

Published 28 May 2019



CrossMark

Abstract

Materials exhibiting a substitutional disorder such as multicomponent alloys and mixed metal oxides/oxyfluorides are of great importance in many scientific and technological sectors. Disordered materials constitute an overwhelmingly large configurational space, which makes it practically impossible to be explored manually using first-principles calculations such as density functional theory due to the high computational costs. Consequently, the use of methods such as cluster expansion (CE) is vital in enhancing our understanding of the disordered materials. CE dramatically reduces the computational cost by mapping the first-principles calculation results on to a Hamiltonian which is much faster to evaluate. In this work, we present our implementation of the CE method, which is integrated as a part of the atomic simulation environment (ASE) open-source package. The versatile and user-friendly code automates the complex set up and construction procedure of CE while giving the users the flexibility to tweak the settings and to import their own structures and previous calculation results. Recent advancements such as regularization techniques from machine learning are implemented in the developed code. The code allows the users to construct CE on any bulk lattice structure, which makes it useful for a wide range of applications involving complex materials. We demonstrate the capabilities of our implementation by analyzing the two example materials with varying complexities: a binary metal alloy and a disordered lithium chromium oxyfluoride.


Keywords: cluster expansion, Monte Carlo, disordered materials, battery material, alloys

(Some figures may appear in colour only in the online journal)

1. Introduction

Computational modeling of materials with a substitutional disorder such as multicomponent alloys and mixed metal oxides is said to have a *configurational* problem [1–5]. The

vast configurational space of these materials makes it practically impossible to explore directly using first-principles calculations such as density functional theory (DFT). A quantitative method capable of establishing the relationship between the structure and property of materials is therefore essential. Cluster Expansion (CE) [1, 5–9] is a method that has been used successfully in the past few decades to parameterize and express the configurational dependence of physical properties. The most widely parameterized physical property

 Original content from this work may be used under the terms of the [Creative Commons Attribution 3.0 licence](https://creativecommons.org/licenses/by/3.0/). Any further distribution of this work must maintain attribution to the author(s) and the title of the work, journal citation and DOI.

is energy computed using first-principles methods, but CE can also be used to parameterize other quantities such as band gap [10, 11] and density of states [12].

Despite its success and usefulness in predicting physical properties of crystalline materials, CE remains as a niche tool used in a small subfield within the computational materials science, primarily used by specialists. On the other hand, the research fields in which CE is becoming relevant is on the rise; one such example is the use of disordered materials for battery applications [13–17]. The objective of our work is to make cluster expansion more accessible for a broad range of computational scientists who do not necessarily possess expertise in cluster expansion. Our approach to achieving such a goal is to implement CE as a part of a widely used, open-source atomic simulation environment (ASE) package [18]. Henceforth, we refer to our implementation as CLEASE, which stands for cluster expansion in atomic simulation environment.

Having CE as a part of a widely used package with interfaces to a multitude of open-source and commercial atomic-scale simulation codes accompanies several practical benefits: (1) a large existing user base does not need to install or learn a new program as the CE module is a part of ASE and inherits its syntax and code style, and (2) all of the atomic-scale simulation codes supported by ASE are also automatically supported by the implemented module. In addition, CLEASE utilizes the database management feature implemented in ASE, which provides an efficient way to store, maintain and share both DFT and CE results. Therefore, the implementation presented in this article appeals to a significant portion of computational materials science community as a versatile and easy-to-learn package, thereby lowering the barrier to incorporate cluster expansion as a part of their research methods to accelerate computational materials prediction and design.

The rest of the paper is organized as follows. A brief overview of cluster expansion formalism and other important concepts are provided in section 2 in order to aid the readers who are not familiar with the cluster expansion method. The implementation of CLEASE is described in section 3. Section 4 contains two application examples with different levels of complexities, namely a binary metal alloy and a lithium metal oxyfluoride. The computational settings and technical details for the examples are provided in section 5.

2. Theory

2.1. Cluster expansion formalism

The core concept of the cluster expansion is to express the scalar physical quantity of a material, $q(\boldsymbol{\sigma})$, to its configuration, $\boldsymbol{\sigma}$, where a crystalline system is represented with a fixed underlying grid of atomic sites. In such a representation, any configuration with the same underlying topology can be completely specified by the atomic occupation of each atomic site. For the case of a crystalline material with N atomic sites, any configuration can be specified by an N -dimensional vector $\boldsymbol{\sigma} = \{s_1, s_2, \dots, s_N\}$, where s_i is a site variable that specifies which type of atom occupies the atomic site i (also referred to as an occupation variable [4, 19, 20] or pseudospin [1, 10,

21–23]). It is noted that the terms configuration and structure are often used interchangeably.

For the case of multinary systems consisting of M different atomic species, s_i takes one of M distinct values. The original formulation of Sanchez *et al* [6] specifies the s_i to take any values from $\pm m, \pm(m-1), \dots, \pm 1$ for $M = 2m$ (for the case where there is an odd number of element types, an additional value of 0 should be included in the possible values of s_i , and the relation between M and m becomes $M = 2m - 1$). Other choices of s_i are also commonly used such as values ranging from 0 to $M - 1$ by van de Walle [20] and from 1 to M by Mueller and Ceder [24]. Based on the original formalism by Sanchez *et al*, single-site basis functions are determined through an orthogonality condition

$$\frac{1}{M} \sum_{s_i=-m}^m \Theta_n(s_i) \Theta_{n'}(s_i) = \delta_{nn'}, \quad (1)$$

where $\Theta_n(s_i)$ is the n th single-site basis function (e.g. Chebyshev polynomials) for i th site and $\delta_{nn'}$ is a Kronecker delta.

The configuration is decomposed into a sum of clusters as shown in figure 1. Each cluster has a set of associated cluster functions, which are defined as

$$\Phi_{\mathbf{n}}(\mathbf{s}) = \prod_i \Theta_{n_i}(s_i), \quad (2)$$

where \mathbf{n} and \mathbf{s} are vectors specifying the order of the single-site basis function and the site variables in the cluster, respectively. n_i and s_i specify the i th element of the respective vectors. The use of orthogonal basis functions guarantees that the cluster functions defined in (2) are also orthogonal. The symmetrically equivalent clusters are classified as the same cluster, and the collection of all symmetrically equivalent clusters are denoted with an α .

The average value of the cluster functions in cluster α is referred to as a correlation function, ϕ_α . The physical quantity, $q(\boldsymbol{\sigma})$, normalized with the number of atomic sites N is then expressed as

$$q(\boldsymbol{\sigma}) = \sum_{\alpha} m_{\alpha} J_{\alpha} \phi_{\alpha}, \quad (3)$$

where m_{α} is the multiplicity factor indicating the number of cluster α per atom and J_{α} is the effective cluster interaction (ECI) per occurrence, which needs to be determined. It is noted that the cluster α includes the cluster of size zero, which have $m_{\alpha} \phi_{\alpha} = 1$. Alternatively, (3) can be written in a more explicitly form,

$$q(\boldsymbol{\sigma}) = J_0 + \sum_{\alpha} m_{\alpha} J_{\alpha} \phi_{\alpha}, \quad (4)$$

where J_0 is the ECI of an empty cluster while α in this case corresponds to the clusters of size one and higher. It is often more practical and convenient to express the ECI per atom rather than per occurrence [5], in which case m_{α} and J_{α} are combined into one term, $\tilde{J}_{\alpha} = m_{\alpha} J_{\alpha}$ and (3) becomes

$$q(\boldsymbol{\sigma}) = \sum_{\alpha} \tilde{J}_{\alpha} \phi_{\alpha}. \quad (5)$$

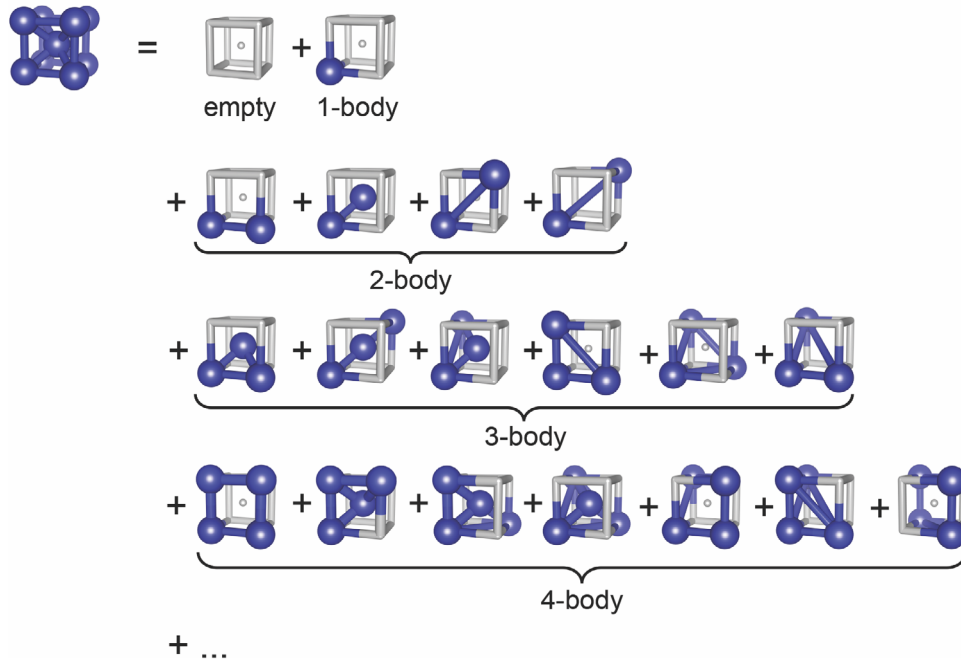


Figure 1. A simplified illustration of the decomposition of a body-centered cubic lattice.

CLEAVE uses the ECI per atom (\tilde{J}_α), but interested users can determine the value of J_α based on the values of m_α and \tilde{J}_α .

Theoretically, there is an infinite number of terms in (5) for an infinite crystal, and the resulting expression can represent any scalar function $q(\sigma)$ given that appropriate ECI values are found. In practice, sufficient accuracy is often reached with clusters with small number of atoms (e.g. one-, two- and three-body clusters) that are relatively compact in size (e.g. 5–7 Å in diameter).

2.2. Cluster selection & determination of ECI values

A crucial element of CE approach is to select relevant clusters from a theoretically infinite number of possible clusters. Many multicomponent systems yield thousands of clusters even when the expansion is limited to relatively compact size and small number of atoms, and they are vastly truncated since only a small fraction of them is needed to achieve the required accuracy. Determining the optimal set of clusters that minimizes the number of clusters without losing its predictive power has been a topic of keen interest in the past decade [3, 19, 25–27], and the cluster selection based on genetic algorithms [3, 25, 26] was considered to be the most robust method.

More recently, the use of compressive sensing [22] was proposed to efficiently select the clusters and determine their ECIs in one shot. The compressive sensing is based on ℓ_1 norm (a special case of ℓ_p norm where $p = 1$), which is defined as

$$\|\mathbf{x}\|_p = \left(\sum_i |x_i|^p \right)^{1/p}, \quad (6)$$

where \mathbf{x} is a vector quantity. It is noted that cluster expansion defined in (5) is in the same form as a linear regression model

in statistics and machine learning. Therefore, one can treat CE as a linear regression problem and apply regularization techniques based not only on ℓ_1 norm but also on any other p values, although ℓ_1 and ℓ_2 norms are most commonly used.

The use of regularization techniques for CE can be illustrated by expressing (5) in a matrix form,

$$\mathbf{q} = \mathbf{X}\boldsymbol{\omega}. \quad (7)$$

\mathbf{X} is a matrix containing the correlation functions of the training data where each element in row i and column α is defined as

$$\mathbf{X}_{i\alpha} = \phi_\alpha(\sigma_i). \quad (8)$$

\mathbf{q} is a column vector in which the i th element is the physical quantity q of the configuration σ_i and $\boldsymbol{\omega}$ is a column vector in which α th element is \tilde{J}_α .

The simplest way of determining $\boldsymbol{\omega}$ is by using ordinary least squares (OLS) method, which minimizes the residual sum of squared errors (RSS). RSS is defined as

$$\text{RSS} = \|\mathbf{X}\boldsymbol{\omega} - \mathbf{q}\|_2^2, \quad (9)$$

and the minimization of the RSS has a unique solution $\hat{\boldsymbol{\omega}}$ where

$$\begin{aligned} \hat{\boldsymbol{\omega}} &= \arg \min_{\boldsymbol{\omega}} \|\mathbf{X}\boldsymbol{\omega} - \mathbf{q}\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{q}. \end{aligned} \quad (10)$$

The OLS has two major drawbacks [22]. The first is the requirement on which the number of configurations in the training set needs to be larger than the number of clusters being considered. The matrix $\mathbf{X}^T \mathbf{X}$ becomes singular in such a case, and the limitations imposed by the first requirement become more severe for systems consisting of many element types since even strict expansion conditions (i.e. small number of atoms per cluster and compact size) can lead to a large

number of clusters. The second drawback is the susceptibility to possible overfitting, which refers to the conditions in which the ECI values are over-tuned to accurately represent $q(\boldsymbol{\sigma})$ of the training set at a cost of losing its predictive power for the new configurations that are not included in the training set. The overfitting also makes the model prone to noise present in the training data because the model attempts to meticulously fit the model to the training data including the noise therein.

Regularization is an efficient technique to address the aforementioned drawbacks of OLS by adding a regularization term to (10). The most common regularization scheme are ℓ_1 and ℓ_2 regularization, which respectively uses ℓ_1 and ℓ_2 norm as a regularization term. For ℓ_1 regularization, the solution becomes

$$\hat{\boldsymbol{\omega}} = \arg \min_{\boldsymbol{\omega}} \|\mathbf{X}\boldsymbol{\omega} - \mathbf{q}\|_2^2 + \lambda \|\boldsymbol{\omega}\|_1, \quad (11)$$

where λ is a regularization parameter that controls the weight given to the regularization term. The main benefit of ℓ_1 regularization is its promotion of sparsity. In context of CE, the sparsity means a selection of a fewer number of clusters, or many clusters with their ECI values set to zero. It is noted that there is no unique analytical solution for (11), and it needs to be solved iteratively. Unlike ℓ_1 regularization, ℓ_2 regularization has a unique analytical solution which is expressed as

$$\begin{aligned} \hat{\boldsymbol{\omega}} &= \arg \min_{\boldsymbol{\omega}} \|\mathbf{X}\boldsymbol{\omega} - \mathbf{q}\|_2^2 + \|\boldsymbol{\omega}\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{q}. \end{aligned} \quad (12)$$

However, ℓ_2 regularization does not promote sparsity, and the resulting solution is likely to contain more clusters than necessary. It is noted that Bayesian compressive sensing [21] scheme is introduced for cluster expansion, which effectively eliminates the parameter λ in ℓ_1 and ℓ_2 regularization schemes while promoting sparsity.

Regardless of the fitting technique used, the predictive power of the expansion needs to be assessed to determine its accuracy and reliability. Cross-validation (CV) is a technique used for assessing the prediction accuracy of the model. A leave-one-out (LOO) scheme is most commonly used in CE community, and the LOOCV score is defined as

$$\text{LOOCV} = \left(\frac{1}{N_{\text{config}}} \sum_{i=1}^{N_{\text{config}}} (\hat{q}_i - q_i)^2 \right)^{1/2}, \quad (13)$$

where N_{config} is the number of configurations in the training set, \hat{q}_i is the physical quantity of a structure i predicted by CE using $N_{\text{config}} - 1$ structures without a structure i and q_i is the calculated physical quantity of structure i . While OLS has only one (likely overfitted) solution, ℓ_1 and ℓ_2 regularization schemes have a solution for each λ value. The solution—a selection of clusters and their ECI values—that yields the lowest LOOCV score is chosen. Although LOO is the most common cross validation scheme in cluster expansion community, k -fold CV is one of the most common schemes used in machine learning community. In a k -fold CV scheme, the pool of configurations are randomly partitioned into k parts of equal size. The structures in $k - 1$ parts are used as training

data while the remaining one part is used as a validation set, and the cross validation is repeated k times.

2.3. Thermodynamics in lattice models

The true benefit of CE is in its ability to predict the expanded scalar quantity $q(\boldsymbol{\sigma})$ based on trained data. An accurate prediction can be made if the CV score of the expanded $q(\boldsymbol{\sigma})$ is sufficiently low, and the prediction speed is very fast on modern computer architecture since it only involves executions of only a small number of simple numerical calculations specified in (5). Such a speed boost allows one to conduct types of analyses that require substantial statistical sampling.

In contrast to zero temperature studies where the system occupies the state with lowest energy, an ensemble of configurations with the lowest free energy are occupied at finite temperature. The free energy G is given by [28]

$$G = -\frac{\ln Z}{\beta}, \quad (14)$$

where $\beta = 1/k_{\text{B}}T$ and Z is the partition function. k_{B} is the Boltzmann constant and T is temperature in Kelvin. It is noted that the DFT energies are obtained for fully relaxed structures without any external forces or pressure. Thus, the resulting thermodynamic quantities are effectively obtained in the NPT ensemble (fixed number of particles, fixed pressure and fixed temperature). However, the energy predicted by CE is only valid for the volume leading to the minimum energy of a particular atomic arrangement, and the volume fluctuations are neglected. The free energy can be calculated by utilizing the exact differential

$$\begin{aligned} d(\beta G) &= -\frac{\partial \ln Z}{\partial \beta} d\beta \\ &= U d\beta \end{aligned} \quad (15)$$

where U is the average internal energy. The free energy can be obtained by a thermodynamic integration from a reference temperature β_{ref} where G is known, which can be written as [29]

$$\beta G = (\beta G)_{\text{ref}} + \int_{\beta_{\text{ref}}}^{\beta} d\beta' U(\beta'). \quad (16)$$

Important information of the materials under study such as the stability of ordered/disordered phases can be determined by comparing the free energy of the material at a given composition with respect to the free energies in the pure phases of its constituents.

3. Implementation

CLEAVE utilizes the existing classes and methods of ASE to perform necessary manipulations and analyses for carrying out CE. Among many adopted features, the most noteworthy are the use of

- an `Atoms` object to represent an atomic configuration ($\boldsymbol{\sigma}$),

- a built-in database to efficiently store, maintain and share settings, atomic configurations of the training set, values of the correlation functions ($\phi_\alpha(\sigma)$) and DFT energies,
- Python programming language and modular design to remove the strict input file/format requirements and to enable easy implementation of new features, and
- a `Calculator` class to determine the physical quantity $q(\sigma)$ of a new configuration based on its correlation functions and their ECI values.

It is noted that the evaluation of correlation functions of a new configuration and the determination of physical quantity, $q(\sigma)$, based on ECI values can be a slow process using Python programming language. It is especially true for carrying out Monte Carlo simulations after the CE model training is complete. CLEAVE includes an optional external module written in C++ programming language that can be installed to accelerate the critical and repetitive calculations, but the usage of the code remains unchanged even when the external module is installed (i.e. CLEAVE automatically determines if the C++ add-on is installed, and uses the C++ version if it is present).

The inheritance of the existing features of ASE allows CLEAVE to be fully integrated to ASE where the users can incorporate CE as a part of their research without losing the continuity with the rest of their workflow. The existing users of ASE do not have to install or learn a new CE program nor select a particular DFT package that a CE code supports. In addition to the benefits of integrating CE as a part of ASE, highlights of the features that makes CLEAVE versatile and user-friendly include:

- a multicomponent cluster expansion that goes beyond binary systems,
- a support for several types of single-site basis functions (e.g. basis functions by Sanchez *et al* [6], Van de Walle [20] and Zhang and Sluiter [5]) for a comparison and compatibility with other CE codes,
- many methods for selecting clusters and determining ECI values such as OLS, ℓ_1 and ℓ_2 regularization schemes, Bayesian compressive sensing and genetic algorithm, and
- both leave-one-out and k -fold cross validation schemes.

A simple flowchart illustrating the procedure for constructing CE using CLEAVE is shown in figure 2. The CLEAVE workflow can be divided in to three main components: definition of CE settings, generation of training structures and evaluation of CE convergence. CLEAVE takes an object-oriented approach where each component has its own class. The modular design approach not only enables easy implementation of new features but also makes the code flexible to use and intuitive to follow the CE construction and evaluation procedure shown in figure 2. A more detailed description of main components of the procedure is provided below.

3.1. Definition of cluster expansion settings

The most fundamental component is to define which underlying crystal structure to use. ASE offers two functions to

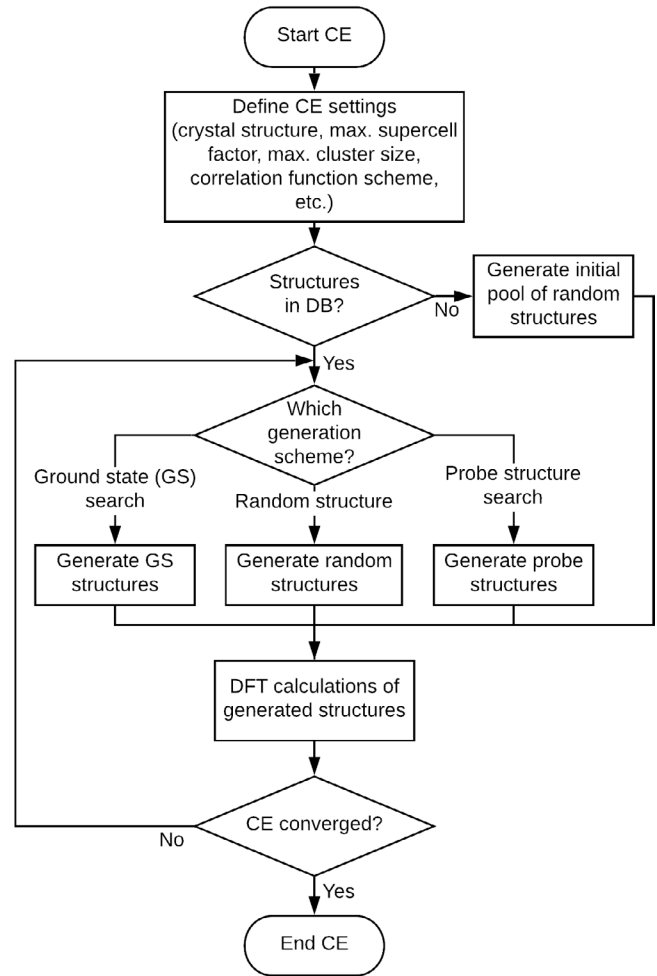


Figure 2. A flowchart of constructing and evaluating CE using CLEAVE.

generate a crystal structure: `bulk` and `crystal`. The `bulk` function provides a simple way of generating common types of crystal structures by specifying the name of the crystal structure and its lattice constant value(s). The crystal structures supported by the `bulk` function are simple cubic, face-centered cubic, body-centered cubic, hexagonal close packed, diamond, zinc blende, rocksalt, cesium chloride, fluorite and wurtzite structures. For more complicated crystal structures, a `crystal` function is used to generate a crystal structure by providing its space group, lattice parameters and scaled coordinates of the unique atomic sites. The definitions of the cluster expansion settings are specified using `CEBulk` and `CECrystal` classes, which respectively calls `bulk` and `crystal` functions to generate an `Atoms` object with the user-specified crystal structure.

The maximum size of the supercell (of the primitive cell) on which the DFT calculations are performed is also defined along with the definition of the underlying crystal structure. The maximum supercell size is specified using a `supercell_factor` parameter, which is an integer corresponding to the product of the absolute values of expansion coefficients (integer weights of a general linear combinations of the unit cell vectors). In other words, if a unit cell has three vectors \vec{a} , \vec{b} and \vec{c} , the configurations in the training set on which the

DFT calculations are performed have cell vectors \vec{a}' , \vec{b}' and \vec{c}' , which are defined as

$$\begin{aligned}\vec{a}' &= i_1\vec{a} + j_1\vec{b} + k_1\vec{c} \\ \vec{b}' &= i_2\vec{a} + j_2\vec{b} + k_2\vec{c} \\ \vec{c}' &= i_3\vec{a} + j_3\vec{b} + k_3\vec{c}\end{aligned}\quad (17)$$

with integer coefficients i_x, j_x and k_x , where $x \in \{1, 2, 3\}$. The `supercell_factor` is then defined as

$$\text{supercell_factor} \geq \prod_{x=1}^3 \prod_{y=1}^3 \prod_{z=1}^3 |i_x| |j_y| |k_z|, \quad (18)$$

and all of the cells used in the training set should have coefficients satisfying the condition in (18). Only unique cell shapes are included in the pool by omitting the cells that can be mapped on to the existing cells in the pool by rotation and reflection. The use of supercells with varying sizes and shapes enables the exploration of a larger configurational space without adding extra computational burden compared to using one fixed supercell size and shape. A set of training structures for CE are later generated iteratively from the pool of possible structures that are realizable using these supercells. To reduce the required computational resources, the structures using smaller supercells are generated (and calculated) first, followed by the larger supercells. The users also have a flexibility to select the supercell size using an optional `size` parameter, which is a 3×3 matrix (or a nested list in Python) specifying the values of the integer coefficients in (17).

Theoretically, an infinite number of clusters can be generated for a given system. The number of clusters is limited to a finite size in practice, and CLEASE takes an approach to generate all possible clusters that are under the truncation threshold (i.e. a maximum number of atoms in clusters and maximum diameter) specified by the user. A whole or subset of the generated clusters is selected during the convergence evaluation process. By default, up to four-body clusters (i.e. empty, one-, two-, three- and four-body clusters) with their diameters up to 5 Å are generated. The users have an option to define their own threshold settings both at the beginning of the CE procedure and at a later stage of the CE iteration cycles. CLEASE also offers `view_clusters` method in `CEBulk` and `CECrystal` classes to visualize the generated clusters in order to assist the user to develop an intuition on the generated clusters.

Within the CE formalism, there does not exist a unique set of definitions for single-site basis functions; the single-site basis functions are considered valid if they form a complete set. Consequently, several definitions are used in practice. The two most widely used definitions are the original definitions by Sanchez *et al* [6] and the one later developed by van de Walle [20], which is used in the Alloy Theoretic Automated Toolkit (ATAT) [20, 30]. The two definitions are equally valid, and both are implemented in CLEASE.

CLEASE offers an option to ignore a set of symmetrically inequivalent atomic sites if they are always occupied by one

element type for all possible configurations. The contributions of these atoms are not explicitly included in the cluster expansion and are automatically included in the constant term (J_0) in (4). For example, lithium metal oxides (LiMO_2) with first-row transition metals ($M = \{\text{Sc}, \text{Ti}, \text{V}, \text{Cr}, \text{Mn}, \text{Fe}, \text{Co}, \text{Ni}, \text{Cu}\}$) have a rocksalt lattice structure [16, 31] with an exception of LiMnO_2 , which is orthorhombic [32]. The rocksalt lattice structure consists of two face-centered cubic sublattices. For the case of the cation-disordered rocksalt lattice LiMO_2 , one sublattice is occupied by lithium and other metal atoms while the other is occupied by oxygen atoms. The complexity of the CE model of such systems can be reduced to a cation sublattice consisting only of two element types (the oxygen sublattice is ignored). As such, an optional Boolean argument is present in CLEASE to enable/disable the reduction of the complexity of the model by ignoring the such atoms if they exist in the system. The reduction of the model complexity leads to a reduction in computational cost as it requires a smaller number of interaction terms to calculate.

A range of compositions (or concentration) of the system to be studied is specified using a `Concentration` class. First, the constituting elements of the system are categorized into the basis which they belong. For example, LiVO_2 in a rocksalt lattice structure is expressed using two lists: ['Li', 'V'] and ['O']. It is noted that CE needs to keep track of the location of vacancies when they are present in the system. The location of vacancies are tracked by treating a vacancy as a regular atom with its atomic symbol set to 'X' or atomic number set to zero. The LiVO_2 with Li vacancies is then expressed using ['Li', 'V', 'X'] and ['O'].

The range of each element (including vacancies) can be specified in one of the two convenient methods built in to the `Concentration` class. The simplest method is to specify the concentration range of each constituting element by calling `set_conc_ranges` method in `Concentration` class. For the cases where concentrations of two or more elements depend on one another, one can specify concentration range using `set_conc_formula_unit` method where the relationships between the concentrations of two or more elements can be expressed in a list of strings. For the example of LiVO_2 with Li vacancies, a list of strings that specifies relationship between the number of Li atoms and the number of vacancies, ['Li<x>V<1-x>', 'O<2>'], is passed as an argument to the `set_conc_formula_unit` method. Another argument specifying the range of the concentration variable, e.g. {'x': (0, 1)}, is also passed to the `set_conc_formula_unit` method in order to specify the concentration range of Li and Li vacancies. The concentration ranges specified by either `set_conc_ranges` or `set_conc_formula_unit` methods are internally interpreted in the `Concentration` class as a list of linear equations that specify (1) the relationships of the concentrations of constituting elements and (2) their upper/lower bounds. The advanced users can alternatively specify the coefficients of the linear equations used in the `Concentration` class if a greater flexibility is needed in specifying the concentration ranges.

3.2. Generation of training structures

CLEAVE uses `NewStructures` class to generate training structures, which provides three different methods perform the task. The first and most trivial method is to generate a set of random structures. This method serves to generate an initial pool consisting of a small number of structures. The random generation method is used in the first iteration cycle of CE construction as shown in figure 2. An initial cluster expansion is capable of making a first set of predictions albeit with a low accuracy. It is noted that all of the generated training structures, along with their correlation function values, are stored in a database file.

Once the initial CE is constructed, the user is given three different choices for introducing an additional set of training structures. The first and most straightforward option is to keep generating random structures. Although trivial, generating random structures is claimed to be the best strategy when compressive sensing is used to select clusters [22]. The second method is to generate ground-state and other low-energy structures based on current cluster expansion (i.e. based on the pool of structures already calculated) [33], which have the enthalpies of formation either on or near the convex hull [34]. The inclusion of ground-state and near-ground-state structures serves an important purpose of increasing the accuracy in predicting the correct ground states. A global minimization technique can be used to generate (near) ground-state configurations, and CLEAVE uses a simulated annealing technique. Simulated annealing takes an initial atomic configuration and cures it at a sequence of decreasing temperatures in order to let the structure evolve towards the ground-state arrangement. In CLEAVE, the user can generate ground-state structure by invoking `generate_gs_structure` method in `NewStructures` class while specifying the initial and final temperatures (intermediate temperatures are interpolated in a logarithmic scheme), number of temperatures and the number of site swaps per temperature.

The last method of generating the training set is referred to as a ‘probe structure’ method [23, 35]. The probe structure method introduces a new structure that minimizes the mean variance of the predicted physical quantity $q(\sigma)$. The mean variance of the predicted quantity q is written as [35]

$$\begin{aligned} \text{Var}[\hat{q}_i] &= \frac{1}{N_{\text{config}}} \sum_{i=1}^{N_{\text{config}}} [\mathbf{X}_i(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}_i^T]e^2 \\ &= \{\text{tr}[(\mathbf{X}^T\mathbf{X})^{-1}\Sigma] + \mu(\mathbf{X}^T\mathbf{X})^{-1}\mu^T\}e^2 \\ &= \Lambda \cdot e^2, \end{aligned} \quad (19)$$

where e^2 is the variance of the error in the training set, Σ is the covariance matrix of the correlation functions of the training set and μ is a vector of the mean correlation functions of the structures in the training set. The probe structure is the one that reduce the value of Λ the most when introduced to the training set, which is found using the simulated annealing procedure. Similar to generating ground-state structures, the probe structures are generated by invoking `generate_probe_structure` method in `NewStructures` class.

The newly generated structures are compared with the existing structures in the training set in order to avoid introducing duplicate structures. We adopted the structure comparison algorithm developed by Lonie and Zurek [36] to identify equivalent structures. It is desirable to have the new structure compared against the existing structures in the training set as efficiently as possible. As a first step, the structures that have different chemical composition than the newly generated structure are filtered, and the new structure is compared only with the remaining structures. Once the candidate transformations for mapping the new structure onto one structure in the database are identified using the algorithm suggested by Lonie and Zurek, we note that exactly the same transformations can be used for the remaining structures in the database. Therefore, the structure comparison algorithm implemented in ASE is optimized for the case where one structure is to be compared against many.

In addition to the aforementioned methods of generating the training structures, CLEAVE also offers a built-in function to import structures to the database. The import function also has an option to specify the calculated q value, which allows users to easily import the previously calculated results.

3.3. Evaluation of cluster expansion convergence

An evaluation process to determine the convergence of CE includes a selection of clusters, a determination of their ECI values and an assessment of the LOO or k -fold CV score using the selected clusters and their ECI values. An entire evaluation process is performed using an `Evaluate` class.

The simplest way to determine the ECI values of the generated clusters is by using OLS to minimize RESS. It is highly likely that the ECI values found using OLS are overfitted. Therefore, Bayesian compressive sensing and ℓ_1 and ℓ_2 regularization methods are implemented, and it is highly recommended to use a regularization methods to select clusters and evaluate their ECI values.

A default option in the `Evaluate` class is to include all of the clusters generated using the cluster truncation conditions specified in `CEBulk` or `CECrystal` class, and either the entire or a subset of these clusters are selected for fitting depending on the method used. The `Evaluate` class provides additional options in which the users can select a subset of the generated clusters to perform any of OLS, Bayesian compressive sensing and ℓ_1 and ℓ_2 regularization. The first option is by manually specifying which clusters to include, while the second option is to provide a stricter truncation conditions than the ones set in the `CEBulk` or `CECrystal` class. The first option allows the `Evaluate` class to be used in conjunction with other cluster selection methods such as genetic algorithm. For example, a user can optionally use genetic algorithm (included in CLEAVE as a separate `GAFit` class) to pre-screen a large cluster pool and subsequently pass a subset of clusters to the `Evaluate` class. The feature to freely select a subset of a large pool of clusters along with the use of OLS, Bayesian compressive sensing and ℓ_1 and ℓ_2 regularization methods allows the users to easily experiment with

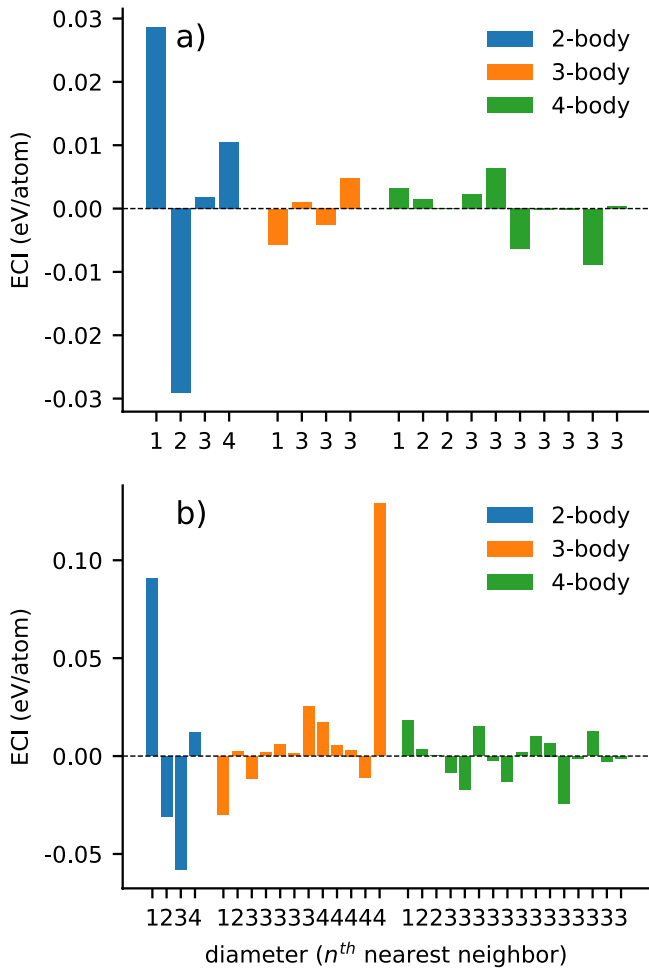


Figure 3. ECIs obtained via (a) ℓ_1 regularization and (b) ℓ_2 regularization.

various settings to understand how the system behaves and to optimize the ECI values for achieving the lowest LOOCV score.

To further assist the evaluation process, the `Evaluate` class contains two built-in methods that automatically determine the LOOCV when a regularization method is used. The first method, `plot_fit`, determines and stores the selected clusters and their ECI values for a value of regularization parameter (λ) specified by the user. It also plots the fit of all data points in the training set to their calculated values and presents the LOO/ k -fold CV score of the specified λ value. Since the most cumbersome task in determining the convergence of CE is finding the optimal λ value that yields the lowest CV score, another method, `plot_cv`, is also implemented. It takes a range and number of λ values to evaluate as inputs and returns the best λ value in the specified range along with its LOO/ k -fold CV score. The `plot_cv` method also plots LOO/ k -fold CV score as a function of λ and provides an option to store the results in a log file such that the users can add more λ values to the list at a later stage without having to re-evaluate the same λ values in the process.

3.4. Metropolis Monte Carlo and simulated annealing

The user can perform statistical sampling of the system on a larger simulation cell once the cluster expansion is constructed. The final selection of cluster and their ECI values can be stored and passed to other classes to conduct statistical analyses. A separate `Calculator` class for cluster expansion is implemented in ASE. The `Cleas` calculator class takes a list of clusters and their ECI values as inputs, and the users can select what type of trial moves are allowed. The sampling in the canonical ensemble allows the swapping two atoms with different constraint conditions (i.e. swap any two atoms, swap any two atoms in the same basis, swap two nearest neighbors, swap two nearest neighbors in the same basis) while the semi-grand canonical ensemble allows changing the type of occupying element at a random site.

The evaluation of the physical quantity $q(\sigma)$ is performed using (5), which is a fast because the `Cleas` calculator keeps track of the changes in the `Atoms` object to update the correlation functions. When the physical quantity being modeled is energy, a trial move of the standard Metropolis algorithm has an acceptance probability [37]

$$P_{\text{acc}} = \min \left\{ 1, \exp \left(\frac{-(E_{\text{new}} - E_{\text{old}})}{k_{\text{B}}T} \right) \right\}, \quad (20)$$

where E_{new} and E_{old} are the energy of the new and old configuration, respectively. As the `Cleas` calculator keeps track of the change in the `Atoms` object after each move, updating the correlation functions is restricted to the contributions of one and two atoms for the semi-grand canonical ensemble and canonical ensemble, respectively.

4. Examples

Here, we present two example systems to illustrate the capabilities of the CLEAS code. The first example illustrates the investigation of a Au–Cu binary alloy. The second example shows the cluster expansion on a more complex $\text{Li}_2\text{CrO}_2\text{F}$ system consisting of four types of elements and vacancy. All of the interactions of cluster expansions are computed from DFT calculations of energies, and the computational settings used for generating the results shown in this section are specified in section 5.

4.1. Au–Cu alloy

The binary Au–Cu alloy system is studied at temperatures ranging from 100 K to 800 K over the entire composition range. The resulting values obtained for both ℓ_1 and ℓ_2 regularization are shown in figure 3. The ECI value of the empty cluster is found to be -3.49 eV/atom for both cases, and the ECI value of the one-body cluster is 0.27 eV/atom and 0.13 eV/atom for ℓ_1 and ℓ_2 regularization, respectively. The ECI values of empty and one-body clusters are not included in figure 3 for better

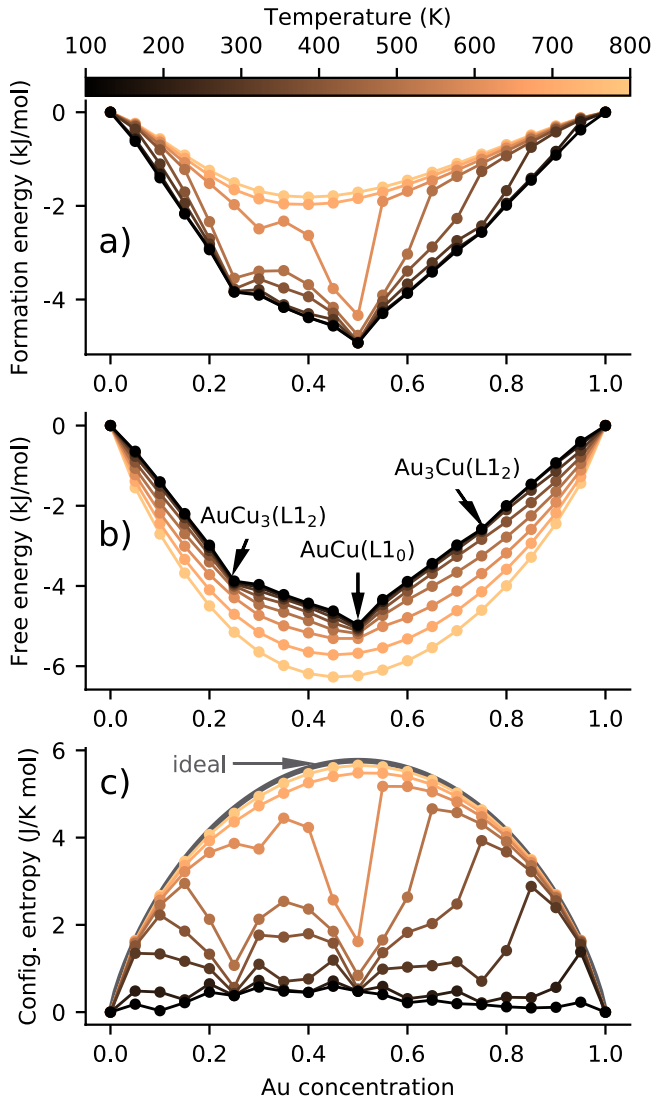


Figure 4. Thermodynamic quantities for the Au–Cu computed for temperatures ranging from 100 K to 800 K over the entire composition range. (a) Formation energy. (b) Free energy of formation. (c) Configurational entropy.

visibility. The LOOCV score for the ℓ_1 - and ℓ_2 -regularized fit were 4.49 meV/atom and 4.67 meV/atom, respectively. The ℓ_1 regularization scheme yields a slightly lower CV score despite having a smaller number of clusters (ℓ_1 -regularized fit has 20 clusters while the ℓ_2 -regularized fit has 34 clusters).

A qualitative information on the thermodynamic behavior of the system can be extracted by inspecting the ECI values for simple binary system. Based on the fact that the energetically favorable configurations have DFT energies that are more negative than less favorable ones and that the two site variables are +1 or -1, one can infer that a positive ECI value of the pair interaction term means that a pair consisting of two different elements is energetically preferred at a low temperature. It can be seen in figure 3 that the ECIs of the nearest-neighbor and second nearest-neighbor pairs are positive and positive, respectively. The signs indicate that the Au–Cu system energetically favors the strong mixing of the constituting elements such that the alternating patterns found

in L1₀- and L1₂-type ordered structures are likely to emerge, which is in a good agreement with experimental and computational observations [33, 38–45].

It is experimentally determined that Au–Cu alloys have three ordered phases at low temperatures [43–45]: AuCu₃, AuCu and Au₃Cu. Furthermore, the transition temperatures for AuCu₃, AuCu and Au₃Cu are reported to be 663 K, 683 K and ~490 K, respectively, and they are often used as reference values for assessing the computational models [33, 38, 39]. The formation energy, free energy of formation and configurational entropy are obtained through Metropolis Monte Carlo simulations and are shown in figure 4. As the CE is trained with fully relaxed structures (zero pressure), the formation energy is determined using

$$\Delta U = U - xU_{\text{Au}} - (1-x)U_{\text{Cu}}, \quad (21)$$

where U is the internal energy of the configuration, x is the gold concentration, U_{Au} is the internal energy of pure gold and U_{Cu} is the internal energy of pure copper. Similarly, the free energy of formation is obtained by subtracting the weighted average of the free energy for the pure phases. The configurational entropy is given by the difference between the internal energy and the free energy, divided by the temperature at which the Monte Carlo is sampled. The three ordered phases (AuCu₃, AuCu and Au₃Cu) are found on the convex hull of the free energy of formation in figure 4(b). Furthermore, the entropy of the ordered phases form local minima as shown in figure 4(c). For comparison, the entropy of an ideal mixture, S_{ideal} , defined as [46, 47]

$$S_{\text{ideal}} = k_B \ln \frac{N!}{n_{\text{Au}}!n_{\text{Cu}}!} \approx -k_B \left[\frac{n_{\text{Cu}}}{N} \ln \frac{n_{\text{Cu}}}{N} + \frac{n_{\text{Au}}}{N} \ln \frac{n_{\text{Au}}}{N} \right], \quad (22)$$

where N , n_{Au} and n_{Cu} are the number of atomic sites, the number of sites occupied by Au and Cu atoms, respectively. The entropy of an ideal mixture is included as a gray line in figure 4(c). The entropy curve resembles that of an ideal mixture at the high-temperature limit at 800 K. The curve starts to deviate from that of the ideal mixture as the temperature is lowered mainly because the entropy drops sharply at concentrations at which the ordered phases (e.g. L1₀ and L1₂) are energetically preferred. As the temperature increases, the free energy becomes a smooth convex curve with a minimum at around 50% composition, and the system is in a random phase with no short-range order.

An accurate estimate of the order/disorder transition temperature can be found by tracking the evolution of an order parameter. The average fraction of sites having a different element than the same site in the ground state, f_{diff} , is tracked as the system evolves. f_{diff} is normalized by the expected fraction of different sites in a random phase, $f_{\text{diff,md}}$, and an order parameter, η , is defined as

$$\eta = 1 - f_{\text{diff}}/f_{\text{diff,md}}. \quad (23)$$

The order parameter is used for detecting the phase transition as shown in figure 5. The computationally predicted order/disorder transition temperature of AuCu₃, AuCu and Au₃Cu are around 600 K, 665 K and 385 K, respectively,

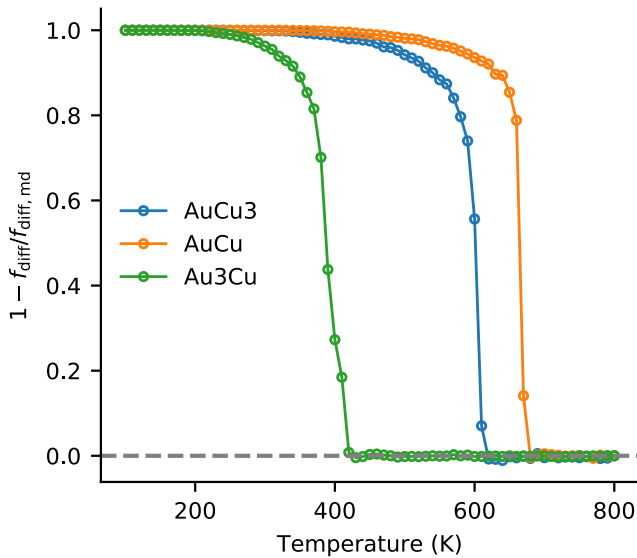


Figure 5. Order parameter as a function of temperature (1 in an ordered phase and 0 in a random phase).

which are in a good agreement with the experimental reference values [43–45].

One of the most common way to describe the characteristics of a binary alloy is by constructing a phase diagram. A phase diagram can be generated computationally using a semi-grand canonical MC where a grand potential is obtained via thermodynamic integration in (16) at fixed chemical potentials. The integration starts from the low temperature limit for the ordered phases and from the high temperature limit for disordered phases where the free energy per atom is given by $k_B T \ln 2$. The phase boundary between two phases is identified by locating the intersection point between the grand potential in the two co-existing phases. The phase diagram generated via semi-grand canonical MC is shown in figure 6. The phase diagram closely resembles the phase diagrams constructed from the experimental measurements [43–45] and is also in a qualitatively agreement with the phase diagrams constructed from computational results [33, 38, 48].

4.2. Lithium chromium oxyfluoride

One of the recent focus areas of lithium-ion battery research is the development of high-capacity cathode materials. Lithium metal oxyfluorides ($\text{Li}_2\text{MO}_2\text{F}$, $\text{M} = \{\text{V}, \text{Cr}, \text{Mn}, \text{Ti}, \text{Ni}, \dots\}$) is a family of materials that is at the forefront of the current research. The challenges for studying $\text{Li}_2\text{MO}_2\text{F}$ is in the vast size of the configurational space, which exhibit not only the cation disorder commonly found in lithium metal oxides [14, 16] but also anion disorder which is also present due to the mixed O/F composition [49, 50]. The fact that the underlying crystal structure of $\text{Li}_2\text{MO}_2\text{F}$ can vary at different lithiation levels [51] adds the complexity to investigate their properties. It is, however, known that the most predominant crystal structure is of disordered rocksalt type [52], particularly at high-lithiation levels. We therefore show an example CE study of $\text{Li}_2\text{CrO}_2\text{F}$ in a rocksalt lattice configuration.

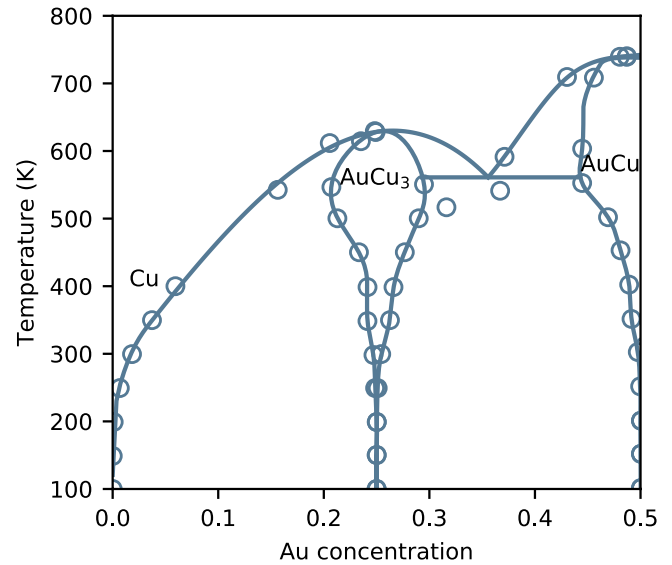


Figure 6. Phase diagram of $\text{Au}_x\text{Cu}_{1-x}$ where $0 \leq x \leq 0.5$. Circles are computed phase boundary points and lines are spline fits of the computed boundary points.

The Monte Carlo annealing study reveals that $\text{Li}_2\text{CrO}_2\text{F}$ (i.e. fully lithiated compound) takes a layered structure at room temperature (293 K) as shown in figure 7(a). The layer structure shows a $\dots\text{Li}-\text{F}-\text{Li}-\text{O}-\text{Cr}-\text{O}-\dots$ pattern, which is similar to a $\dots\text{Li}-\text{O}-\text{M}-\text{O}-\dots$ layered pattern observed in lithium metal oxides [16, 53, 54]. The layered structure is lost upon delithiation, which leads to disordered structures as shown in figure 7(b). The emergence of disordered structures agrees well with the previous experimental observations [50, 52], and it is important to model the disordered atomic arrangement as it has a direct link to the Li transport mechanism (e.g. a presence of zero-transition-metal pathways [16, 31, 55]).

Thermodynamics quantities of $\text{Li}_x\text{CrO}_2\text{F}$ can be extracted with the same procedure described for the Au–Cu system. One of the most crucial thermodynamic parameters for characterizing cathode materials for Li-ion batteries is the free energy as it is directly linked to the operating voltage of the cell. The operating voltage of $\text{Li}_x\text{CrO}_2\text{F}$ is defined as

$$\begin{aligned} \text{voltage} &= -\frac{\mu_{\text{Li}}^{\text{cathode}} - \mu_{\text{Li}}^{\text{anode}}}{e} \\ &= -\frac{\frac{dG_{\text{Li}_x\text{Cr}_2\text{F}}}{dx} - \mu_{\text{Li}}^{\text{anode}}}{e}, \end{aligned} \quad (24)$$

where μ_{Li} is the chemical potential in eV per Li atom, e is an electron charge and $G_{\text{Li}_x\text{Cr}_2\text{F}}$ is the free energy of $\text{Li}_x\text{CrO}_2\text{F}$ in eV per formula unit. Li metal is used as an anode and thus, $\mu_{\text{Li}}^{\text{anode}}$ is constant.

The free energy of $\text{Li}_x\text{CrO}_2\text{F}$ and its voltage profile at 293 K are shown in figure 8. The free energy in figure 8(a) has three parts: free energy values computed from MC simulations, a smooth curve fitted to the computed values (using Redlich–Kister polynomials [56]) and a convex hull of the fitted curve. The curve fit is used for generating the voltage plot because the derivative of the free energy used for calculating the voltage

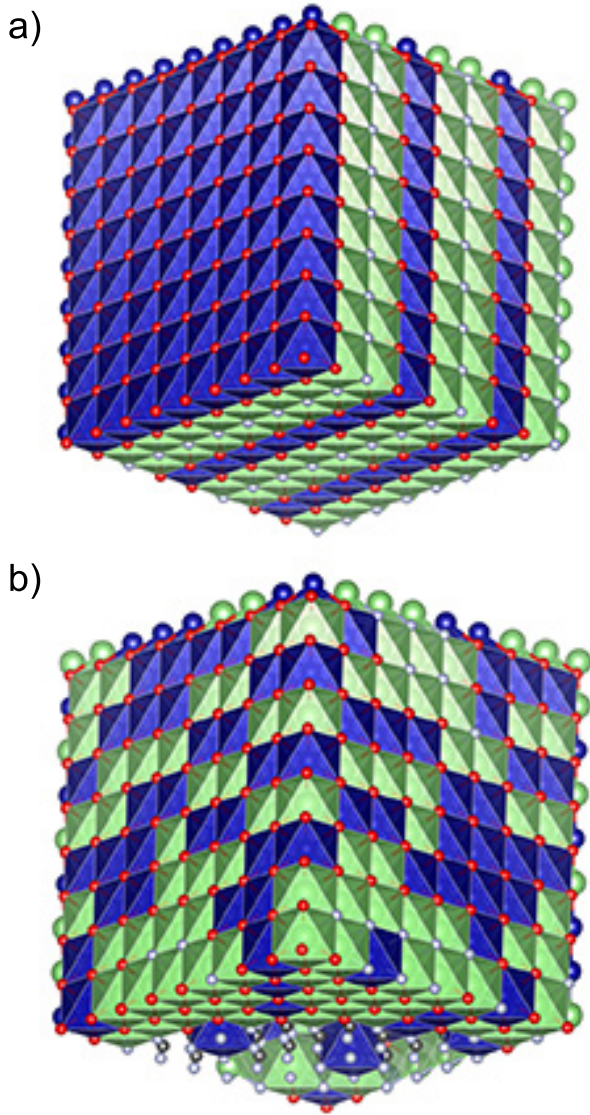


Figure 7. A snapshot of $\text{Li}_x\text{CrO}_2\text{F}$ during the Monte Carlo run at 293 K where x is (a) 2.0, (b) 1.5. The Li atoms are shown in green, the Cr atoms are shown in blue, the oxygen atoms are shown in red and the F atoms are shown in white.

values are susceptible to small noise that are present in the MC simulation results. Furthermore, a range in which the free energy curve is above the convex hull represents the region where a phase transition occurs: the cathode forms a mixture of two phases at which the fitted curve and the convex hull intersect. The voltage profile in figure 8(b) is generating using (24) where the values on the convex hull are used for $G_{\text{Li}_x\text{Cr}_2\text{F}}$. The voltage profile in figure 8(b) is in a good agreement with those observed experimentally [50, 52].

5. Methods

5.1. Density functional theory calculations

All of the calculations are performed with the Vienna *Ab initio* Simulation Package (VASP) [57–60] using the projector augmented-wave (PAW) method [61]. The generalized

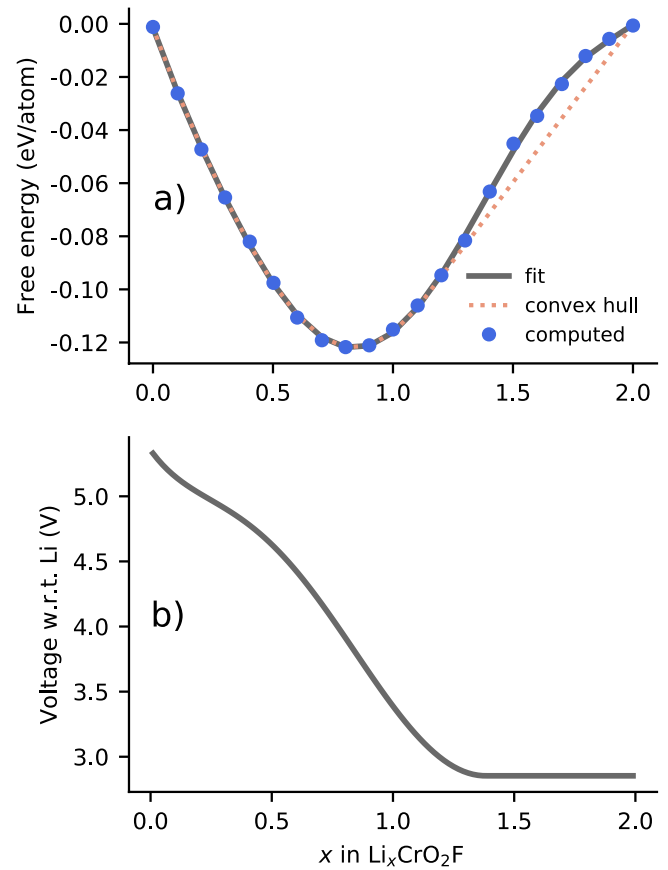


Figure 8. Free energy of formation for $\text{Li}_x\text{CrO}_2\text{F}$ and its voltage with respect to Li metal at 293 K.

gradient approximation as parametrized by Perdew, Burke and Ernzerhof [62] is used as the exchange-correlation functional. It is important to have a consistent and accurate dataset (i.e. DFT calculations with high energy cutoff and k -point mesh density) in order to minimize the numerical noise introduced to the CE training. The plane-wave cutoff of 500 eV is used, and both the cell and atomic positions are fully relaxed such that all the forces are smaller than $0.02 \text{ eV } \text{\AA}^{-1}$. A rotationally invariant Hubbard U correction [63, 64] is applied to the d orbital of Cr with the U value of 3.7 eV. The calculations are performed with supercells containing up to 18 and 54 atoms for Au–Cu alloy and $\text{Li}_2\text{CrO}_2\text{F}$ systems, respectively. Integrations over the Brillouin zone were carried out using the Monkhorst–Pack scheme [65] with a grid with a maximal interval of 0.04 \AA^{-1} .

5.2. Cluster expansion model

The CE model for Au–Cu alloy and $\text{Li}_x\text{CrO}_2\text{F}$ are trained using 34 and 390 DFT calculations, respectively. CE model is trained for the entire composition range of Au–Cu alloy (from pure Au to pure Cu) and $\text{Li}_x\text{CrO}_2\text{F}$ on a rocksalt lattice with x ranges from 0 to 2. Up to four-body clusters with the maximum diameter of 6.0 \AA are generated for Au–Cu alloy. Up to four-body clusters are generated for $\text{Li}_x\text{CrO}_2\text{F}$ with the maximum diameter of 7.0 for two- and three-body clusters and 4.5 \AA for four-body clusters. ℓ_1 and ℓ_2 regularization schemes

with the regularization parameter ranging from 10^{-7} to 10^2 are assessed at various maximum radii to find the optimal setting that leads to the lowest LOOCV score. For the Au–Cu alloy, ℓ_1 regularization with the maximum diameter of 6.0 Å, 5.0 Å and 5.0 Å for 2-, 3- and 4-body clusters, respectively, yields the lowest LOOCV score of 4.49 meV/atom. The minimum LOOCV score achieved using ℓ_2 regularization scheme is 4.67 meV/atom when the maximum diameter is set to 6.0 Å, 6.0 Å and 5.0 Å for 2-, 3- and 4-body clusters, respectively. Similarly, ℓ_1 regularization performed better than ℓ_2 regularization on $\text{Li}_x\text{CrO}_2\text{F}$ with the lowest LOOCV score of 21.38 meV/atom (maximum diameter set to 7.0 Å, 7.0 Å and 4.5 Å for 2-, 3- and 4-body clusters, respectively). It is noted that although the LOOCV of $\text{Li}_x\text{CrO}_2\text{F}$ seems larger compared to that of Au–Cu, it should be taken into account that the cohesive energy of metallic alloys are in general much smaller than those of oxyfluorides.

5.3. Metropolis Monte Carlo simulations

For Au–Cu alloy, Metropolis Monte Carlo simulations are carried out using a $10 \times 10 \times 10$ supercell consisting of 1000 atoms for determining thermodynamic quantities. The system is equilibrated with 100 sweeps, and an average energy is collected through an additional 2000 sweeps at each temperature for determining the thermodynamic quantities. A $30 \times 30 \times 30$ supercell consisting of 27000 atoms is used to determine the transition temperatures and to construct a phase diagram. The transition temperatures are determined by equilibrating the systems with 100 sweeps, followed by sampling the order parameter via an additional 1000 sweeps. A phase diagram is generated by performing a semi-grand canonical MC, where the system is equilibrated using 100 sweeps, followed an additional 1000 sweeps to obtain an average semi-grand canonical energy at each temperature at a fixed chemical potential. A $9 \times 9 \times 9$ cell consisting of 1458 atoms is used for $\text{Li}_x\text{CrO}_2\text{F}$. The temperature is gradually lowered from 10000 K, and the structures are equilibrated at each temperature via 100 sweeps to ensure that the system is equilibrated before sampling. The average energy is then sampled via 1000 sweeps at each temperature.

6. Conclusions

We present the implementation of CLEASE, which fully integrates the cluster expansion method to ASE package. The aim of the developed code is to make cluster expansion more accessible to non-specialists and to incorporate modern machine learning techniques to cluster expansion method in one comprehensive and versatile package. The use of the popular Python programming language and implementing the code as a part of widely used ASE package lowers the barrier for the newcomers to the field to easily learn and use CE as a part of their research methods. By automatically generating clusters and calculating the correlation functions of both semi-automatically generated and user-supplied structures, it minimizes both the possible introduction of user errors and

complicated process of constructing/evaluating the cluster expansion. The capability of CLEASE is presented with two example usage cases with a different level of system complexity. The examples demonstrate that CE can correctly predict the material behavior that require statistical sampling on a large simulation cell.

Acknowledgments

This project has received funding from the European Union's Horizon 2020 research and innovation programme under Grant Agreement No 711792 (FET-OPEN project LiRichFCC). Authors thank valuable discussions with Dr Juhani Teeriniemi and Prof Kari Laasonen. JMGL acknowledges support from the Villum Foundation's Young Investigator Programme (4th round, project: *in silico design of efficient materials for next generation batteries*. Grant Number: 10096).

ORCID iDs

Jin Hyun Chang  <https://orcid.org/0000-0003-0668-4530>

Juan Maria Garcia-Lastra  <https://orcid.org/0000-0001-5311-3656>

Tejs Vegge  <https://orcid.org/0000-0002-1484-0284>

References

- [1] Fontaine D D 1994 Cluster approach to order-disorder transformations in alloys *Solid State Physics* vol 47, ed H Ehrenreich and D Turnbull (London: Academic) pp 33–176
- [2] Zunger A, Wang L G, Hart G L W and Sanati M 2002 *Modelling Simul. Mater. Sci. Eng.* **10** 685–706
- [3] Lerch D, Wieckhorst O, Hart G L W, Forcade R W and Müller S 2009 *Modelling Simul. Mater. Sci. Eng.* **17** 055003
- [4] Meng Y S and Arroyo-de Dompablo M E 2009 *Energy Environ. Sci.* **2** 589
- [5] Zhang X and Sluiter M H F 2016 *J. Phase Equilib. Diffus.* **37** 44–52
- [6] Sanchez J M, Ducastelle F and Gratias D 1984 *Phys. A: Stat. Mech. Appl.* **128** 334–50
- [7] Zunger A 1994 First-principles statistical mechanics of semiconductor alloys and intermetallic compounds *Statics and Dynamics of Alloy Phase Transformations* ed P E A Turchi and A Gonis (Boston, MA: Springer) pp 361–419
- [8] Asta M, Ozolins V and Woodward C 2001 *JOM* **53** 16–9
- [9] van de Walle A 2008 *Nat. Mater.* **7** 455–8
- [10] Magri R and Zunger A 1991 *Phys. Rev. B* **44** 8672–84
- [11] Franceschetti A and Zunger A 1999 *Nature* **402** 60–3
- [12] Geng H Y, Sluiter M H F and Chen N X 2005 *J. Chem. Phys.* **122** 214706
- [13] Wang R, Li X, Liu L, Lee J, Seo D H, Bo S H, Urban A and Ceder G 2015 *Electrochem. Commun.* **60** 70–3
- [14] Abdellahi A, Urban A, Dacek S and Ceder G 2016 *Chem. Mater.* **28** 3659–65
- [15] Abdellahi A, Urban A, Dacek S and Ceder G 2016 *Chem. Mater.* **28** 5373–83
- [16] Urban A, Matts I, Abdellahi A and Ceder G 2016 *Adv. Energy Mater.* **6** 1600488
- [17] Kitchaev D A *et al* 2018 *Energy Environ. Sci.* **11** 2159–71

- [18] Hjorth Larsen A *et al* 2017 *J. Phys.: Condens. Matter* **29** 273002
- [19] Zarkevich N A and Johnson D D 2004 *Phys. Rev. Lett.* **92** 255702
- [20] van de Walle A 2009 *Calphad* **33** 266–78
- [21] Nelson L J, Ozoliš V, Reese C S, Zhou F and Hart G L W 2013 *Phys. Rev. B* **88** 155105
- [22] Nelson L J, Hart G L W, Zhou F and Ozoliš V 2013 *Phys. Rev. B* **87** 035125
- [23] Seko A and Tanaka I 2014 *J. Phys.: Condens. Matter* **26** 115403
- [24] Mueller T and Ceder G 2010 *Phys. Rev. B* **82** 184107
- [25] Blum V, Hart G L W, Walorski M J and Zunger A 2005 *Phys. Rev. B* **72** 165113
- [26] Hart G L W, Blum V, Walorski M J and Zunger A 2005 *Nat. Mater.* **4** 391–4
- [27] Díaz-Ortiz A, Dosch H and Drautz R 2007 *J. Phys.: Condens. Matter* **19** 406206
- [28] Andersen J O 2012 *Introduction to Statistical Mechanics* (Trondheim: Akademica Publishing)
- [29] Tuckerman M E 2010 *Statistical Mechanics: Theory and Molecular Simulation* (Oxford: Oxford University Press)
- [30] van de Walle A, Asta M and Ceder G 2002 *Calphad* **26** 539–53
- [31] Urban A, Lee J and Ceder G 2014 *Adv. Energy Mater.* **4** 1400478
- [32] Hewston T and Chamberland B 1987 *J. Phys. Chem. Solids* **48** 97–108
- [33] van de Walle A and Ceder G 2002 *J. Phase Equilib.* **23** 348
- [34] Urban A, Seo D H and Ceder G 2016 *npj Comput. Mater.* **2** 16002
- [35] Seko A, Koyama Y and Tanaka I 2009 *Phys. Rev. B* **80** 165122
- [36] Lonie D C and Zurek E 2012 *Comput. Phys. Commun.* **183** 690–7
- [37] Thijsen J 2007 *Computational Physics* (Cambridge: Cambridge University Press)
- [38] Wei S H, Mbaye A A, Ferreira L G and Zunger A 1987 *Phys. Rev. B* **36** 4163–85
- [39] Ozoliš V, Wolverton C and Zunger A 1998 *Phys. Rev. B* **57** 6427–43
- [40] Ozoliš V, Wolverton C and Zunger A 1998 *Phys. Rev. B* **58** R5897(R)
- [41] Wolverton C, Ozoliš V and Zunger A 1998 *Phys. Rev. B* **57** 4332–48
- [42] Lysgaard S, Mýrdal J S G, Hansen H A and Vegge T 2015 *Phys. Chem. Chem. Phys.* **17** 28270–6
- [43] Massalski T B, Murray J L, Bennett L H and Baker H 1986 *Binary Alloy Phase Diagrams* (Metals Park, OH: American Society for Metals)
- [44] Hultgren R, Hawkins D T and Desai P D 1973 *Selected Values of the Thermodynamic Properties of Binary Alloys* (Metals Park, OH: American Society for Testing and Materials)
- [45] Fedorov P P and Volkov S N 2016 *Russ. J. Inorg. Chem.* **61** 772–5
- [46] Reuter K and Scheffler M 2003 *Phys. Rev. B* **68** 045407
- [47] Murty B S, Yeh J W and Ranganathan S 2014 *High-Entropy Alloys* (Oxford: Butterworth-Heinemann)
- [48] Walle A V D and Asta M 2002 *Modelling Simul. Mater. Sci. Eng.* **10** 521–38
- [49] Chen R, Ren S, Knapp M, Wang D, Witter R, Fichtner M and Hahn H 2015 *Adv. Energy Mater.* **5** 1401814
- [50] Chen R, Ren S, Mu X, Maawad E, Zander S, Hempelmann R and Hahn H 2016 *ChemElectroChem* **3** 892–5
- [51] Cambaz M A, Vinayan B P, Clemens O, Munnangi A R, Chakravadhanula V S K, Kübel C and Fichtner M 2016 *Inorg. Chem.* **55** 3789–96
- [52] Ren S, Chen R, Maawad E, Dolotko O, Guda A A, Shapovalov V, Wang D, Hahn H and Fichtner M 2015 *Adv. Sci.* **2** 1500128
- [53] Mizushima K, Jones P, Wiseman P and Goodenough J 1981 *Solid State Ion.* **3–4** 171–4
- [54] Van der Ven A and Ceder G 2004 *Electrochem. Commun.* **6** 1045–50
- [55] Lee J, Urban A, Li X, Su D, Hautier G and Ceder G 2014 *Science* **343** 519–22
- [56] Redlich O and Kister A T 1948 *Ind. Eng. Chem.* **40** 345–48
- [57] Kresse G and Hafner J 1993 *Phys. Rev. B* **47** 558–61
- [58] Kresse G and Hafner J 1994 *Phys. Rev. B* **49** 14251–69
- [59] Kresse G and Furthmüller J 1996 *Comput. Mater. Sci.* **6** 15–50
- [60] Kresse G and Furthmüller J 1996 *Phys. Rev. B* **54** 11169–86
- [61] Blöchl P E 1994 *Phys. Rev. B* **50** 17953–79
- [62] Perdew J P, Burke K and Ernzerhof M 1996 *Phys. Rev. Lett.* **77** 3865–8
- [63] Anisimov V I, Zaanen J and Andersen O K 1991 *Phys. Rev. B* **44** 943–54
- [64] Cococcioni M and de Gironcoli S 2005 *Phys. Rev. B* **71** 035105
- [65] Monkhorst H J and Pack J D 1976 *Phys. Rev. B* **13** 5188–92