# Active learning of linearly parametrized interatomic potentials

Evgeny V. Podryabinkin, Alexander V. Shapeev *

*Skolkovo Institute of Science and Technology, Moscow, Russia*

## ARTICLE INFO

## ABSTRACT

This paper introduces an active learning approach to the fitting of machine learning interatomic potentials. Our approach is based on the D-optimality criterion for selecting atomic configurations on which the potential is fitted. It is shown that the proposed active learning approach is highly efficient in training potentials on the fly, ensuring that no extrapolation is attempted and leading to a completely reliable atomistic simulation without any significant decrease in accuracy. We apply our approach to molecular dynamics and structure relaxation, and we argue that it can be applied, in principle, to any other type of atomistic simulation. The software, test cases, and examples of usage are published at http://gitlab.skoltech.ru/shapeev/mlip/.

## 1. Introduction

Many research areas in materials science, molecular physics, chemistry, and biology involve atomistic modeling. For example, in molecular dynamics (MD), as a rule, one of the following two classes of interatomic interaction models is used. The first class is the empirical interatomic potentials—they are very computationally efficient and allow for simulating large atomistic systems for microseconds of simulation time. However, they typically yield only qualitative accuracy. The other class is quantum-mechanical (QM) models, such as the density functional theory (DFT). They are very accurate, but computationally expensive. Their applicability is typically limited to hundreds of atoms and hundreds of picoseconds of simulation time.

Several directions of developing the models that would be both accurate and computationally efficient have been pursued. They include the so-called linear scaling DFT [1–3] that ensures that the algorithmic complexity grows linearly when the size of the atomistic system increases beyond hundreds of atoms. Another direction is the development of semi-empirical models, such as the tight-binding model [4], whose accuracy and efficiency is between those of the empirical potentials and DFT. In this paper we pursue a more recent approach based on machine learning.

### 1.1. Machine learning interatomic potentials

Application of machine learning (ML) has recently been put forward as a promising idea that would combine the accuracy of the QM models and the efficiency of the interatomic potentials [5–18]. Such machine-learning interatomic potentials (MLIPs) postulate a partitioning of the interatomic interaction energy into individual contributions of the atoms (and sometimes bonds, bond angles, etc.) and assume a very flexible functional form for such a contribution, making it a function of the positions of the neighboring atoms, typically with hundreds or more parameters. These parameters are found by requiring the energy, forces and/or stresses predicted by a MLIP to be close to those obtained by a QM model on some atomic configurations. These configurations are called the *training set*, and finding the parameters of a MLIP is known as *training* or *fitting*. One of the important features of MLIPs are their ability to approximate potential energy surfaces with arbitrary accuracy (at least theoretically) by increasing the number of parameters and the training set. It should be noted that there are other, ML-based atomistic models of solids, including those predicting the energy directly without partitioning it [19,20], or constructing a density functional in a DFT with machine learning [21]. A recent overview of ML-based models of materials can be found in [22].

Each of the existing MLIPs has a nontrivial functional form accounting for the physical symmetries of interatomic interaction. Namely, a MLIP should be invariant with respect to translation, rotation, and reflection of the space, and also permutation of chemically equivalent atoms. In addition, the potential should have a

* Corresponding author.
*E-mail addresses:* e.podryabinkin@skoltech.ru (E.V. Podryabinkin), a.shapeev@skoltech.ru (A.V. Shapeev).

local support (i.e., depend on surrounding atoms only within a finite cut-off radius) and be smooth with respect to atoms coming and leaving the support. In many instances, it is achieved by designing a fixed number of descriptors [23,24]—scalar functions that satisfy all the symmetries and uniquely encode each atomic environment, and assuming that a MLIP is an arbitrary function (which we call the *regression model*) of these descriptors. This idea was first put forward by Behler and Parrinello [10] proposing an ML model which they called a neural network potential (NNP), based on their descriptors and neural networks as the regression model. Since then, there has been many works on NNPs, see the review papers [8,9] and references therein, and also more recent works [5,12–15,25,26]. Another group of authors adopted the Gaussian process regression framework [7]. They used the coefficients of spherical harmonics expansion of the smeared atomic positions as descriptors and used the kernel-based ML model, where the kernel was based on the distance between the vectors comprised of those coefficients. In a follow-up paper, [17], the authors refined this idea by proposing the smooth overlap of atomic positions kernel, bypassing the step of designing the descriptors. For other examples of using Gaussian process regression for constructing interatomic potentials refer to [6,27]. Three closely related works, [28–30], use Gaussian process regression to predict the forces on atoms directly, without predicting the energy and taking its gradient. Finally, [18] proposes a linear regression model with spherical harmonics coefficients as the basis functions. In the present work, we use the moment tensor potentials (MTPs) [16]. These potentials adopt a linear regression model with polynomial-like functions of atomic coordinates as the basis functions. The MTPs can be interpreted as having descriptors which are based on tensors of inertia of atomic environments.

The MLIPs described above allow for improving their accuracy through increasing the number of the fitting parameters. However, the approximation properties of ML potentials depend not only on their algebraic form, but also on the training set used to fit them. Choosing a good training set for a potential with many parameters (say, more than ten) proves to be a highly nontrivial practical problem. Indeed, all the existing MLIPs are interpolative, they fail to give reasonable answers outside their training domain. Therefore, a good training set should make a MLIP to be interpolative over all the relevant configurations. Obviously, the more parameters a MLIP involves, the larger and more diverse the training set is required in order to fit such a MLIP.

The problem of choosing a proper training set for the fitting of a reliable MLIP is related to the problem of transferability—the ability of interatomic potentials to extrapolate, i.e., give reasonable predictions outside the training domain (e.g., predict the double vacancy formation energy if only single vacancies are present in the training set). It is hardly expected that a MLIP can extrapolate beyond the training domain, but even developing a reliable problem-specific MLIP that would accurately interpolate within the training domain is nontrivial, as pointed out, for instance, by Behler [9, Section 4]. As an illustration of this, the authors of [17] sampled gamma surfaces (by shifting, in different ways, a part of a crystal along a glide plane) and included them in the training set, which allowed them to compute the properties of dislocations accurately with the exception of their Peierls barrier. To accurately reproduce the latter they devised a more complicated scheme of generating configurations from the MD trajectories using one version of their potential in order to fit a better version of their potential.

An attractive idea is to attempt to sample the entire space of atomic environments within, for example, a constraint on the minimal interatomic distance. It is, however, not clear how to do this with sufficient accuracy due to extremely high dimensionality of the space of atomic neighborhoods. Therefore, in practice, the

training set is usually generated by specially designed sampling procedures such as, for example, random perturbations of ideal crystalline configurations [18], sampling from an ab initio MD, or a classical MD with empirical potential or another (already fitted) MLIP [17]. These sampling procedures, however, do not ensure that the training set covers fully, without "gaps", the region in the configuration space required for training MLIPs reliably. In other words, a potential resulted from such a training procedure may later encounter configurations on which this potential will have to extrapolate.

## 1.2. Active learning and learning on the fly

The problem of extrapolation could be resolved if a MLIP were able to detect extrapolative configurations, obtain the QM data for those configurations, and be re-trained. In this scenario, the extrapolation problem (or the transferability problem) would be solved by reliably predicting on the fly whether a potential is extrapolating on a given configuration. Alternatively, in the case when learning on the fly cannot be done, the selection of extrapolative configurations can be done offline yielding the training set that improves the transferability of the fitted potential.

Both scenarios are related to a set of ML techniques called active learning (AL). In contrast to passive learning in which a potential learns every configuration in the training set, in AL a potential is trained only on a set of selected configurations. The key component of any AL method is, thus, its *query strategy*—an algorithmic criterion for deciding whether a given configuration can be treated reliably by an ML model, or we need to re-train our model by querying the QM data for this configuration. If such decision can be made reliably then, as we show in this paper, we do not have to ensure that the training set generated offline has all the representative configurations.

A general overview of AL approaches can be found in [31]. In the context of interatomic potentials, the first work that proposed AL was [32] putting forward a Bayesian query-by-committee strategy. AL was applied by Behler to the neural network potentials [9, Section 4], using the query by committee-type AL strategy. Finally, the authors of [33,34] train a machine learning model predicting the force errors based on the distance between a given atomic configuration and the training set. A very natural AL approach applicable to force fields based on Gaussian process regression [7,17,6,27,29,30], which has not yet been implemented in practice, would be to use the Bayesian predictive variance, shown to correlate with the actual error, e.g., in [21].

In this paper we propose another AL approach for MLIPs based on the D-optimality criterion [31, Section 3.5] allowing for detecting the configurations on which a MLIP extrapolates. This criterion was chosen because there exists an efficient algorithm for checking for D-optimality [35]. Also, as will be discussed in this paper, D-optimality has appealing mathematical interpretations, such as decreasing the uncertainty in determining the parameters or maximizing the volume spanned by the training set in the space of configurations, thus avoiding extrapolation. We apply our AL approach to the fitting of MTPs, however, it is easily generalizable to a any other linear potential, i.e., a potential whose energy depends linearly on the parameters, such as SNAP [18] or GAP. In principle, we can apply AL to atomistic systems with any number of chemically different types of atoms, however, most linearly parametrized potentials developed to date are only applicable to systems with a single type of atoms. We demonstrate that our AL approach allows one to train potentials on the fly with a limited number of QM calculations (occurring, typically, in the initial stage of MD or another atomistic simulation) without loss in accuracy. In addition, we show that even without learning on the fly, AL can "optimize" the training set, in the sense of extracting a significantly smaller