

INTERATOMIC-POTENTIAL-FREE, DATA-DRIVEN MOLECULAR DYNAMICS

J. BULIN¹, J. HAMAEEKERS¹, M. P. ARIZA² AND M. ORTIZ^{3,4}

ABSTRACT. We present a Data-Driven (DD) paradigm that enables molecular dynamics calculations to be performed directly from sampled force-field data such as obtained, e. g., from *ab initio* calculations, thereby eschewing the conventional step of modeling the data by empirical interatomic potentials entirely. The data required by the DD solvers consists of local atomic configurations and corresponding atomic forces and is, therefore, *fundamental*, i. e., it is not beholden to any particular model. The resulting DD solvers, including a fully explicit DD-Verlet algorithm, are provably convergent and exhibit robust convergence with respect to the data in selected test cases. We present an example of application to C_{60} buckminsterfullerenes that showcases the feasibility, range and scope of the DD molecular dynamics paradigm.

1. INTRODUCTION

Present-day computational chemistry methodology, including all-electron calculations, Density Functional Theory (DFT), and other *ab initio* paradigms, enables the generation of vast sets of parameter-free force-field data for complex molecular systems (cf., e. g., [1] and references therein). However, for large systems of atoms molecular dynamics based on empirical interatomic potentials remains the only method of choice. In that setting, the empirical potentials required in the calculations are, inevitably, a main source of empiricism and error. Indeed, at present there does not exist a rigorous theoretical means of generating sequences V_h of approximate interatomic potentials that are guaranteed to converge to the underlying—and unknown—exact potential $V_{ab\ initio}$ in a manner that additionally ensures convergence of trajectories. Instead, empirical potentials based on *ad hoc* assumptions and parametrizations are contrived and fitted to data without global control of errors or convergence guarantees (cf., e. g., [2] for the state of the art of empirical interatomic potentials).

The availability of vast quantities of high-fidelity force-field data (cf., e. g., [3, 4, 5]) suggests a game-changing paradigm shift whereby the data themselves are the sole basis of molecular dynamics calculations and the conventional interatomic-potential modeling step is eschewed altogether. Evidently, such a strict Data-Driven (DD) paradigm, if feasible, would have the immediate beneficial consequences of eliminating the biases, loss of information, empiricism and error that inevitably afflict empirical interatomic potentials.

In this work we present a provably convergent strict DD paradigm for molecular dynamics based solely on force-field data and demonstrate the feasibility of the paradigm with the aid of selected examples. The resulting DD solvers return approximate trajectories that converge to those of the underlying—and unknown—exact potential when computed from force-field data sets of increasing fidelity. Such data sets can be generated *ab initio*, stored, reused, merged and adapted to specific application domains along the tenets of *active learning* [6].

The specific DD paradigm considered in this work builds on similar approaches for static problems [7, 8] and dynamic problems [9]. Specifically, the trajectories of the system are thought to take place in a *phase space* of atomic configurations and force fields, cf. Section 2.2. The objective is then to determine trajectories that satisfy Newton’s laws of motion exactly, viewed as a constraint on forces and accelerations, while remaining as close as possible to a given force-field data set in the sense of some suitable distance, cf. Section 3. We show that the resulting problem for the trajectories can be given the structure of an optimal control problem, cf. Section 3.1, or a game-theoretical problem, cf. Section 3.2. In this latter formulation, we also show that the DD problem can be expressed in terms of an effective, or *learned*, force field. However, the effective force field is only implied and need not be computed explicitly in calculations.

We note, cf. Section 3.4, that a natural distance between local clusters of atoms that is invariant under relabeling of the atoms is provided by the Wasserstein distance of optimal transport [10]. In addition, the force-field data is invariant under the action of the Euclidean group of translations and orthogonal transformation, which requires the evaluation of distances between entire orbits of

the Euclidean group. The calculation of distances between local atomic configurations is one of the main computational bottlenecks of the DD solver and, for large atomic clusters, requires the use of relaxation techniques *à la* Kantorovich [10] and possibly interior-point regularizations such as max-ent [11]. In addition, the DD solver entails frequent searches in large force-field data sets. Here again, efficient search algorithms originally developed for Big Data applications are in existence (cf., e. g., [12] and references therein) and can be deployed as part of the DD solver.

The convergence of the DD solvers with respect to the data, including a fully-explicit DD-Verlet algorithm derived by time discretization, cf. Sections 3.5 and 3.6, can be verified mathematically under simple data sampling scenarios, cf. Section 3.3, or numerically based on selected test cases, cf. Section 4.1. As a proof of concept, we also present a simple application of the DD-Verlet solver to C_{60} buckminsterfullerenes based on synthetic data sampled from the Stillinger-Weber potential [13], cf. Section 4.2. In all these cases, the ability of the DD solvers to compute qualitatively correct trajectories from relatively small data sets is remarkable.

2. CLASSICAL MOLECULAR DYNAMICS

In this section, we define the class of problems under consideration and set forth notational conventions. We begin by considering classical dynamics without any assumptions regarding the structure and properties of the force field such as locality, Euclidean invariance or other symmetries. The consequences of such additional structures are elucidated in subsequent sections. We also concern ourselves with the reformulation of classical dynamics as a variational problem for trajectories in phase space, which sets forth a natural framework for the Data-Driven paradigm developed subsequently.

2.1. General interatomic potentials. We consider a system of N atoms adopting configurations described by coordinates $r_i = (r_i^\alpha)_{\alpha=1}^3$, $i = 1, \dots, N$, in three-dimensional Euclidean space, which we identify with \mathbb{R}^3 equipped with the standard Euclidean metric. For shorthand, we denote by $r \equiv \{r_i\}_{i=1}^N \in \mathbb{R}^{3N}$ the collection of all the position vectors of the atoms in the system and refer to \mathbb{R}^{3N} as the *configuration space*.

The motion of the atoms obeys Newton's second law

$$(2.1a) \quad m_i \ddot{r}_i(t) + f_i(r(t)) = f_i^{\text{ext}}(t),$$

$$(2.1b) \quad r(0) = r_0, \quad \dot{r}(0) = v_0,$$

where t denotes time, m_i is the mass of atom i , $f_i(r)$ is the force on atom i due to atomic interactions, $f_i^{\text{ext}}(t)$ are applied forces, $r_0 \in \mathbb{R}^{3N}$ is the initial configuration, $v_0 \in \mathbb{R}^{3N}$ is the initial velocity field and a superimposed dot denotes time differentiation. The *force field* of the system is the function $f(r) \equiv \{f_i(r)\}_{i=1}^N$ from configuration space \mathbb{R}^{3N} to *force-field space* \mathbb{R}^{3N} .

We note that the configuration space, identified with \mathbb{R}^{3N} , is to be regarded as an affine space of points, whereas the force-field space, while also identified with \mathbb{R}^{3N} , is to be regarded as a vector space. This distinction is consequential when metrizing configuration space, where a distance between sets of points must be defined, cf. Section 3.4¹.

The interaction forces are *conservative* if there exists an interatomic potential $V(r)$ such that

$$(2.2) \quad f_i(r) = \frac{\partial V}{\partial r_i}(r).$$

Examples of classical interatomic potentials commonly used in practice are presented, e. g., in [15]. Recent work aimed at formulating a framework for developing physics-based and machine learning interatomic potentials is exemplified by [16].

2.2. Phase-space reformulation. An alternative set-oriented representation of the force field $f(r)$ is to view it as a graph D in *phase space* $Z = \mathbb{R}^{3N} \times \mathbb{R}^{3N}$, namely, the space of all pairs (r, f) of system configurations and forces². In this representation, the force field is regarded as a material-specific $3N$ -dimensional manifold, or graph, in $6N$ -dimensional phase space Z characterizing the entire range of possible atomic interactions.

¹More general formulations allow for the configuration space to be a smooth manifold with force fields taking values in the corresponding cotangent spaces [14], but that degree of generality is not required here.

²It should be noted that the term *phase space* in classical dynamics is often applied to the space of positions and momenta, whereas in this work we use the term to signify the space of positions and forces.

We shall additionally denote by \mathcal{Z} the linear space of all trial trajectories $z(\cdot) \equiv (r(\cdot), f(\cdot))$ of the system over a given time interval $[0, T]$. Thus, the elements of \mathcal{Z} are curves $z(\cdot)$ in Z parameterized by time, not necessarily satisfying the equations of motion or compatible with the force field of the system.

NB: In order to carefully differentiate between points and trajectories, henceforth we shall denote by r, f and $z = (r, f)$ points in configuration space \mathbb{R}^{3N} , force space \mathbb{R}^{3N} and phase space $Z = \mathbb{R}^{3N} \times \mathbb{R}^{3N}$, respectively; and we shall denote by $r(\cdot), f(\cdot)$ and $z(\cdot) = (r(\cdot), f(\cdot))$ trajectories in the same spaces defined over a given time interval $[0, T]$. In particular, $r(t), f(t)$ and $z(t) = (r(t), f(t))$ denote the values of trajectories $r(\cdot), f(\cdot)$ and $z(\cdot) = (r(\cdot), f(\cdot))$ at time $t \in [0, T]$, respectively.

For a specific material, the phase-space trajectories of the system must take values in the force-field graph D at all times, i. e., they must be contained in the set

$$(2.3) \quad \mathcal{D} = \{(r(\cdot), f(\cdot)) \in \mathcal{Z} : (r(t), f(t)) \in D, t \in [0, T]\}.$$

Thus, \mathcal{D} is the set of phase-space trajectories compatible with the force field of the system and may thus be regarded as a material-specific trajectory set, or *material set* for short. Euclidean invariance requires that $(r, f) \in D$ if and only if $(Q(r - c), Qf) \in D$ for all translations $c \in \mathbb{R}^3$ and orthogonal transformations $Q \in O(3)$. Thus, D must contain entire $E(3)$ -orbits and be invariant under the action of the *Euclidean group* $E(3)$. In addition, we recall that the configuration space consists of point sets and should, therefore, be invariant under relabeling of the atoms, i. e., $(r_i, f_i)_{i=1}^N \in D$ if and only if $(r_{\sigma(i)}, f_{\sigma(i)})_{i=1}^N \in D$ for all permutations σ of the index set $\{1, \dots, N\}$.

The physically admissible trajectories of the system are additionally subject to the constraint set forth by the equations of motion (2.1). The collection of all such admissible trajectories defines the admissible-trajectory set, or *admissible set* for short,

$$(2.4) \quad \mathcal{E} = \{z(\cdot) \equiv (r(\cdot), f(\cdot)) \in \mathcal{Z} : m_i \ddot{r}_i(t) + f_i(t) = f_i^{\text{ext}}(t), t \in [0, T]; r(0) = r_0, \dot{r}(0) = v_0\}.$$

We note that the admissible set \mathcal{E} depends parametrically on the applied force field $f^{\text{ext}}(t)$ and the initial conditions (r_0, v_0) . However, the trial trajectories $(r(\cdot), f(\cdot))$ in \mathcal{E} need not be solutions of the initial-value problem (2.1), since we do not require $r(t)$ and $f(t)$ to be related by the force field $f(r)$ of the material.

Evidently, the actual phase-space trajectories $(r(\cdot), f(\cdot))$ of the system, if they exist, lie in the intersection $\mathcal{D} \cap \mathcal{E}$, i. e., are the admissible trajectories that are consistent with the force field of the material, or, equivalently, the material trajectories that are consistent with the equations of motion and initial conditions.

3. DATA-DRIVEN REFORMULATION

Suppose that, as is often the case in practice, the graph D of the force field is not known in its entirety, but only through an approximating sequence of data sets D_h , $h = 0, 1, \dots$. For instance, the sequence (D_h) may consist of increasing collections of points (r, f) in phase space Z obtained by means of ancillary *ab initio* calculations or by some other means. As noted before, D_h must contain entire $E(3)$ -orbits in order to be invariant under the action of the Euclidean group $E(3)$. In addition, the configurations sampled in D_h are clusters of points and, therefore, their representation must be invariant under relabeling of the atoms.

Example 3.1 (Molecular oxygen O_2). A simple notional example concerns molecular oxygen O_2 , also known as dioxygen or diatomic oxygen, vibrating in a bound configuration without rotation. The ground state of O_2 has a molecular weight of 31.9988, bond length of 1.21 Å and a cohesive energy of -2.58 eV/atom [17]. The oxygen molecule is held together by a strong $\text{O}=\text{O}$ double covalent bond, each oxygen atom sharing two of its outer-shell electrons with the other atom. A sequence of three force-field data sets, containing $\#D_1 = 10$, $\#D_2 = 34$ and $\#D_3 = 100$ points, computed using the Vienna *ab initio* Simulation Package VASP is shown in Fig. 1. The calculations use PAW-PBE as the pseudo-potential and exchange-correlation functional [18, 19]. An energy cutoff of 400 eV is employed for the plane wave basis, and reciprocal space is sampled using a Γ -centered Monkhorst Pack grid [20]. \square

The data sets D_h define approximate material trajectory sets of the form

$$(3.1) \quad \mathcal{D}_h = \{(r_h(\cdot), f_h(\cdot)) \in \mathcal{Z} : (r_h(t), f_h(t)) \in D_h, t \in [0, T]\}.$$

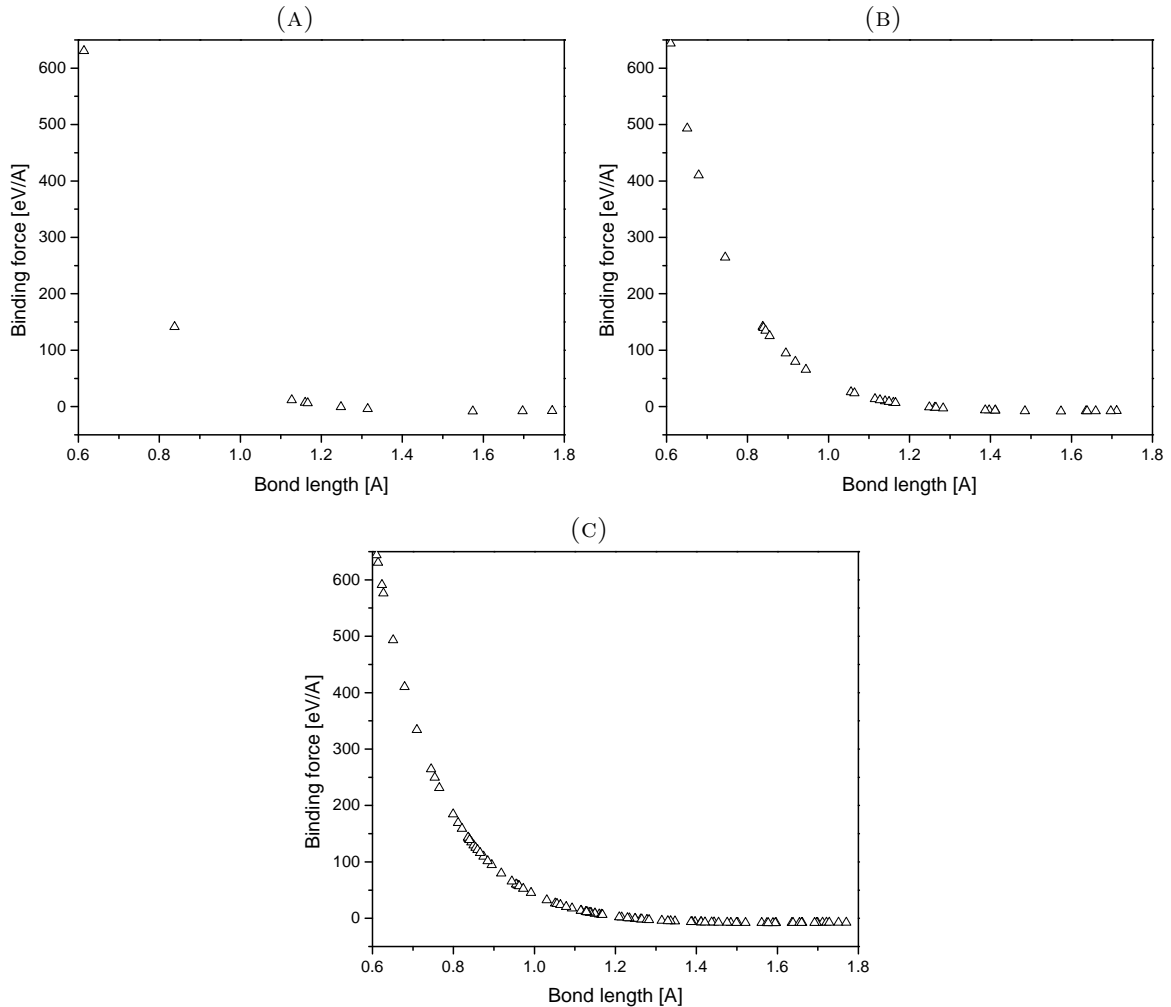


FIGURE 1. Sequence of force-field data sets of molecular oxygen O_2 computed using the Vienna *ab initio* Simulation Package (VASP). a) $\#D_1 = 10$ points; b) $\#D_2 = 34$ points; c) $\#D_3 = 100$ points.

Thus, if D_h is a sample of points in phase space, the corresponding material trajectory set \mathcal{D}_h consists of piecewise trajectories that 'visit' the orbits in D_h in turn.

It is clear that, in general, the intersection between the admissible trajectory set \mathcal{E} and the approximating material trajectory sets \mathcal{D}_h may be empty, in which case no classical approximating solution exists. One way to circumvent this difficulty is to relax the notion of 'approximating solutions' generated from the sequence of data sets D_h . In particular, we would like such Data-Driven solutions to converge to the exact trajectory as the data sets D_h sample the exact force-field graph D with increasing fidelity.

3.1. Optimal-control Data-Driven reformulation. In order to derive approximations directly from a sequence of data sets D_h , we relax the notion of solution in the spirit of [7, 8, 9, 21] to mean a pair $(y_h(\cdot), z_h(\cdot))$ of trajectories in phase space such that $y_h(\cdot) \in \mathcal{D}_h$ is *closest to* \mathcal{E} , and, conversely, $z_h(\cdot) \in \mathcal{E}$ is *closest to* \mathcal{D}_h , in some metric to be defined. Thus, $y_h(\cdot)$ is the material trajectory that is closest to being admissible and $z_h(\cdot)$ is the admissible trajectory that is closest to being material. Formally,

$$(3.2) \quad (y_h(\cdot), z_h(\cdot)) \in \operatorname{argmin} \left\{ \operatorname{dist}(y(\cdot), z(\cdot)) : y(\cdot) \in \mathcal{D}_h, z(\cdot) \in \mathcal{E} \right\},$$

where dist is some suitable distance between phase-space trajectories, to be defined.

We remark that problem (3.2) does not identify, or 'learn', an effective force-field from the data. Instead, it endeavors to minimize discrepancy between admissible solutions, satisfying the equations of motion (2.1), and the data, as measured by a suitably defined distance.

Evidently, if the material data set D_h coincides with the entire force-field graph D the solutions of the Data-Driven problem coincide with the exact classical solutions of the initial-value problem (2.1), as required. Most importantly, the Data-Driven problem continues to make sense—and returns approximating solutions—in the case of partial point data, thus setting forth an approximation scheme for trajectories. We expect the approximate Data-Driven trajectories to converge to exact classical trajectories as the data samples the force field with increasing fidelity.

We note that problem (3.2) has the structure of an *optimal control problem* with differential constraints [22]. To exhibit this structure, suppose, for definiteness, that the distance between trajectories in \mathcal{Z} is of the integral form

$$(3.3) \quad \text{dist}(y(\cdot), z(\cdot)) = \left(\int_0^T \text{dist}^2(y(t), z(t)) dt \right)^{1/2},$$

for some local distance $\text{dist}(y, z)$ over Z , not renamed. Begin by minimizing out the material trajectories $y(\cdot)$ in (3.2) to define a *cost function*

$$(3.4) \quad F_h(z(\cdot)) = \inf_{y(\cdot) \in \mathcal{D}_h} \frac{1}{2} \text{dist}^2(y(\cdot), z(\cdot)) = \int_0^T \frac{1}{2} \text{dist}^2(z(t), D_h) dt,$$

where

$$(3.5) \quad \text{dist}(z, D_h) = \inf_{y \in D_h} \text{dist}(y, z),$$

is the distance from state $z \in Z$ to the data set D_h . Evidently, the cost function $F_h(z(\cdot))$ measures the discrepancy between a trial trajectory $z(\cdot)$ and the force-field data. With these definitions, the abstract Data-Driven problem (3.2) reduces to

$$(3.6a) \quad \text{Minimize: } F_h(z(\cdot)) \text{ in } \mathcal{Z},$$

$$(3.6b) \quad \text{subject to: Eqs. of motion (2.1).}$$

As announced, this problem has the structure of an optimal control problem constrained by ordinary-differential equations if, with $z(\cdot) = (r(\cdot), f(\cdot))$, we identify $f(\cdot)$ as the *control* and $r(\cdot)$ as the *state variable* (cf. [22], Section 3.5). The problem is then to minimize the cost $F(r(\cdot), f(\cdot))$ with respect to the control $f(\cdot)$, with $r(\cdot)$ taken as a solution of the initial-value problem (2.1) for any trial force-field trajectory $f(\cdot)$. General conditions for existence and convergence of solutions for such problems are given in [22], Section 3.5.

In practice, the constraint (3.6b) can be enforced by means of Lagrange multipliers $w(t)$, subject to initial conditions

$$(3.7) \quad w(0) = 0, \quad \dot{w}(0) = 0,$$

resulting in the Lagrangian

$$(3.8) \quad L(y(\cdot), z(\cdot), w(\cdot)) = \frac{1}{2} \text{dist}^2(y(\cdot), z(\cdot)) + \int_0^T \left(m\ddot{r}(t) + f(t) - f^{\text{ext}}(t) \right) w(t) dt,$$

to be rendered stationary. Assuming, for definiteness, an Euclidean distance of the form

$$(3.9) \quad \text{dist}^2(y, z) = \sum_{i=1}^N \left(|r_i - s_i|^2 + \kappa^2 |f_i - g_i|^2 \right),$$

with $y = (s, g)$, $z = (r, f)$, the Euler-Lagrange equations corresponding to variations in $z(\cdot)$ and $w(\cdot)$ are

$$(3.10a) \quad r_h(t) - s_h(t) + m\ddot{w}_h(t) = 0,$$

$$(3.10b) \quad \kappa^2 (f_h(t) - g_h(t)) + w_h(t) = 0,$$

$$(3.10c) \quad m\ddot{r}_h(t) + f_h(t) - f_h^{\text{ext}}(t) = 0.$$

where we write $y_h(\cdot) = (s_h(\cdot), g_h(\cdot))$ and $z_h(\cdot) = (r_h(\cdot), f_h(\cdot))$ for the resulting Data-Driven solutions. In (3.10), $(s_h(t), g_h(t))$ is to be chosen at all times $t \in [0, T]$ as the point in D_h closest to $(r_h(t), f_h(t))$. We note that, by this choice, problem (3.10) is reduced to two coupled second-order ordinary differential equations in time in the unknowns $r_h(\cdot)$ and $w_h(\cdot)$. This duplicate structure is reminiscent of

static Data-Driven problems in which the equilibrium constraint is enforced by means of Lagrange multipliers [7].

Suppose that $D_h = D$, i. e., the data supply a full representation of the force field of the system. Then, in (3.10a) and (3.10b) we have $r_h(t) = s_h(t)$ and $f_h(t) = g_h(t)$, which, together with (3.7) give $w_h(t) = 0$. In addition, (3.10c) reduces to (2.1a), as required. In general, for underlying force-field graphs of sufficient regularity we expect the Lagrange multipliers $w_h(\cdot)$ to tend to zero and the Data-Driven trajectories to converge to exact classical trajectories when the density of sampling increases and the data sets D_h approximate D with increasing fidelity.

3.2. Game-theoretical Data-Driven reformulation. An alternative Data-Driven paradigm that does learn an effective force-field from the data consists of recasting the Data-Driven problem (3.2) as a game-theoretical problem. To this end, we regard the functional $-F_h(r(\cdot), f(\cdot))$, eq. 3.4 with $z(\cdot) = (r(\cdot), f(\cdot))$, as the payoff for the *force player*. Evidently, $-F_h(r(\cdot), f(\cdot))$ which measures the agreement between the trajectory $z(\cdot) = (r(\cdot), f(\cdot))$ and the data set. In this reinterpretation, the *objective* of the force player is to determine a *strategy* $f(\cdot)$ that maximizes its payoff, or, equivalently, minimizes $F_h(r(\cdot), \cdot)$, for given $r(\cdot)$. The force player competes against a second *position player*, whose *objective* is to determine a *strategy* $r(\cdot)$ that satisfies the initial-value problem (2.1) for given $f(\cdot)$. The objective of the position player can be expressed variationally by means of the action functional

$$(3.11) \quad G(r(\cdot), f(\cdot)) = \int_0^T \sum_{i=1}^N \left(\frac{m_i}{2} |\dot{r}_i(t)|^2 - (f_i(t) - f_i^{\text{ext}}(t)) \cdot r_i(t) \right) dt,$$

to be minimized with respect to $r(\cdot)$ at fixed $f(\cdot)$, with initial conditions (2.1b) replaced by the boundary conditions

$$(3.12) \quad r(0) = r_0, \quad r(T) = r_T,$$

with r_0 and r_T given. The minimizing property of the action functional, which normally only attains stationarity in general Lagrangian mechanics, is remarkable and owes to the independence of the force field. The Data-Driven problem then becomes

$$(3.13a) \quad f_h(\cdot) \in \operatorname{argmin} F_h(r_h(\cdot), \cdot),$$

$$(3.13b) \quad r_h(\cdot) \in \operatorname{argmin} G(\cdot, f_h(\cdot)),$$

with defines a non-cooperative game between the force and position players [23].

The difference between the game (3.13) and the optimal control problem (3.2) is subtle but significant. Thus, in the optimal control problem the trial position histories $r(\cdot)$ are tied to the trial control histories $f(\cdot)$ through the initial-value problem (2.1). By virtue of this constraint, the cost functional $F_h(r(\cdot), f(\cdot))$ has a double dependence on the trial control histories $f(\cdot)$, once through its direct dependence and twice implicitly through $r(\cdot)$. By contrast, in the game problem (3.13) the same functional $F_h(r(\cdot), f(\cdot))$ is minimized with respect to $f(\cdot)$ at fixed $r(\cdot)$. Therefore, the Data-Driven approximations $(r_h(\cdot), f_h(\cdot))$ generated by the two procedures are different in general.

Specifically, suppose that we fix a trial trajectory $r(\cdot)$ in configuration space. From (3.4) and (3.13a) it follows that the optimal force strategy $f(\cdot)$ of the force player is given at every time by the *effective force field*

$$(3.14) \quad f(r(t); D_h) = f(t) : (s(t), f(t)) \in D_h, \operatorname{dist}(s(t), r(t)) \rightarrow \min!,$$

which assigns to every system position $r(t)$ the force $f(t)$ such that $(s(t), f(t)) \in D_h$ and $s(t)$ is nearest to $r(t)$. The optimal position strategy $r_h(\cdot)$ then follows by solving the initial-value problem

$$(3.15a) \quad m_i \ddot{r}_i(t) + f_i(r(t); D_h) = f_i^{\text{ext}}(t),$$

$$(3.15b) \quad r(0) = r_0, \quad \dot{r}(0) = v_0,$$

formally identical to (2.1) with the exact, but unknown, force-field $f(r)$ replaced by the 'learned' Data-Driven force field $f(r; D_h)$.

It bears emphasis that the 'learned' Data-Driven force field $f(r; D_h)$ is not the product of modeling, but rather a result of analysis. Piecewise constant approximations of the force field have been analyzed by Gonzalez *et al.* [24], who termed the approximation scheme *force stepping* and analyzed its convergence properties. As in the optimal control formulation, we expect the approximate trajectories $(r_h(\cdot), f_h(\cdot))$ obtained from (3.15) to converge to the exact solution as the data sets D_h approximate

the graph D with increasing fidelity. General conditions for existence and convergence of solutions of general game-theoretical problems are given in [23].

3.3. Convergence with respect to the data. The convergence properties of the Data-Driven solutions $(r_h(\cdot), f_h(\cdot))$ to the exact trajectory $(r(\cdot), f(\cdot))$ as the material data set D_h samples the force-field graph D with increasing fidelity can easily be verified for simple scenarios. The analysis is greatly simplified by assuming long-term existence of trajectories $(r(\cdot), f(\cdot))$ and then focusing on the convergence of approximations $(r_h(\cdot), f_h(\cdot))$ thereof, in keeping with the main focus of this work.

We recall that a function $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is Lipschitz continuous if there is a constant $L > 0$ such that

$$(3.16) \quad \|f(r') - f(r'')\| \leq L \|r' - r''\|.$$

for all $r', r'' \in \mathbb{R}^m$. The space of Lipschitz-continuous functions of time over $[0, T]$ with values in \mathbb{R}^n is denoted $W^{1,\infty}((0, T); \mathbb{R}^n)$; it is a Banach space with norm

$$(3.17) \quad \|r\|_{1,\infty} = \max \left\{ \operatorname{ess\,sup}_{t \in (0, T)} |r(t)|, T \operatorname{ess\,sup}_{t \in (0, T)} |\dot{r}(t)| \right\},$$

where $\operatorname{ess\,sup}$ denotes the essential supremum (cf., e. g., [25]).

Proposition 3.2 (Convergence of optimal-control Data-Driven problem). *Suppose that $f(r)$ is Lipschitz continuous and the exact initial-value problem (2.1) has solutions $r(\cdot)$ in $W^{1,\infty}((0, T); \mathbb{R}^{3N})$. Suppose that the data sets D_h are generated by sampling the exact force field $f(r)$ at points $s \in \mathbb{R}^{3N}$ and that there is $\epsilon_h \downarrow 0$ such that for every $r \in \mathbb{R}^{3N}$ there is $(s, g) \in D_h$ with $|r - s| \leq \epsilon_h$. Then, the optimal-control Data-Driven solutions $r_h(\cdot)$ converge to the exact solution $r(\cdot)$ strongly in $W^{1,\infty}((0, T), \mathbb{R}^{3N})$.*

The assumption of Lipschitz continuity places restrictions on the force field. For instance, for a dimer such as molecular oxygen O_2 , Example 3.1, Lipschitz continuity and inversion symmetry require $f(r)$ to be continuous through the origin $r = 0$ with $f(0) = 0$, i. e., the interaction force must vanish when the two particles merge. We also recall that, if $f(r)$ is differentiable, then the Lipschitz constant L is given by the maximum value of $|Df(r)|$, where $Df(r)$ is the matrix of partial derivatives of $f(r)$, or Hessian. Convergence in $W^{1,\infty}$ means, in particular, that trajectories $r(\cdot)$ in configuration space are Lipschitz continuous and the velocities $\dot{r}(\cdot)$ exist almost everywhere in $[0, T]$ and are essentially bounded.

The convergence of the game-theoretical Data-Driven solutions can be verified likewise.

Proposition 3.3 (Convergence of game-theory Data-Driven problem). *Under the assumptions of Prop. 3.2, the game-theoretical Data-Driven solutions $r_h(\cdot)$ converge to the exact solution $r(\cdot)$ strongly in $W^{1,\infty}((0, T), \mathbb{R}^{3N})$.*

The proofs of the preceding propositions are straightforward and are presented in Appendix A for completeness.

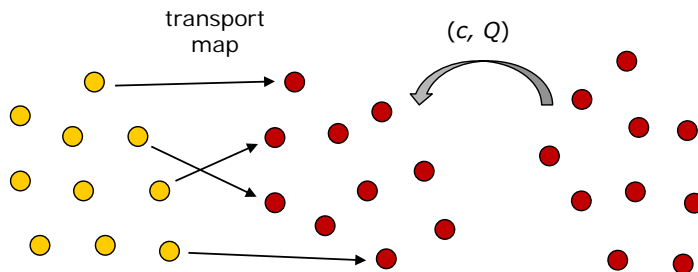


FIGURE 2. Schematic representation of the Wasserstein distance between the $E(n)$ -orbits of two clusters of points (yellow and red).

3.4. The distance between $E(3)$ -orbits of point sets. Recall that the force-field graph D and material data sets D_h thereof are invariant under the action of the Euclidean group $E(3)$, therefore contain entire $E(3)$ -orbits, and are additionally invariant under permutations of the atoms. However, in practice actual material data sets are likely to consist of representatives of the $E(3)$ -orbits pegged

to a particular numbering of the atoms. Details of the calculation of distances between $E(3)$ -orbits of point sets from representatives thereof are presented in this section.

Distances between point sets are naturally measured by the discrete Wasserstein distance [10], i. e., the cost of transportation of one point set to another. Thus, let $r = \{r_i\}_{i=1}^N$ and $s = \{s_i\}_{i=1}^N$ be two clusters of points in configuration space \mathbb{R}^{3N} (e. g., red and yellow in Fig. 2) and let $f = \{f_i\}_{i=1}^N$ and $g = \{g_i\}_{i=1}^N$ be force fields attached to them, respectively. We write $y = (s, g)$ and $z = (r, f)$ for the corresponding points in phase space. For simplicity, we consider clusters with the same number of points, though the optimal transport framework can be extended to clusters of different size [10]. The discrete 2-Wasserstein distance between the $E(3)$ -orbits of the point sets $\cup_{i=1}^N \{y_i\}$ and $\cup_{i=1}^N \{z_i\}$ in $\mathbb{R}^3 \times \mathbb{R}^3$ is

$$(3.18) \quad \text{dist}^2(y, z) = \min_{\sigma \in S_N, (c, Q) \in E(3)} \sum_{i=1}^N \left(|r_i - Q(s_{\sigma(i)} - c)|^2 + \kappa^2 |f_i - Qg_{\sigma(i)}|^2 \right),$$

where the minimum is sought over the group S_N of all permutations σ of the index set $\{1, \dots, N\}$, translations $c \in \mathbb{R}^3$ and orthogonal transformations $Q \in O(3)$, Fig. 2. Fixing $\sigma \in S_N$ and taking variations of (3.18) with respect to (c, Q) gives the Euler-Lagrange equations

$$(3.19a) \quad \sum_{i=1}^N (Q^{-1}r_i - (s_{\sigma(i)} - c)) = 0,$$

$$(3.19b) \quad \sum_{i=1}^N Q^{-1}(r_i \otimes (s_{\sigma(i)} - c) + \kappa^2 f_i \otimes g_{\sigma(i)}) = U = U^T.$$

From (3.19a), we obtain

$$(3.20) \quad c = \frac{1}{N} \sum_{i=1}^N (Q^{-1}r_i - s_{\sigma(i)}),$$

which, inserted into (3.19b), gives

$$(3.21) \quad \sum_{i=1}^N Q^{-1} \left(r_i \otimes s_{\sigma(i)} + r_i \otimes \left(\frac{1}{N} \sum_{j=1}^N s_{\sigma(j)} \right) + \kappa^2 f_i \otimes g_{\sigma(i)} \right) = U = U^T,$$

upon a redefinition of U . Evidently, the orthogonal transformation Q results from the polar decomposition of a 3×3 matrix that is explicit in the data. Once, Q is determined, the translation c follows from (3.20), and the corresponding trial distance as

$$(3.22) \quad \text{dist}_\sigma^2(y, z) = \sum_{i=1}^N \left(|r_i - Q(s_{\sigma(i)} - c)|^2 + \kappa^2 |f_i - Qg_{\sigma(i)}|^2 \right).$$

The sought distance (3.18) then follows by minimizing over the trial distances from all permutations of the indices.

Evidently, the exact evaluation of the distance (3.18) is of combinatorial complexity due to the need to examine all permutations of the atoms and it is suitable only for small clusters. An alternative form of the distance is the Kantorovich reformulation [10]

$$(3.23) \quad \text{dist}^2(r, s) = \inf_{\substack{\mu \in \mathbb{R}^{N \times N} \\ (c, Q) \in E(3)}} \left\{ \sum_{i=1}^N \sum_{j=1}^N \mu_{ij} \left(|r_i - Q(s_j - c)|^2 + \kappa^2 |f_i - Qg_j|^2 \right) : \mu_{ij} \geq 0, \sum_{j=1}^N \mu_{ij} = 1 \right\},$$

which reduces the computation to the solution linear programming problem for the weights μ_{ij} .

Explicit approximations for the weights can further be obtained by recourse to interior point methods [11]. For instance, a maximum-entropy (max-ent) regularization [26] consists of replacing (3.23)

by

$$(3.24) \quad \text{dist}^2(r, s) = \inf_{\substack{\mu \in \mathbb{R}^{N \times N} \\ (c, Q) \in E(3)}} \left\{ \sum_{i=1}^N \sum_{j=1}^N \mu_{ij} (|r_i - Q(s_j - c)|^2 + \kappa^2 |f_i - Qg_j|^2 + \frac{1}{\beta} \log(\mu_{ij})) : \sum_{j=1}^N \mu_{ij} = 1 \right\},$$

which can be readily solved for the approximate weights μ_{ij}^β . A further minimization with respect to (c, Q) , exploiting the optimality of μ_{ij}^β , results in Euler-Lagrange equations identical to (3.19). The 2-Wasserstein distance $\text{dist}(y, z)$ then follows from the regularized, or *thermalized*, distance $\text{dist}_\beta(y, z)$ in the limit of $\beta \rightarrow +\infty$.

We note that the game-theoretical Data-Driven problem requires distances in configuration space only. Such distances follow from the preceding expressions as a special case by formally setting $\kappa = 0$.

3.5. Time discretization. Suppose now that we wish to approximate trajectories at discrete points $t_0 = 0, t_1, \dots, t_n, t_{n+1} = t_n + \tau, \dots, T$, where, for simplicity, we consider a constant time step τ with integer T/τ . The discrete trajectories are then sequences $(r, f) \equiv (r_n, f_n)_{n=0}^{T/\tau}$ of points in phase space. We denote by \mathcal{Z}_τ the space of such discrete phase-space trajectories.

NB: As in the time-continuous case, Section 2.2, in order to carefully differentiate between points and discrete trajectories, henceforth we shall denote by r, f and $z = (r, f)$ trajectories in \mathcal{Z}_τ . In particular, r_n, f_n and $z_n = (r_n, f_n)$ denote the values of trajectories r, f and $z = (r, f)$ at time $t_n = n\tau \in [0, T]$, respectively.

In the time-discrete setting, we may approximate the equations of motion by means of a general difference formula. For definiteness, we consider three-point formulae of the form

$$(3.25) \quad m_i \frac{r_i^{n+1} + r_i^{n-1} - 2r_i^n}{\tau^2} + \alpha_{-1} (f_i^{n-1} - f_i^{\text{ext}}(t_{n-1})) + \alpha_0 (f_i^n - f_i^{\text{ext}}(t_n)) + \alpha_1 (f_i^{n+1} - f_i^{\text{ext}}(t_{n+1})) = 0,$$

with

$$(3.26) \quad \alpha_{-1} + \alpha_0 + \alpha_1 = 1,$$

and denote by \mathcal{E}_τ the set of discrete trajectories satisfying the discrete equations of motion (3.25) and initial conditions, namely,

$$(3.27) \quad \mathcal{E}_\tau = \{z = (r, f) \in \mathcal{Z}_\tau : (3.25), (r_0, r_1) \text{ given}\}.$$

Likewise, we identify the set of discrete material trajectories as

$$(3.28) \quad \mathcal{D}_\tau = \{y = (r, f) \in \mathcal{Z}_\tau : (r_n, f_n) \in D, n = 1, \dots, T/\tau\},$$

and the set of discrete material-data trajectories as

$$(3.29) \quad \mathcal{D}_{h,\tau} = \{y = (r, f) \in \mathcal{Z}_\tau : (r_n, f_n) \in D_h, n = 1, \dots, T/\tau\}.$$

The time-discrete versions of the optimal-control and game-theoretical Data-Driven formulations defined in the foregoing follow now *mutatis mutandi* through the introduction of a sui metric in the space \mathcal{Z}_τ of discrete phase-space trajectories. For instance, we may replace (3.3) by

$$(3.30) \quad \text{dist}_\tau(y, z) = \left(\sum_{n=0}^{T/\tau} \text{dist}^2(y_n, z_n) \tau \right)^{1/2}.$$

Then, the time-discrete version of the optimal-control Data-Driven problem (3.2) is

$$(3.31) \quad (y^{h,\tau}, z^{h,\tau}) \in \text{argmin} \left\{ \text{dist}_\tau(y, z) : y \in \mathcal{D}_{h,\tau}, z \in \mathcal{E}_\tau \right\}.$$

As in the time-continuous case, the constraint set forth by the discrete equations of motion (3.25) can be enforced by means of a discrete Lagrange multiplier w satisfying homogeneous initial conditions. Stationarity then results in two coupled second-order recurrence relations for r and w .

Likewise, a time-discrete version of the game-theoretical Data-Driven problem (3.13) can be set forth by introducing the discrete cost function, cf. eq. (3.4),

$$(3.32) \quad F_{h,\tau}(z(\cdot)) = \inf_{y \in \mathcal{D}_{h,\tau}} \frac{1}{2} \text{dist}_\tau^2(y, z) = \sum_{n=0}^{T/\tau} \frac{1}{2} \text{dist}^2(z_n, D_h) \tau.$$

The objective of the force player is to determine a strategy f that maximizes its payoff, or, equivalently, minimizes $F_{h,\tau}(r, \cdot)$, for given r . The force player competes against a second *position player*, whose objective is to determine a strategy r that satisfies the discrete equations of motion (3.25) for fixed f . The objective of the position player can be expressed variationally by means of the discrete action functional

$$(3.33) \quad G_\tau(r, f) = \sum_{n=0}^{T/\tau} \sum_{i=1}^N \left(\frac{m_i}{2} |\dot{r}_i^n|^2 - (f_i^n - f_i^{\text{ext}}(t_n)) \cdot r_i^n \right) \tau$$

to be minimized with respect to r at fixed f with given r_0 and $r_{T/\tau}$.

The treatment of the time-discrete problems is otherwise identical to that of the time-continuous problems and details are therefore omitted in the interest of brevity.

Assuming that the time discretization scheme is strongly convergent, i. e., time-discrete trajectories converge to the exact time-continuous trajectories strongly in $W^{1,\infty}((0, T), \mathbb{R}^{3N})$ as $\tau \downarrow 0$, the joint convergence of the time-discrete Data-Driven trajectories with respect to data and time discretization follows directly from Props. 3.2 and 3.3 by passing to diagonal sequences.

Algorithm 1 Data-driven Verlet algorithm

Require: Number of atoms N ; r_{n-1} , r_n and f_n ; force-field data set D_h ; applied loads f_n^{ext} . Then:

- i) Compute r_{n+1} from (3.34).
 - iii) Find (s_{n+1}, g_{n+1}) in D_h such that s_{n+1} is closest to r_{n+1} .
 - iii) Set $f_{n+1} = g_{n+1}$.
- Return r_n , r_{n+1} and f_{n+1} .
-

3.6. Data-driven Verlet algorithm. An explicit time-discretization scheme commonly used in practice is *Verlet's algorithm*

$$(3.34) \quad m_i \frac{r_i^{n+1} + r_i^{n-1} - 2r_i^n}{\tau^2} + f_i^n = f_i^{\text{ext}}(t_n),$$

which is a special case of (3.25) with $\alpha_{-1} = \alpha_1 = 0$. It is readily verified that, in this case, the optimal control and game-theoretical Data-Driven problems coincide.

The resulting time-stepping scheme is summarized in Algorithm 1. It bears emphasis that the Data-Driven Verlet algorithm retains the explicit character of the classical Verlet algorithm. The determination of (r_{n+1}, f_{n+1}) requires a search over the force-field data set D_h , which is the main computational bottleneck of the algorithm. Fast algorithms for searching large data sets may be found in [12] and references therein.

3.7. Short-ranged force fields. As already remarked, the computation of distances between $E(3)$ -orbits of point sets is computationally intensive for large systems. Conveniently, many force fields of practical interest are *short-ranged*, which allows distance calculations to be restricted to small clusters defined according to the range of interaction.

Suppose, for simplicity, that the system under consideration is conservative and the attendant global force field $f(r)$ derives from a global potential $V(r)$. Suppose, in addition, that $V(r)$ is short-ranged. In order to exploit this property, we define the local interatomic potentials

$$(3.35) \quad V_i(r) = \int f_i(r) \cdot dr_i = \int \frac{\partial V}{\partial r_i}(r) \cdot dr_i,$$

in terms of indefinite integrals defined modulo additive constants. Thus, $V_i(r)$ is the primitive function of $f_i(r)$ with respect to r_i with the coordinates of all other atoms held constant. We note that, by the definition of the local interatomic potentials, we have

$$(3.36) \quad f_i(r) = \frac{\partial V_i}{\partial r_i}(r),$$

i. e., the force at atom i follows as the derivative of the local interatomic potential V_i with respect to r_i , with the coordinates of all other atoms held constant.

It bears emphasis that definition (3.35) applies to arbitrary interatomic potentials without loss of generality. In addition, from (3.36) it follows that the collection of local interatomic potentials jointly supply the exact global force field of the system and, therefore, it affords an equivalent representation of said global force field.

For short-ranged interatomic potentials, we may expect the number of atoms involved in the definition of each local interatomic potentials to be much smaller than N . Two atoms i and j are *non-interacting* if

$$(3.37) \quad \frac{\partial f_i}{\partial r_j}(r) = \frac{\partial^2 V}{\partial r_i \partial r_j}(r) = 0,$$

and *interacting* otherwise. We denote by $N_i \subset \{1, \dots, N\}$ the subset of atoms that interact with r_i , not including r_i itself, and by $\#N_i$ the number of atoms in N_i . It then follows from (3.36) that $f_i(r)$ depends only on r_i and its cohort N_i of interacting atoms and that $V_i(r)$ can be chosen to have the same dependence up to inconsequential additive terms.

Example 3.4 (Monatomic chain). Consider a monatomic chain consisting of atoms interacting through the harmonic potential

$$(3.38) \quad V(r) = \sum_i \frac{m\omega^2}{2} (r_{i+1} - r_i)^2.$$

The local potentials satisfy the relations

$$(3.39) \quad \frac{\partial V_i}{\partial r_i}(r) = \frac{\partial V}{\partial r_i}(r) = m\omega^2(r_i - r_{i-1}) + m\omega^2(r_i - r_{i+1}).$$

Integrating, we obtain

$$(3.40) \quad V_i(r) = \frac{m\omega^2}{2} (r_i - r_{i-1})^2 + \frac{m\omega^2}{2} (r_{i+1} - r_i)^2,$$

up to inconsequential additive constants. We note that, whereas $V(r)$ depends on the coordinates of all atoms, $V_i(r)$ depends only on the coordinate r_i of the central atom and the coordinates of its nearest-neighbors $N_i = \{i-1, i+1\}$. \square

We proceed to verify that, if the interatomic potential $V(r)$ is invariant under the action of the Euclidean group $E(3)$, then so are the local interatomic potentials $V_i(r)$. To this end, we note from definition (3.35) that, modulo inconsequential additive constants,

$$(3.41) \quad V_i(r) = \int_0^1 \frac{\partial V}{\partial r_i}(\gamma(s)) \cdot \gamma'(s) ds$$

for any curve $\gamma(s)$ joining a reference configuration r_0 to r , i. e., such that $r_0 = \gamma(0)$ and $r = \gamma(1)$. Since, for any $c \in \mathbb{R}^3$ and $Q \in O(3)$, the curve $Q\gamma(s) + c$ joints the reference configuration $Qr_0 + c$ to $Qr + c$, it follows that, modulo inconsequential additive constants,

$$(3.42) \quad V_i(Qr + c) = \int_0^1 \frac{\partial V}{\partial r_i}(Q\gamma(s) + c) \cdot (Q\gamma'(s)) ds.$$

Suppose now that $V(r)$ is invariant under rigid-body mappings. Then,

$$(3.43) \quad V_i(Qr + c) = \int_0^1 \frac{\partial V}{\partial r_i}(\gamma(s)) \cdot \gamma'(s) ds = V_i(r),$$

modulo inconsequential additive constants, which establishes the invariance of $V_i(r)$.

Finally, we note that the global initial-value problem (2.1) may be regarded as a *game* involving N players, one per atom in the system, with strategies $r_i(t)$ whose payoff is to satisfy their corresponding equations of motion (2.1a) and initial conditions. In variational form, each player seeks to render stationary its action functional

$$(3.44) \quad A_i(r_i(\cdot)) = \int_0^T \left(\frac{m_i}{2} |\dot{r}_i(t)|^2 - V_i(r(t)) + f_i^{\text{ext}}(t) \cdot r_i(t) \right) dt,$$

with appropriate boundary conditions. The objective of the i th player in the game is then to determine a strategy $r_i(\cdot)$ that renders $A_i(r_i(\cdot))$ stationary. We note that, by (3.36), the Euler-Lagrange equations of the game so defined are identical to (2.1a), which shows the equivalence between the two paradigms.

Evidently, the Data-Driven reformulation of the global molecular-dynamics problem set forth in the foregoing can be applied, *mutatis mutandi*, to each of the local problems. In order to formulate the local Data-Driven problems, it suffices to supply *local* force-field data sets $D_{h,i}$ collecting corresponding values of $(r_i, r_{|N_i}, f_i)$, i. e., local position r_i of atom i , positions $r_{|N_i}$ of the atoms in the local neighborhood N_i of interaction and corresponding interaction force f_i on atom i . The dimensionality of the local material data sets $D_{h,i}$ is thus much reduced with respect to that of the global data set D_h and, in particular, is independent of the global size N of the system.

4. NUMERICAL EXAMPLES

Finally, we present examples of application aimed at demonstrating the feasibility, range and scope of the Data-Driven molecular dynamics paradigm.

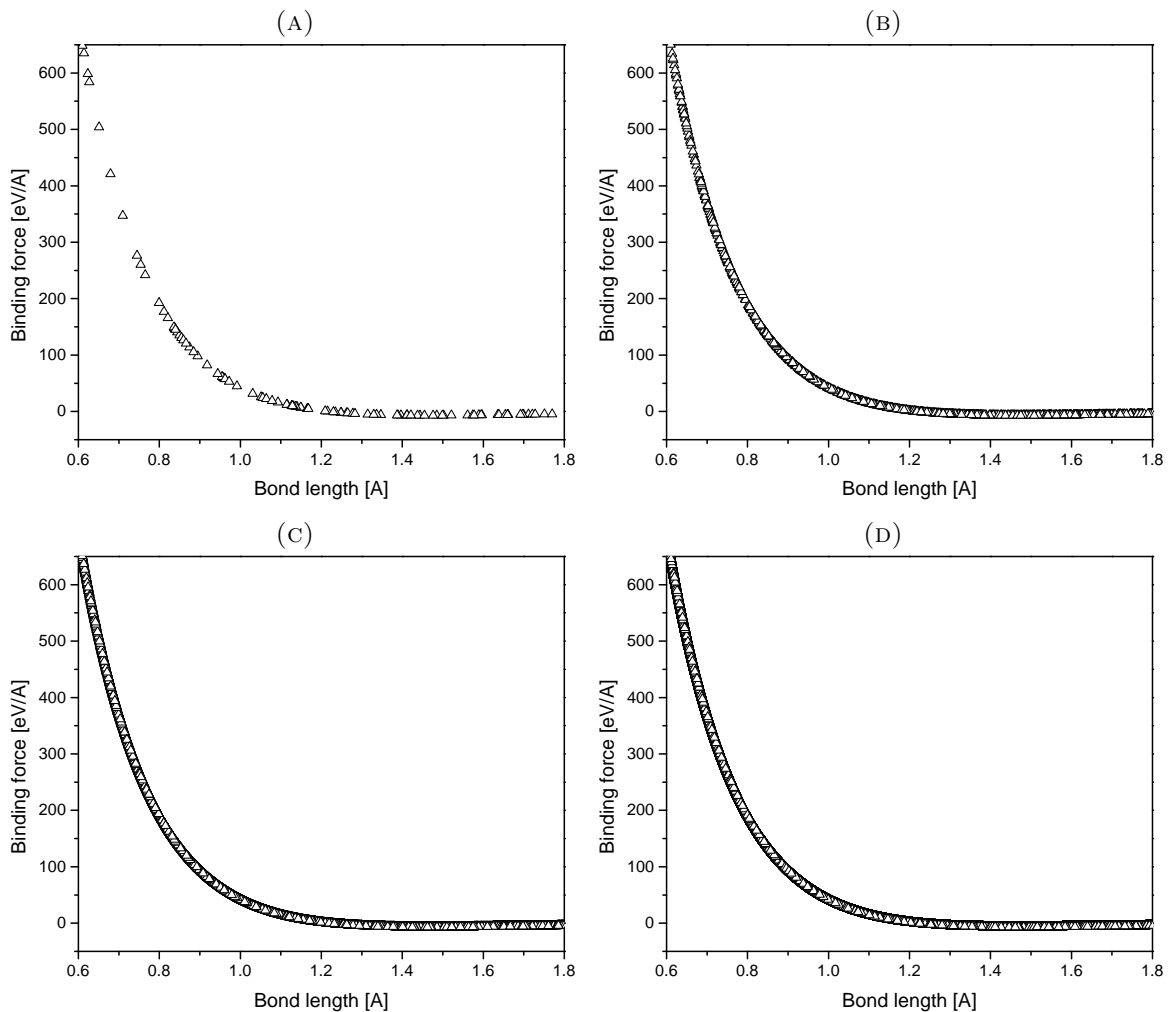


FIGURE 3. Molecular oxygen O_2 . Sequence of data sets sampled from a Morse potential [27] with constants fitted to experimental molecular spectra data [28, 29], cf. Table 1.. a) $\#D_1 = 100$ points; b) $\#D_2 = 1000$ points; c) $\#D_3 = 10000$ points.; d) $\#D_4 = 100000$ points.

4.1. Diatomic oxygen O_2 . The computation of the vibrational spectrum of diatomic oxygen O_2 is an example of a class of problems that arises in chemical physics in connection with the characterization of the molecular spectra of small polyatomic molecules [30]. Because of the complexity of the quantum many-body problem, the accurate calculation of high-resolution vibrational spectra of even small molecules remains a challenge (cf., e. g., [31]). The calculations presented here are elementary and

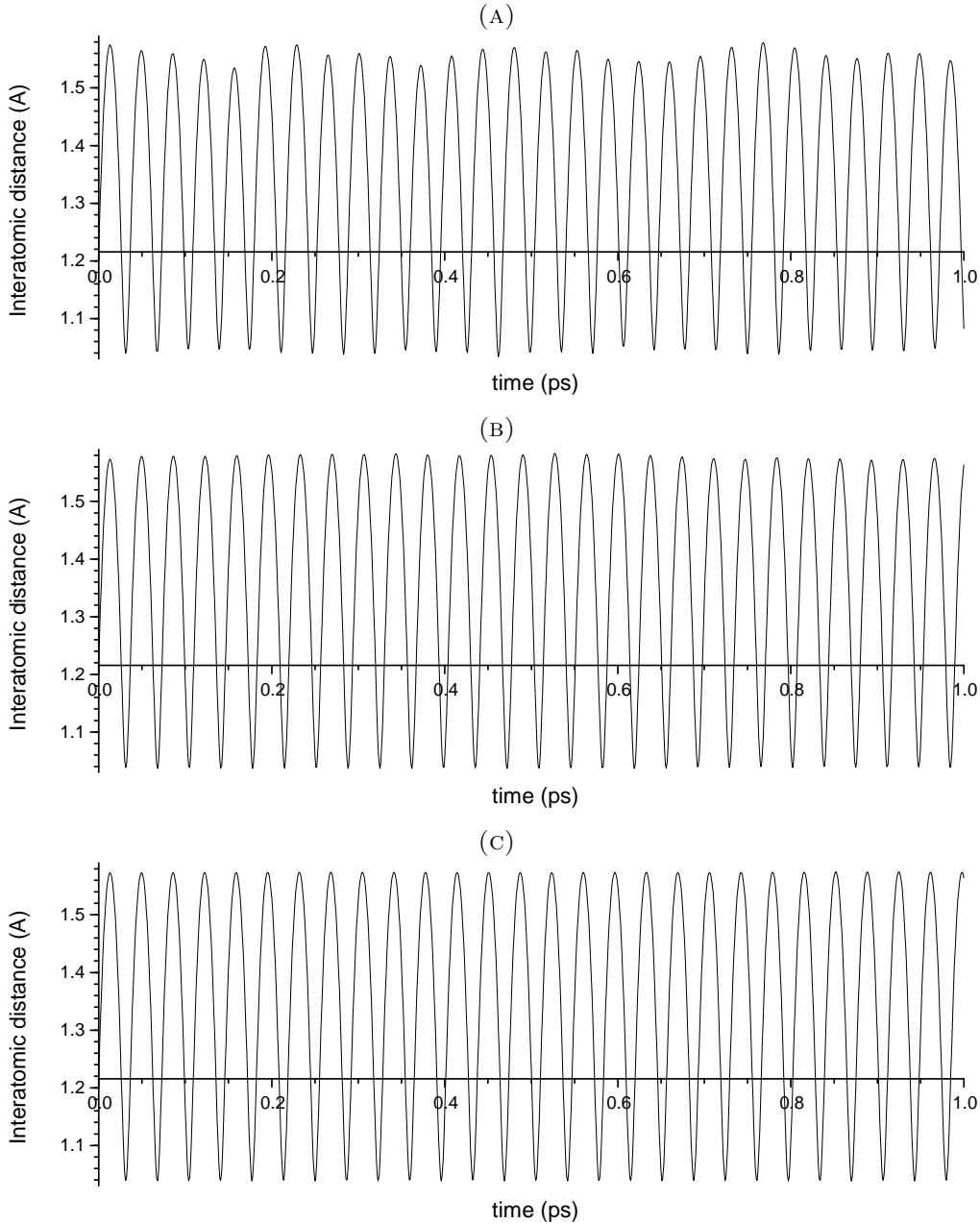


FIGURE 4. Molecular oxygen O_2 . Sequence of DD solutions computed from the data sets in Fig. 3. a) $\#D_1 = 100$ points; b) $\#D_2 = 1000$ points; c) $\#D_3 = 10000$ points.

intended as a simple numerical example aimed at demonstrating the convergence properties of the DD-Verlet algorithm, Section 3.6.

D_e (eV)	r_e (Å)	a (1/Å)
5.12931	1.21560	2.75911

TABLE 1. Molecular oxygen O_2 . Morse potential constants fitted to experimental molecular spectra data [29].

For simplicity, we assume that the O_2 molecule obeys the Morse potential [27]

$$(4.1) \quad V(r) = D_e \left(e^{2a(r_e-r)} - 2e^{a(r_e-r)} \right).$$

with constants fitted to experimental molecular spectra data [28, 29], cf. Table 1. The molecule is assumed to undergo anharmonic vibration resulting, e. g., from a resonance induced by an external

magnetic or optical field. In calculations, the molecule is initially at the equilibrium bond length $r_0 = r_e$ and the atoms are imparted an initial relative velocity $v_0 = 50 \text{ \AA}/\text{ps}$. The time step is set to $\tau = 1 \text{ fs}$ throughout.

Synthetic force-field data sets are generated by sampling the Morse force field derived from (4.1) over the range of bond lengths from $0.5r_e$ to $1.5r_e$. The resulting force-field data sets are shown in Fig. 3. In view of the simplicity of the Morse potential, the agreement with *ab initio* data, Fig. 1, is quite remarkable.

The corresponding DD-Verlet solutions are shown in Fig. 4. The 100-data point solution shows ostensible deviations from the Morse solution. Such deviations notwithstanding, it is remarkable that a qualitatively correct solution, showing the expected smoothness, anharmonicity and periodicity, is obtained from such a small data set. The 1000-data point and subsequent DD-Verlet solutions are ostensibly converged.

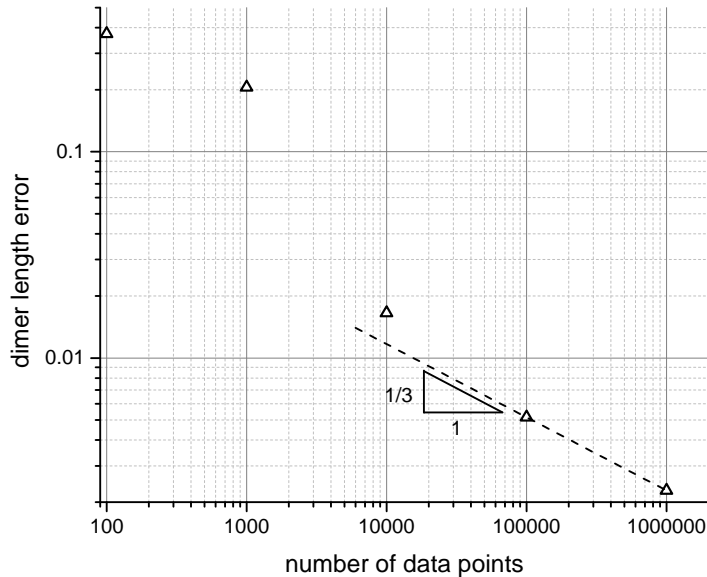


FIGURE 5. Molecular oxygen O_2 . Convergence of DD solutions with respect to the data.

In order to make the convergence analysis quantitative, we measure the discrepancy between the DD-Verlet and Morse solutions using the weighted L^2 -norm [32]

$$(4.2) \quad \|u(\cdot)\| = \left(\int_0^{+\infty} |u(t)|^2 \frac{dt}{t^2} \right)^{1/2}.$$

The norm quantifies both amplitude and frequency errors in infinite wave-train solutions. The simple identities

$$(4.3a) \quad \|A' \sin \omega t - A'' \sin \omega t\| = \frac{\pi \omega}{2} |A' - A''|,$$

$$(4.3b) \quad \|A \sin \omega' t - A \sin \omega'' t\| = \sqrt{\frac{\pi}{2}} |A| |\omega' - \omega''|,$$

exemplify those properties.

Fig. 5 shows a log-log plot of the error incurred by the DD-Verlet solutions with respect to the Morse solutions as a function of the size of the force-field data set. As expected from Props. 3.2 and 3.3, the DD-Verlet solutions exhibit robust convergence towards the solution of the underlying Morse force field, whence the data is sampled. It bears emphasis that the underlying Morse force-field is assumed unknown for purposes of the DD-Verlet calculations. In addition, all DD-Verlet results stem directly from the sampled force-field data and no intermediate interatomic-potential modeling step is required at any stage of the calculations.

4.2. Buckminster fullerene C_{60} . We close with an example concerned with the radiation-induced fragmentation of C_{60} Buckminster fullerene, or *buckyball* for short. The aim of the example is to

demonstrate the ability of the DD-Verlet algorithm to navigate complex conformational transitions using relatively small *ab initio* data sets.

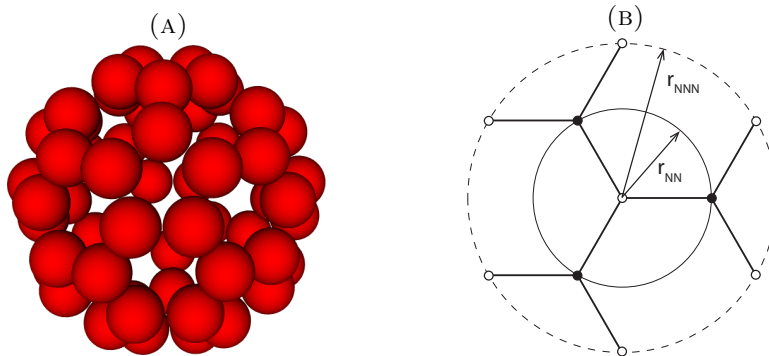


FIGURE 6. a) Relaxed ground-state configuration of the buckminsterfullerene C_{60} using the Stillinger-Weber (SW) [13] potential parametrized as in Table 2. b) Local environment of an atom in the SW ground-state configuration of C_{60} comprising: i) 3 nearest neighbors (black), at an equilibrium distance $r_{NN} < r_c \equiv$ SW cutoff radius, and ii) 6 connected next-to-nearest neighbors (white) at a distance $r_{NNN} > r_c$.

Buckminster fullerene C_{60} , discovered in 1985 by Robert Curl, Harold Kroto, and Richard Smalley [33], is an allotrope of carbon consisting of 60 carbon atoms in a spatial arrangement of positions and bonds in the form of a dodecahedron (32-sided figure), Fig. 6a. The term fullerene more generally refers to compounds C_n consisting of n carbon atoms, each of which is bonded to three other carbons through two single bonds and one double bond, forming a closed surface.

Interactions of intense ultrashort laser pulses with C_{60} fullerene, and the ensuing competition between ionization and fragmentation, have attracted considerable attention [34, 35, 36] as a means of studying the mechanisms of molecular energy deposition. At low intensity irradiation, the vibrational modes of C_{60} are excited, especially its fundamental A_g , or *breathing*, mode [37, 38]. Beyond a certain intensity, C_{60} undergoes *fragmentation* and decomposes into C_n fragments [39, 36].

For purposes of demonstration, we carry out direct and DD-Verlet calculations based on the Stillinger-Weber (SW) potential [13]

$$(4.4) \quad V_i = \sum_{j: r_{ij} \leq r_c} V_2(r_{ij}) + \sum_{k > j: r_{ij} \leq r_c, r_{ik} \leq r_c} V_3(r_{ij}, r_{ik}),$$

where the two-body and three-body interaction energies are chosen of the form

$$(4.5a) \quad V_2(r_{ij}) = A \left[B \left(\frac{|r_{ij}|}{\sigma} \right)^{-4} - 1 \right] \exp \left[\left(\frac{|r_{ij}|}{\sigma} - b \right)^{-1} \right],$$

$$(4.5b) \quad V_3(r_{ij}, r_{ik}) = \lambda \exp \left[\gamma \left(\frac{|r_{ij}|}{\sigma} - b \right)^{-1} + \gamma \left(\frac{|r_{ik}|}{\sigma} - b \right)^{-1} \right] \cos^2(\theta_{ijk} - \kappa),$$

with

$$(4.6) \quad \theta_{ijk} = \arccos \frac{r_{ij} \cdot r_{ik}}{|r_{ij}| |r_{ik}|}$$

representing the angle subtended by the vectors r_{ij} and r_{ik} . In calculations, we use the parametrization of [40], Table 2.

We note that the SW potential has finite range with cutoff radius $r_c = b\sigma = 2.686\text{\AA}$ and an equilibrium bond length $r_e = 1.59169\text{\AA}$. From these values, an examination of the local interaction relations (3.37) for an atom i in the SW ground-state configuration of C_{60} determines the local interaction neighborhood N_i to be of the form shown in Fig. 6b. The local interaction neighborhood comprises: 3 nearest neighbors (black), at an equilibrium distance $r_{NN} < r_c \equiv$ SW cutoff radius, and ii) 6 connected next-to-nearest neighbors (white) at a distance $r_{NNN} > r_c$, for a total $\#N_i = 9$. The actual calculation of the Wasserstein distance and the associated transformations is carried using the comparison algorithm of [41].

Fig. 7 shows the time evolution of the average radius of the C_{60} molecule after irradiation as computed directly from the SW potential. Fig. 7a corresponds to an initial radial velocity $v_0 =$

A (eV)	B	b	λ (eV)	γ	κ	σ (\AA)
5.3790	0.5082	1.8945	18.7079	1.2	-0.5	1.41800

TABLE 2. Stillinger-Weber (SW) potential parameters for carbon [40].

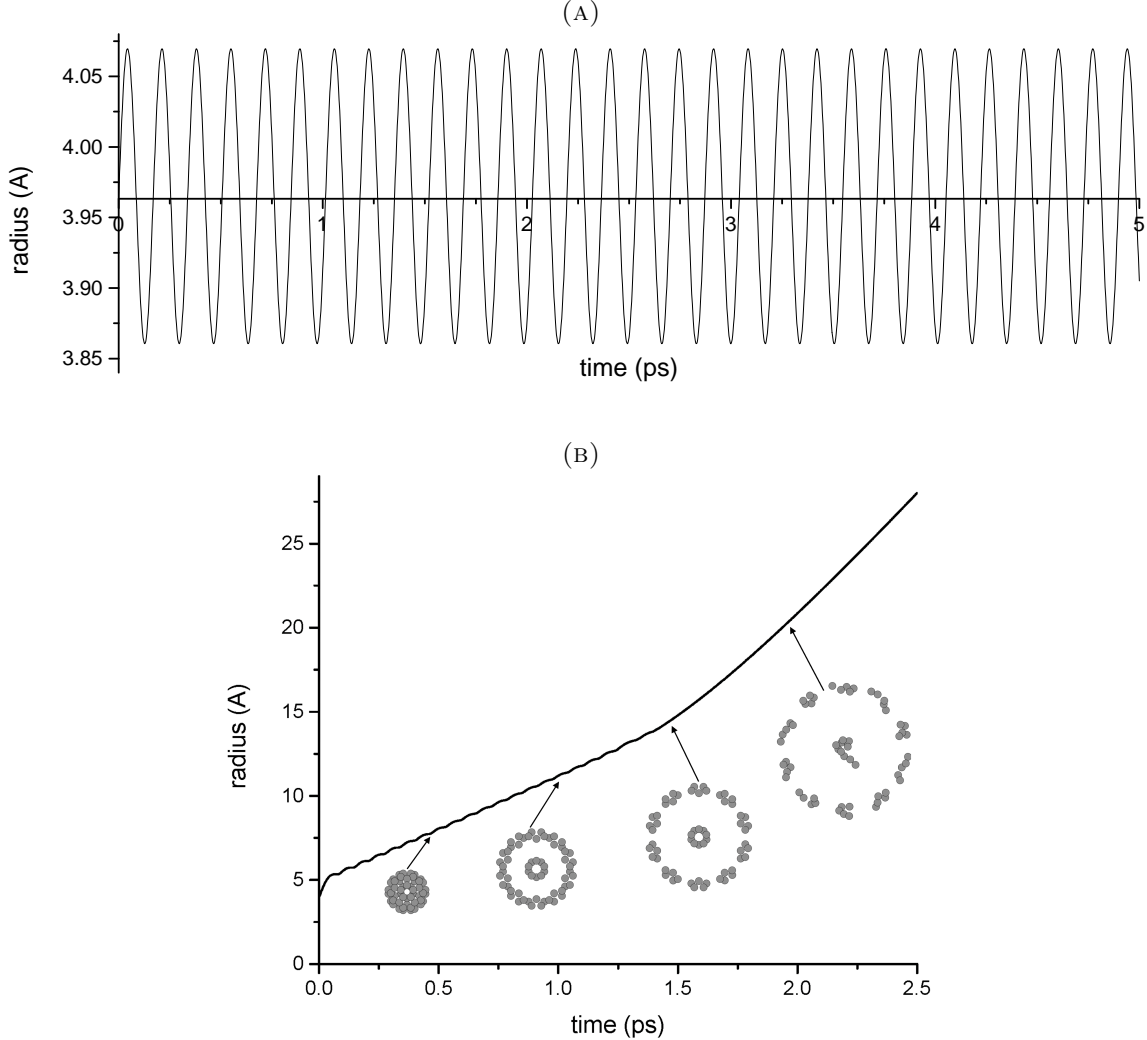


FIGURE 7. Buckminsterfullerene C₆₀. Time evolution of the average radius predicted by the SW potential after: a) low-intensity irradiation; b) high-intensity irradiation. Inset: computed configurations at 0.5ps intervals (plotted to scale).

0.00382Å/fs, below the fragmentation threshold, and Fig. 7b to an initial radial velocity $v_0 = 0.0382\text{\AA}/\text{fs}$, above the fragmentation threshold. In the former case, the fullerene undergoes an anharmonic vibration in the breathing mode, whereas in the latter case the molecule initially dissociates into identical pentarings C₅ that scatter ballistically. A sequence of snapshots showing the scattering of the pentarings following fragmentation is shown inset in Fig. 7b.

We wish to ascertain whether these behaviors can be reproduced by DD-Verlet calculations based on force-field data sets sampled from the SW potential. To this end, we generate synthetic force-field data sets that tabulate local interaction neighborhoods N_i and corresponding forces on the central atom r_i . Since the dimensionality is too large for local configurations to be sampled uniformly, we adopt a *goal-oriented*, or *importance-sampling*, sampling strategy, whereby the goal is to bias the sampling towards local configurations that are likely to arise under the specific conditions of the application of interest. Since the trajectories followed by the C₆₀ molecule after low-intensity irradiation are expected to consist predominantly of radial expansion/contraction, we generate local configurations by randomly expanding/contracting the ground-state configuration, relaxing it and adding random

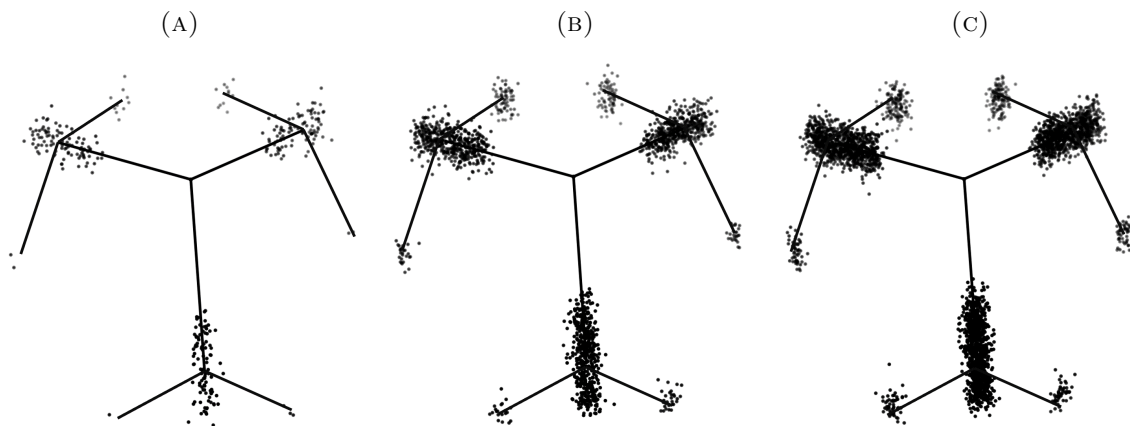


FIGURE 8. Local configuration data sets for assumed cutoff radius $r_c = 2\text{\AA}$. Number of points: a) $\#D_{\text{loc}} = 100$; b) $\#D_{\text{loc}} = 500$; c) $\#D_{\text{loc}} = 1000$.

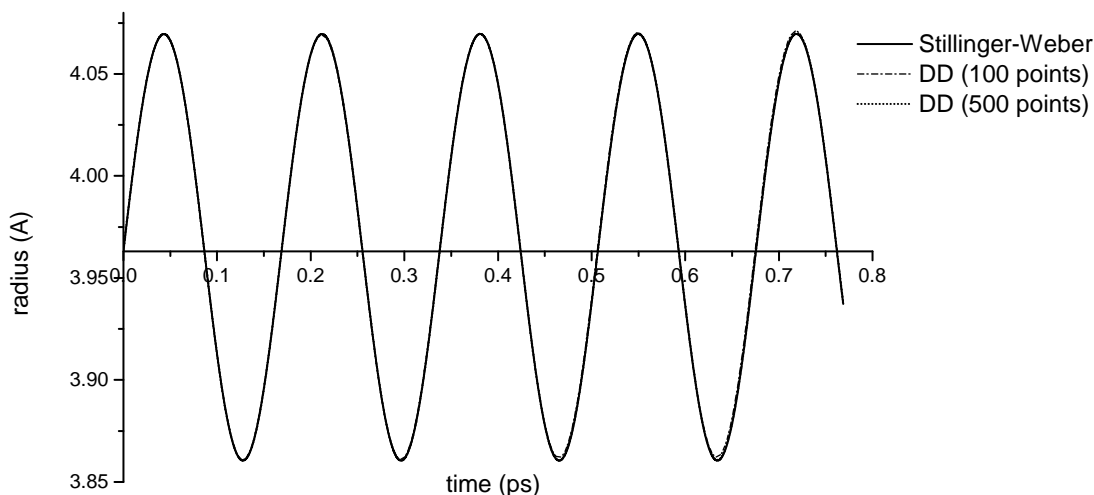


FIGURE 9. DD-Verlet simulation of C_{60} ($r_c = 2\text{\AA}$). Time evolution of the average radius after low-intensity irradiation. Comparison of trajectories obtained directly from the SW potential and from DD-Verlet for local force-field data sets of size 100 and 500.

noise to the positions of the atoms. Since the effective cutoff radius r_c , if any, is presumed unknown *a priori*, we generate data sets over a range of assumed values of r_c . Examples of local configuration sets generated by this procedure assuming a cutoff radius $r_c = 2\text{\AA}$ are shown in Fig. 8.

Fig. 9 shows results of low-intensity irradiation DD-Verlet calculations for data sets generated assuming $r_c = 2\text{\AA}$ compared to the direct calculations. As may be seen from the figure, the DD-Verlet trajectories do closely match the results of the direct SW calculation even for the small data set of size 100.

Generating good force-field data sets for the high-intensity irradiation case *a priori* is challenging, as many different structural changes must be covered by the data in order to reproduce the fragmentation behavior of the fullerene. We may circumvent this difficulty by exploiting the variational structure of the DD-Verlet solver in order to adapt the force-field data sets *on-the-fly*, in the spirit of *active learning* [6]. Thus, we recall that DD-Verlet uses the Wasserstein distance $\text{dist}(r, s)$ from a local configuration r to find similar local configurations s in the force-field data set. Therefore, situations in which the trajectory wanders into poorly sampled regions of phase space are immediately identified, since they result in large distances to the data. To correct for such sampling deficiencies, we perform an explicit Stillinger-Weber calculation whenever in the calculation we encounter a local configuration r such that no local configuration s in the force-field data set satisfies the condition $\text{dist}(r, s) \leq \epsilon$, for some prespecified tolerance ϵ , i. e., if the distance from r to the data set exceeds the tolerance.

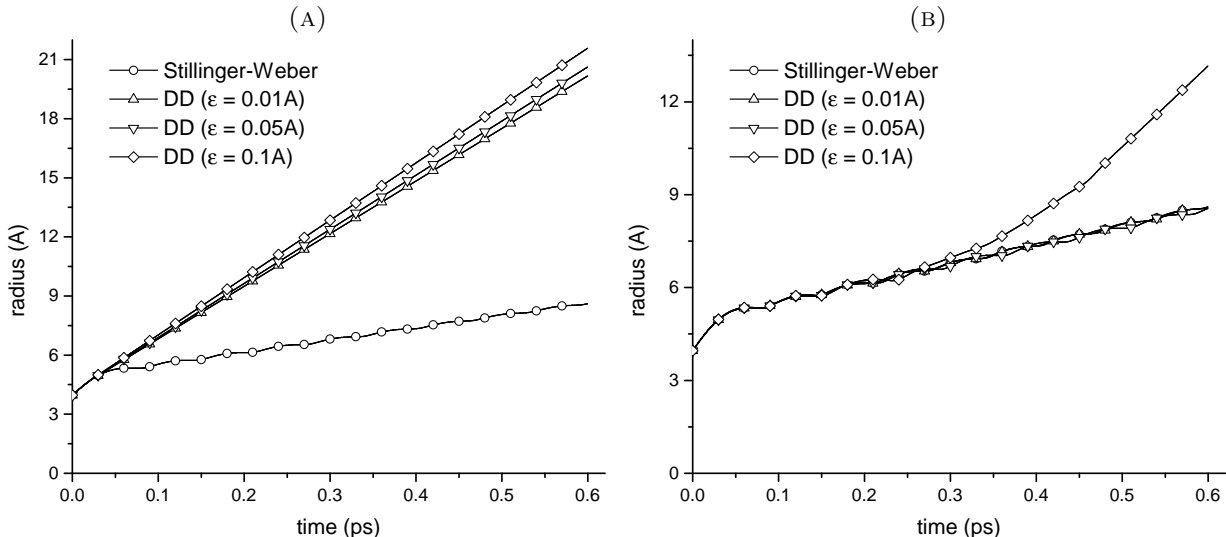


FIGURE 10. DD-Verlet simulation of C_{60} . Time evolution of the average radius after high-intensity irradiation. Comparison of trajectories obtained directly from the SW potential and from DD-Verlet with sampling tolerance $\epsilon = 0.1\text{\AA}$, 0.05\AA and 0.01\AA . a) Local force-field data sampled using a cutoff radius $r_c = 2\text{\AA}$; b) Local force-field data sampled using a cutoff radius $r_c = 4\text{\AA}$.

Stillinger-Weber	DD ($\epsilon = 0.1\text{\AA}$)	DD ($\epsilon = 0.05\text{\AA}$)	DD ($\epsilon = 0.01\text{\AA}$)
6000	269	370	1460

TABLE 3. DD-Verlet simulation of C_{60} . Response up to 0.6ps following high-intensity irradiation. Number of force-evaluations for different values of the distance tolerance ϵ .

As may be seen from Fig. 10a, for a neighborhood cutoff radius of $r_c = 2\text{\AA}$, DD-Verlet exhibits robust convergence with respect to the data-set tolerance (and thus the data-set size) but converges to a trajectory that is at variance with SW. Evidently, in hindsight this behavior is to be expected since the data is generated assuming a cutoff radius that underestimates the actual value $r_c = 2.686\text{\AA}$. This underestimation results in local force-field data sets with poor coverage of the actual SW force-field. By contrast, Fig. 10b collects similar curves obtained from force-field data sets generated assuming a cutoff radius $r_c = 4\text{\AA}$. As may be seen from the figure, in this case the DD-Verlet trajectories converge robustly to the correct SW trajectory as the sampling tolerance ϵ is decreased. Remarkably, convergence is achieved for modest force-field data set sizes, Table 3.

5. CONCLUDING REMARKS

We have presented a paradigm that enables molecular dynamics calculations to be performed directly from sample data, e. g., obtained from *ab initio* calculations, thereby eschewing the conventional modeling of the data by empirical interatomic potentials entirely. The DD paradigm is *lossless*, i. e., it uses the available data in its entirety; *unbiased* in the sense of being free of modeling assumptions, *ad hoc* representational 'features', parametrization and regression choices, and other extraneous elements; and, by its direct and exclusive reliance on data, it is *modeling-error free*.

The data used by the DD paradigm specifically samples local configurations and corresponding atomic forces and is, therefore, *fundamental*, i. e., it is not beholden to any particular model. This strict reliance on fundamental data is potentially far-reaching as it facilitates data generation, management and sharing. Thus, fundamental data is by its very nature *fungible*, i. e., data sets of diverse provenance can be *merged* directly to obtain larger data sets. This property opens the way to the standing of publically-editable, public-domain, force-field data sets that can be built collaboratively worldwide.

Considerations of sampling error, reliability and data provenance can easily be incorporated into the DD paradigm as quality-control measures [8].

The computation of distances between clusters of atoms and the need to execute fast searches in large data sets are among the main computational bottlenecks of the DD paradigm. For large local atomic clusters, it becomes impractical to compute the discrete Wasserstein distance by examining exhaustively all the rearrangements of the atoms, which is an operation of combinatorial complexity. When the sampling is sufficiently dense, the ideas in [41] can be used to limit the combinatorial complexity of the distance calculation to cubic order. In other cases, a Kantorovich relaxation, which reduces the operation to the solution of a linear programming problem, interior-point regularizations thereof, and generally considerations of efficiency in the computation of the distance become of the essence. In addition, fast search algorithms for purposes of Big Data analysis have been extensively investigated and enable the efficient utilization of billion-size data sets (cf., e. g., [12] and references therein).

Evidently, the local force-field data sets used in the DD calculations need to provide adequate coverage of the regions of phase space traversed by the trajectories of interest. However, data sets covering the entirety of phase space uniformly are not possible and the sampling must of necessity be *goal-oriented*, i. e. suited to a particular class of trajectories. For instance, in the application to irradiated C₆₀ fullerenes presented above, the expectation that the molecules undergo predominantly radial expansion/contraction at low irradiation intensities has been used in order to fashion *a priori* goal-oriented data sets delivering fast convergence.

In more general settings, it is not always obvious at the outset whether a given data set is likely to supply adequate coverage. Conveniently, the DD paradigm is *self-correcting* in that regard. Thus, a trajectory that wanders into a poorly sampled region of phase space inevitably results in large distances to the data set being recorded during the calculations, which immediately flags the data as insufficient for the intended purpose. Furthermore, the poorly sampled regions of phase space thus identified can then be sampled by carrying out additional *ab initio* calculations, with the process repeated until the distance between the computed trajectory and the data set is below a prespecified tolerance at all times. The resulting iteration, something referred to to as *active learning* [6], sets forth an approximating scheme with respect to the data itself. If convergent, the scheme results in trajectories that are close to those of the underlying—and unknown—force field sampled by the data, a remarkable feat indeed.

These and other questions and improvements of the DD paradigm presented here suggest themselves as worthwhile directions for further research.

ACKNOWLEDGMENTS

This work has been funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) *via* project 211504053 - SFB 1060; project 441211072 - SPP 2256; and project 390685813 - GZ 2047/1 - HCM. Hausdorff Center for Mathematics (HCM). P. Ariza gratefully acknowledges financial support from the Consejería de Economía y Conocimiento of the Junta de Andalucía, Spain, under grant number P18-RT-1485; and from Ministerio de Ciencia, Innovación y Universidades under grant number RTI2018-094325-B-I00.

APPENDIX A. PROOFS OF PROPOSITIONS 3.2 AND 3.3

Proof. (of Prop. 3.2) For simplicity, we assume that all atoms have identical mass $m_i = m$. We work directly with the Euler-Lagrange equations (3.10). Integrating (3.10a) twice using (3.7) we obtain

$$(A.1) \quad w_h(t) = - \int_0^t \left(\int_0^{t'} \frac{1}{m} (r_h(t'') - s_h(t'')) dt'' \right) dt'.$$

Estimating,

$$(A.2) \quad |w_h(t)| \leq \int_0^t \left(\int_0^{t'} \frac{1}{m} |r_h(t'') - s_h(t'')| dt'' \right) dt'.$$

By the assumption on the density of sampling, we have $|r_h(t) - s_h(t)| \leq \epsilon_h$, whereupon (A.2) gives

$$(A.3) \quad |w_h(t)| \leq \frac{T^2}{m} \epsilon_h,$$

which shows that the Lagrange multiplier $w_h(\cdot)$ tends to zero uniformly over $[0, T]$. Let $r(\cdot)$ be the exact solution of the initial-value problem (2.1). Subtracting equations of motion and initial conditions, we obtain

$$(A.4a) \quad m(\ddot{r}_h(t) - \ddot{r}(t)) + g_h(t) - f(r(t)) = \kappa^{-2}w_h(t).$$

$$(A.4b) \quad r(0) - r_h(0) = 0, \quad \dot{r}(0) - \dot{r}_h(0) = 0.$$

Integrating with respect to time and using the second of (A.4b),

$$(A.5) \quad m(\dot{r}_h(t) - \dot{r}(t)) = \int_0^t (\kappa^{-2}w_h(t') - (g_h(t') - f(r(t')))) dt'.$$

Estimating with the aid of (A.3),

$$(A.6) \quad m|\dot{r}_h(t) - \dot{r}(t)| \leq \int_0^t |g_h(t') - f(r(t'))| dt' + \frac{T^3}{m\kappa^2}\epsilon_h.$$

Furthermore, by the Lipschitz continuity and sampling assumptions,

$$(A.7) \quad \begin{aligned} |g_h(t) - f(r(t))| &= |f(s_h(t)) - f(r(t))| \leq L|s_h(t) - r(t)| \leq \\ L(|r_h(t) - s_h(t)| + |r_h(t) - r(t)|) &\leq L\epsilon_h + L|r_h(t) - r(t)|. \end{aligned}$$

Inserting into (A.6),

$$(A.8) \quad m|\dot{r}_h(t) - \dot{r}(t)| \leq L \int_0^t |r_h(t') - r(t')| dt' + \left(LT + \frac{T^3}{m\kappa^2} \right) \epsilon_h.$$

By Pointcaré's inequality and the first of (A.4b), there is a constant $C > 0$ such that

$$(A.9) \quad m|\dot{r}_h(t) - \dot{r}(t)| \leq CL \int_0^t |\dot{r}_h(t') - \dot{r}(t')| dt' + \left(LT + \frac{T^3}{m\kappa^2} \right) \epsilon_h.$$

Finally, by the integral form of Gronwall's inequality [25],

$$(A.10) \quad |\dot{r}_h(t) - \dot{r}(t)| \leq \left[1 + \frac{CL}{m}t \exp\left(\frac{CL}{m}t\right) \right] \left(\frac{LT}{m} + \frac{T^3}{m^2\kappa^2} \right) \epsilon_h,$$

whence we conclude that the velocity error $|\dot{r}_h(t) - \dot{r}(t)|$ and, together with the first of (A.4b), the position error $|r_h(t) - r(t)|$ converge to zero uniformly over $[0, T]$, as advertised. \square

Proof. (of Prop. 3.3) Again we work directly with the Euler-Lagrange equations (3.15). Subtracting equations of motion and initial conditions, we obtain

$$(A.11a) \quad m(\ddot{r}_h(t) - \ddot{r}(t)) + f(r_h(t), D_h) - f(r(t)) = 0.$$

$$(A.11b) \quad r(0) - r_h(0) = 0, \quad \dot{r}(0) - \dot{r}_h(0) = 0.$$

Integrating with respect to time and using the second of (A.11b),

$$(A.12) \quad m(\dot{r}_h(t) - \dot{r}(t)) = - \int_0^t (f(r_h(t'), D_h) - f(r(t'))) dt'.$$

Estimating,

$$(A.13) \quad m|\dot{r}_h(t) - \dot{r}(t)| \leq \int_0^t |f(r_h(t'), D_h) - f(r(t'))| dt'.$$

By the Lipschitz continuity and sampling assumptions, there is $s_h \in \mathbb{R}^{3N}$ such that

$$(A.14) \quad \begin{aligned} |f(r_h(t), D_h) - f(r(t))| &= |f(s_h(t)) - f(r(t))| \leq L|s_h(t) - r(t)| \leq \\ L(|r_h(t) - s_h(t)| + |r_h(t) - r(t)|) &\leq L\epsilon_h + L|r_h(t) - r(t)|. \end{aligned}$$

Inserting into (A.13),

$$(A.15) \quad m|\dot{r}_h(t) - \dot{r}(t)| \leq L \int_0^t |r_h(t') - r(t')| dt' + LT\epsilon_h.$$

By Pointcaré's inequality and the first of (A.11b), there is a constant $C > 0$ such that

$$(A.16) \quad m|\dot{r}_h(t) - \dot{r}(t)| \leq CL \int_0^t |\dot{r}_h(t') - \dot{r}(t')| dt' + LT\epsilon_h.$$

Finally, by the integral form of Gronwall’s inequality [25],

$$(A.17) \quad |\dot{r}_h(t) - \dot{r}(t)| \leq \left[1 + \frac{CL}{m}t \exp\left(\frac{CL}{m}t\right) \right] \frac{LT}{m} \epsilon_h,$$

whence we conclude that the velocity error $|\dot{r}_h(t) - \dot{r}(t)|$ and, together with the first of (A.11b), the position error $|r_h(t) - r(t)|$ converge to zero uniformly over $[0, T]$, as advertised. \square

REFERENCES

- [1] P. Xu, E. B. Guidez, C. Bertoni, and M. S. Gordon. Perspective: *ab initio* force field methods derived from quantum mechanics. *Journal of Chemical Physics*, 148:090901, 2018.
- [2] The Knowledgebase of Interatomic Models. <https://openkim.org/>.
- [3] The NoMaD Repository. <http://nomad-repository.eu/cms/>.
- [4] The Materials Project. <https://materialsproject.org/>.
- [5] The NIST Materials Genome Initiative. <https://mgi.nist.gov/materials-data-repository/>.
- [6] B. Settles. *Active Learning*, volume 18 of *Synthesis Lectures on Artificial Intelligence and Machine Learning*. Morgan & Claypool, 2012.
- [7] T. Kirchdoerfer and M. Ortiz. Data-Driven computational mechanics. *Computer Methods in Applied Mechanics and Engineering*, 304:81–101, 2016.
- [8] T. Kirchdoerfer and M. Ortiz. Data-Driven computing with noisy material data sets. *Computer Methods in Applied Mechanics and Engineering*, 326:622–641, 2017.
- [9] T. Kirchdoerfer and M. Ortiz. Data-Driven computing in dynamics. *International Journal for Numerical Methods in Engineering*, 113(11):1697–1710, 2018.
- [10] L. C. Evans. *Partial Differential Equations and the Monge-Kantorovich Mass Transfer*. Current Developments in Mathematics. International Press of Boston, Somerville, MA, 1997.
- [11] F. Shu-Cherng and J. H. S. Tsao. *Entropy optimization: interior point methods*, volume 2, pages 544–548. Springer, Boston, MA, 2001.
- [12] R. Eggersmann, L. Stainier, M. Ortiz, and S. Reese. Efficient data structures for model-free Data-Driven computational mechanics. *Computer Methods in Applied Mechanics and Engineering*, 382:113855, 08 2021.
- [13] F. H. Stillinger and T. A. Weber. Computer simulation of local order in condensed phases of silicon. *Physical Review B*, 31:5262–5271, 1985.
- [14] R. Abraham, J. E. Marsden, and T. S. Ratiu. *Manifolds, Tensor Analysis, and Applications*, volume 75 of *Applied Mathematical Sciences*. Springer-Verlag, New York, NY, second edition, 1988.
- [15] M. Finnis. *Interatomic forces in condensed matter*, volume 1 of *Oxford Series on Materials Modelling*. Oxford University Press, New York, NY, 2003.
- [16] Mingjian Wen, Yaser Afshar, Ryan S. Elliott, and Ellad B. Tadmor. Kliff: A framework to develop physics-based and machine learning interatomic potentials. *Computer Physics Communications*, 272:108218, 2022.
- [17] David R. Lide, editor. *CRC Handbook of Chemistry and Physics: A Ready-Reference of Chemical and Physical Data*. CRC Press LLC, Boca Raton, FL, 85th edition, 2005.
- [18] P. E. Blöchl. Projector augmented-wave method. *Physical Review B*, 50:17953–17979, Dec 1994.
- [19] John P. Perdew, Kieron Burke, and Matthias Ernzerhof. Generalized gradient approximation made simple. *Physical Review Letters*, 77:3865–3868, Oct 1996.
- [20] Hendrik J. Monkhorst and James D. Pack. Special points for Brillouin-zone integrations. *Phys. Rev. B*, 13:5188–5192, Jun 1976.
- [21] Sergio Conti, Stefan Müller, and Michael Ortiz. Data-driven problems in elasticity. *Archive for Rational Mechanics and Analysis*, 229(1):79–123, 2018.
- [22] D. Bucur and G. Buttazzo. *Variational methods in some shape optimization problems*. Birkhäuser, Boston, MA, 2005.
- [23] T. Roubicek. *Relaxation in Optimization Theory and Variational Calculus*, volume 4 of *Nonlinear Analysis and Applications*. De Gruyter, Berlin/Boston, second edition, 2020.
- [24] M. Gonzalez, B. Schmidt, and M. Ortiz. Force-stepping integrators in Lagrangian mechanics. *International Journal for Numerical Methods in Engineering*, 84(12):1407–1450, 2010.
- [25] L. C. Evans. *Partial Differential Equations*, volume 19 of *Graduate Studies in Mathematics*. American Mathematical Society, Providence, RI, 1998.
- [26] M. Arroyo and M. Ortiz. Local maximum-entropy approximation schemes: a seamless bridge between finite elements and meshfree methods. *International Journal for Numerical Methods in Engineering*, 65(13), 2006.
- [27] P. M. Morse. Diatomic molecules according to the wave mechanics. ii. Vibrational levels. *Physical Review*, 34(1):57–64, 1929.
- [28] T. J. Zielinski. Exploring the Morse potential. *Journal of Chemical Education*, 75(9):1191, September 1998.
- [29] K. P. Huber and G. H. Herzberg. *NIST Chemistry WebBook*, chapter Constants of Diatomic Molecules. Number 69 in NIST Standard Reference Database. National Institute of Standards and Technology, Gaithersburg MD, 20899, May 2022. data prepared by J. W. Gallagher and R. D. Johnson, III.
- [30] G. Herzberg. *Molecular spectra and molecular structure*, volume 1. van Nostrand, 1950.
- [31] O. B. Gadzhiev, S. K. Ignatov, M. Y. Kulikov, A. M. Feigin, A. G. Razuvaev, P. G. Sennikov, and O. Schrems. Structure, energy, and vibrational frequencies of oxygen allotropes on ($n \leq 6$) in the covalently bound and van der Waals forms: *ab initio* study at the CCSD(T) level. *Journal of Chemical Theory and Computation*, 9:247–262, 2013.

- [32] Ortiz M., E. D. Sotelino, and B. Nour-Omid. Efficiency of group implicit concurrent algorithms for transient finite element analysis. *International Journal for Numerical Methods in Engineering*, 28(12):2761–2776, 1989.
- [33] H. W. Kroto, J. R. Heath, S. C. O'Brien, R. F. Curl, and R. E. Smalley. C60: Buckminsterfullerene. *Nature*, 318(6042):162–163, 1985.
- [34] I. V. Hertel, T. Laarmann, and C. P. Schulz. Ultrafast excitation, ionization, and fragmentation of C60. *Advances In Atomic, Molecular, and Optical Physics*, 50:219–286, 2005.
- [35] J. Kou, V. Zhakhovskii, S. Sakabe, K. Nishihara, S. Shimizu, S. Kawato, M. Hashida, K. Shimizu, S. Bulanov, Ya. Izawa, et al. Anisotropic coulomb explosion of C60 irradiated with a high-intensity femtosecond laser pulse. *The Journal of Chemical Physics*, 112(11):5012–5020, 2000.
- [36] Z. Z. Lin and X. Chen. Ultrafast dynamics and fragmentation of C60 in intense laser pulses. *Physics Letters A*, 377:797–800, 2013.
- [37] W. Krätschmer, L. D. Lamb, K. Fostiropoulos, and D. R. Huffman. Solid C60: a new form of carbon. *Nature*, 347(6291):354–358, 1990.
- [38] R. Meilunas, R. P. H. Chang, S. Z. Liu, M. Jensen, and M. Kappes. Infrared and Raman spectra of C60 and C70 solid films at room temperature. *Journal of Applied Physics*, 70(9):5128–5130, 1991.
- [39] S. C. O'Brien, J. R. Heath, R. F. Curl, and R. E. Smalley. Photophysics of buckminsterfullerene and other carbon cluster ions. *Journal of Chemical Physics*, 88(1):220–230, 1988.
- [40] F. F. Abraham and I. P. Batra. Theoretical interpretation of atomic-force microscope images of graphite. *Surface Science Letters*, 209:L125–L129, 1989.
- [41] Johannes Bulin and Jan Hamaekers. Similarity of particle systems using an invariant root mean square deviation measure, 2021.

¹FRAUNHOFER INSTITUTE FOR ALGORITHMS AND SCIENTIFIC COMPUTING, SCHLOSS BIRLINGHOVEN, 53757 SANKT AUGUSTIN, GERMANY,

²ESCUELA TÉCNICA SUPERIOR DE INGENIERÍA, UNIVERSIDAD DE SEVILLA, SEVILLA 41092, SPAIN,

³HAUSDORFF CENTER FOR MATHEMATICS, UNIVERSITÄT BONN, ENDENICHER ALLEE 60, 53115 BONN, GERMANY,

⁴DIVISION OF ENGINEERING AND APPLIED SCIENCE, CALIFORNIA INSTITUTE OF TECHNOLOGY, 1200 E. CALIFORNIA BLVD., PASADENA, CA 91125.