

## INTRODUCTION

We are excited by the study section's enthusiasm for our Focused Technology Research & Development R01 (PAR-19-253) proposal, and our "very accomplished team with excellence in all areas of hardware design, software engineering, and molecular simulation methodology needed to advance the objectives" [Summary Statement].

**Quantum machine learning (QML) potentials will be a critical part of the future of biomolecular modeling, but only if we fund the tools needed to build and apply them.** The primary source of skepticism expressed in the critiques was whether there is a significant role for QML potentials in the next five years. We have updated the proposal with compelling evidence of the growing importance of QML, as well as evidence of our original work in this important area that will be enabled by the Focused Technology R&D this proposal mechanism supports:

Preliminary data for hybrid QML/MM potentials: Hybrid QML/MM potentials—where part of the system (e.g. ligand or reactive center) is treated with QML and the rest with MM—will be a key driver of near-term applications and can be accelerated to near-MM speeds. Because they only treat a subsystem with quantum chemical accuracy, they do not suffer from any of the issues raised in the critiques. **We present the first protein-ligand alchemical free energy calculations using hybrid QML/MM potentials** (where only the ligand is treated with QML) (**Figure 2**) and **demonstrate the error in binding free energies can be reduced from 1 kcal/mol (MM) to 0.5 kcal/mol (QML/MM) [1]—sufficient for transformative impact on drug discovery.**

Preliminary data for full QML potentials: While it will take several years of work—supported by our infrastructure—for QML potentials to reliably and efficiently model entire biomolecular systems in a manner that satisfyingly includes long-range physical interactions, we demonstrate that alchemical free energy calculations can easily be applied to systems fully treated with QML potentials (**Figure 3**) [2]. Furthermore, we show these QML potentials can be fit to experimental free energies and other properties (in addition to quantum chemical data) to significantly improve their ability to generalize from limited data [2].

Publications: This field is growing rapidly: 7,030 new papers have been published in this area in 2020 alone [3].

We cannot comprehensively review this literature, but we cite several key papers and reviews.

User demand: A recent survey [4] shows  $84\pm 3\%$  of our users would like to use QML potentials in their work, but only  $14\pm 3\%$  use them now due to technical limitations this proposal specifically addresses ( $n = 166$  respondents).

Letters of Support: This resubmission includes 38 external Letters of Support, with the vast majority articulating strong interest in our proposed R&D efforts to advance the training, assessment, and use of QML potentials.

**Machine learning of collective variables (CVs) and integrators:** The critiques articulated concerns about how meaningful collective variables should be constructed, which we wholeheartedly agree with. We now emphasize that we aim only to better support the large communities developing collective variable based sampling methods (e.g. metadynamics; 14,400 citations [5]) and advanced integrators (38,000 citations [5]) by enabling these functions to be defined in convenient machine learning frameworks (PyTorch, TensorFlow, JAX) and easily used within OpenMM. Our proposed work focuses on *enabling* research in this area by providing the means to construct CVs in machine learning frameworks and apply them on the GPU within the OpenMM ecosystem.

**Increased capacity for generating new quantum chemical data:** Since the original submission, the capacity of Folding@home for generating quantum chemical data for QCArchive to accelerating R&D of QML potentials has increased to over 1M actively computing CPU cores available for computation [6].

**Clarity:** We have attempted to improve clarity in all aspects of the proposal following feedback in the critiques.

**Deployment and user support:** OpenMM has been downloaded >385,000 times from Anaconda cloud (one of several distribution mechanisms), +115,000 times since our first submission. This resubmission includes 38 external Letters of Support from group leaders, hardware vendors, and biotech companies representing hundreds of researchers who find the free and open source licensing of OpenMM a significantly enabling technology in their work. Since the original submission, OpenMM has also been integral to powering the worlds first exaFLOP/s distributed computing infrastructure [7] in the hunt for new therapeutics for COVID-19.

**CZI Award:** Thanks to a one-year Chan Zuckerberg Initiative Essential Open Source Software for Science award (\$180K total costs, not renewable), the lead developer of OpenMM was able to extend his contract until early 2021.

**Concerns about the OpenMM "business model":** Following discussions with the NIH Program Officer regarding a puzzling critique that "The business model of OpenMM needs to be improved", we remind the study section that *discussions of whether a proposal merits NIH funding are not germane to the role of the study section, which must evaluate proposals on the basis of the significance and impact of the proposed work and whether the proposed budget is appropriate for the work proposed.* We note, however, that the NIH currently supports the development of legacy codes such as AMBER, CHARMM, GROMACS, NAMD, and Tinker [8], some of which are decades old.

## SPECIFIC AIMS

Predictive, quantitative modeling of biomolecular systems is critical to understanding fundamental molecular mechanisms underlying biological function and disease, as well as developing small molecule and biological therapeutics. Two distinct approaches have emerged to model the chemical interactions in these simulations: **molecular mechanics (MM) models**, in which physical interactions are modeled at the atomic level, and **machine learning (ML) models**, in which a flexible function approximator fit to data makes quantitative predictions. While these approaches have traditionally been disjoint, recent work to bring these two together to develop **quantum machine learning (QML) potentials** has demonstrated the enormous potential for differentiable many-body potentials learned by deep neural networks to substitute for or supplement traditional physical force fields or expensive quantum chemical calculations. Despite their enormous promise to power a new decade of discovery in biomolecular modeling and simulations, a host of technical limitations stand in the way of practical widespread use.

Here, we propose to both sustain and significantly extend our popular **OpenMM Toolkit** to realize this goal. OpenMM is the most widely-used GPU-accelerated framework for biomolecular simulation (>1300 citations, >385,000 downloads, >1M deployed instances). First, we will transition OpenMM to a distributed development and governance model that will ensure it can continue to meet the needs of the biomolecular simulation community for years to come. Second, we add flexible support for plug-in machine learning models to enable multiple uses: (1) as quantum machine learning (QML) potentials able to substitute for traditional MM potentials while delivering extraordinary accuracy in new and existing applications, especially in hybrid QML/MM contexts (where a subsystem, such as a ligand or reactive center, is treated with QML and the rest with MM); (2) as collective variables (CVs) or studying rare events or conformational changes; and (3) as advanced integrators or samplers for simulating long timescales or multiscale dynamics. These plug-ins will be hardware accelerated and fully supported by the existing OpenMM ecosystem, which includes numerous tools for enhanced sampling, protein-ligand free energy calculations, and rare event simulations. Third, to accelerate development of next-generation hybrid QML models that blend traditional QML approaches with physical long-range potentials, we will develop an open portable library of hardware-accelerated atomic features and OpenMM physical potentials within popular machine learning frameworks (TensorFlow, PyTorch, JAX). Together with continued optimizations to exploit accelerated hardware architectures (TPUs, tensor cores, etc), these technological advances will enable a variety of new scientific investigations that were previously computationally intractable or limited by force field accuracy.

**Aim 1: Transition OpenMM from a single-laboratory single-developer-led code to a community governance and distributed developer model to ensure long-term sustainability.** Originally developed in a single laboratory, OpenMM has grown so quickly and is integrated so deeply within the biomolecular simulation community that a distributed development model is necessary to ensure a sustainable future. A new OpenMM Consortium will ensure OpenMM continues to meet the most urgent needs for the biomolecular modeling and simulation community.

**Aim 2: Develop a next-generation hardware-accelerated molecular simulation platform that will enable direct interoperability with machine learning frameworks.** We will extend OpenMM to allow facile integration of QML potentials and machine learning models created within popular frameworks such as TensorFlow, PyTorch, and JAX for use as new forces, collective variables, or integrators. In addition, we will develop custom hardware-accelerated kernels to enable these models to run at high speed within OpenMM by exploiting accelerator hardware such as GPUs, TPUs, and tensor cores. Driving applications include the use of hybrid physical/ML potentials to develop ultra-high accuracy protein-ligand free energy calculations generalizable to new chemistries and the study of enzyme mechanisms or covalent inhibition with enormous speed advantages over traditional QM/MM methods.

**Aim 3: Create a software infrastructure to enable hybrid machine learning / physical model development.** We will develop a library of hardware-accelerated kernels for QML atomic features and physical force terms to be easily used within popular machine learning ecosystems (such as TensorFlow and PyTorch) as custom operations ("ops"). This will enable developers to rapidly develop new QML potentials and hybrid physical/ML models able to deliver quantum chemical accuracy at molecular mechanics speeds, as well as enable high-throughput machine learning for structure-enabled biomolecular datasets in a manner that exploits both physical and ML features.

**Aim 4: Develop community repositories of large quantum chemical datasets and models to accelerate development and deployment of QML potentials.** We will work with MolSSI to create an online QML model repository and standardized QML model specification to allow developers to easily share, update, and deploy their models to OpenMM and other codes. To accelerate the development of QML models with increased accuracy and greatly expanded coverage, we will use the planetary-scale Folding@home distributed computing framework (>1M CPU cores) to create large public quantum chemical datasets made available through the MolSSI QCArchive.

Together, these aims will ensure OpenMM will significantly advance the capabilities of predictive modeling and simulation of biomolecular systems and produce new insight into the modeling of small molecule therapeutics.

## SIGNIFICANCE

**Biomolecular modeling and simulation is a key technology for leveraging the \$16B global investment in structural data in the protein databank.** Molecular mechanics (MM) force fields model biomolecular systems in atomistic detail, allowing researchers to probe a variety of phenomena that involve energetics, forces, and dynamics of biomolecular systems [9–11]. A wide variety of methodologies have been developed for making use of molecular mechanics models, from simple methods that use energy minimization or molecular dynamics to advanced calculations that involve algorithms for enhanced or targeted sampling [9]. Biomolecular modeling and simulation has a wide variety of applications, from probing the molecular mechanism of biological function [10, 12] and the molecular underpinnings of disease [13–15] to the interpretation of experiments [16] to quantitative predictions that guide the development of new small molecule therapeutics [17–19] and biologics [20,21]. Because these investigations only require access to now-ubiquitous computational resources to access microsecond timescales [22], biomolecular modeling has become a facile and indispensable tool in the biosciences.

**OpenMM is a critical technology for the biomolecular simulation and modeling community.** OpenMM [23–28] is a biomolecular modeling and simulation code that is both fast and flexible, enabling researchers to use it directly as a research tool for accessing state-of-the-art algorithms on essentially any computing hardware or as a library to build more complex modeling tools that take advantage of these algorithms in new and advanced ways. OpenMM's key papers [23–28] have been cited over 1300 times<sup>1</sup>. The OpenMM conda package (one of multiple distribution routes) has been downloaded over 385,000 times [31]. OpenMM has been deployed on more than one million GPUs [6] and powered the world's first exascale distributed computing platform [7].

While relatively young, OpenMM provides a number of critical modern features not present in other biomolecular modeling and simulation packages that endow it with numerous advantages compared to legacy biomolecular modeling and simulation packages for a number of important application domains:

- **Speed:** OpenMM was one of the first<sup>2</sup> molecular simulation packages to exploit inexpensive consumer-grade graphics processors (GPUs) [26], which brought a 100× increase in performance/price ratio over CPUs and drove a revolution in biomolecular simulations. OpenMM features both CUDA and OpenCL computational platform backends that allow it to achieve extremely high performance on a wide variety of inexpensive consumer-grade GPUs, a multithreaded CPU backend to enable it to run on ubiquitous CPUs, and a plugin architecture that speeds the development of new kernels to exploit new accelerated hardware.
- **Productivity:** OpenMM provides a simple, expressive API [33] that makes it incredibly easy to build new applications that make use of molecular modeling operations in Python, C++, C, or Fortran without the need to wrap command-line tools. By using an object-oriented model that provides useful abstractions for biomolecular modeling, biosciences researchers can achieve high productivity just as modern machine learning researchers can with modern frameworks such as Tensorflow, PyTorch, and JAX.
- **Python friendly:** The Python API allows scientists to either use OpenMM directly or build applications or libraries in Python, such as the openmmtools library [34] (downloaded >150,000 times [35]). The enormous number of scientific libraries available for Python makes it easy to quickly build new, complex applications or research projects in Python. Educational tutorials can also exploit modern Jupyter notebooks and run in the cloud to make it easy for new researchers to get started.
- **Ubiquity:** OpenMM is conda-installable [31], allowing it to be easily used as a dependency in other tools that can be automatically installed without the user even knowing they are using OpenMM under the hood.
- **Modularity:** Developers and users can easily add new forces, barostats, thermostats, collective variables, and integrators without modifying the core code (and often in pure Python!). New features and kernels can easily be added via a simple plugin architecture that can be built and distributed separately.
- **Extensibility:** Research and experimentation with new potential forms is incredibly easy with Custom Forces, which enable researchers to write algebraic expressions that are automatically optimized, differentiated, and compiled into GPU-accelerated kernels.
- **Interoperability:** OpenMM is force field agnostic, supporting all major molecular mechanics (MM) models. It provides a Python layer to make it easy to import biomolecular systems from other major simulation packages (such as AMBER [36], CHARMM [29], NAMD [37], or GROMACS [38]).

<sup>1</sup>The number of OpenMM citations is likely an underestimate of works using OpenMM because multiple legacy codes—such as CHARMM [29] and TINKER-OpenMM [30]—make use of OpenMM "under the hood", so users may not even know they are using OpenMM.

<sup>2</sup>The other major GPU-accelerated code, ACEMD [32]—created by co-I Fabritius—now uses OpenMM under the hood.

- **Legacy support:** OpenMM is the “webkit for biomolecular simulation,” in that it is ubiquitously found powering molecular simulation packages—with C++, C, and Fortran APIs, OpenMM has been used under the hood in major legacy codes (including CHARMM, Tinker, and ACEMD) to enable these packages to provide high performance on modern computing hardware like GPUs without a complete rewrite.
- **Robustness:** OpenMM drives the GPU simulation core in the Folding@home worldwide distributed computing project, and has run on all major commercially-available GPUs and **over a million unique computer systems**, making it the most widely-tested and widely-deployed biomolecular simulation package. In Mar–May 2020, it powered the rise of the world’s first exascale distributed computing platform [7]. It also powers pugrid.net through ACEMD, a biomolecular simulation package for drug discovery.
- **Simplicity:** With an extensive user guide, clear API docs, a simple script generator, and an entire ecosystem of tools built around it, OpenMM makes it easy for new users to dive in.
- **Free, open source, and permissively licensed:** Unlike many legacy biomolecular simulation packages, OpenMM is available free of charge and without the need to sign a license agreement. OpenMM is licensed under the permissive MIT License (with some of the CUDA and OpenCL components under the LGPL License) to enable it to have maximum impact in the biosciences, and allowing it to be used “under the hood” within other software applications, often without the user even knowing OpenMM is the underlying simulation engine. The MIT License allows OpenMM to be used, modified, and redistributed without restriction, provided attribution is maintained, ensuring that even commercial entities can use OpenMM.

**OpenMM needs sustainable federal support to continue to power the biomolecular modeling and simulation community into the next decade.** OpenMM now needs sustainable funding to continue to meet the needs of the biomolecular modeling and simulation community, and grow to deliver the technology needed to drive new innovations in the field as machine learning becomes more tightly integrated at every level. Previously, OpenMM development was funded via NIH support to Vijay Pande, formerly a professor at Stanford University; Pande also founded the Folding@home Distributed Computing project (now governed by the Folding@home Consortium [39], of which co-I Chodera is a member). Specifically, OpenMM development was supported by Simbios, funded via NIH Roadmap for Medical Research Collaborative Grant U54 GM072970 (which ended in 2013) and subsequently by NIH grant R01 GM062868 (which ended in 2018 when Vijay Pande transitioned to a General Partner at Andressen Horowitz, stepping down from his position to become an Adjunct Professor at Stanford); these grants funded a single software scientist (Eastman) throughout its development history. This has left OpenMM without long-term financial support; funding for the current lead developer (Peter Eastman) will expire in early 2021 when a Chan-Zuckerberg Initiative Essential Open Source Software for Science (\$180K total costs for one year, non-renewable) grant ends.

**Without NIH support, OpenMM development and maintenance will cease**, and the large biomolecular modeling and simulation community that depends on OpenMM will be left without a vital tool for quantitative bioscience. While other MD packages have recently added GPU support, no existing alternative code can fulfill all the roles that OpenMM currently does in enabling research, allowing legacy codes to exploit modern hardware, and enabling the rapid development of new molecular modeling and simulation applications, and development of a new equivalent framework that can fill these roles from scratch would cost millions of dollars take many years to develop.

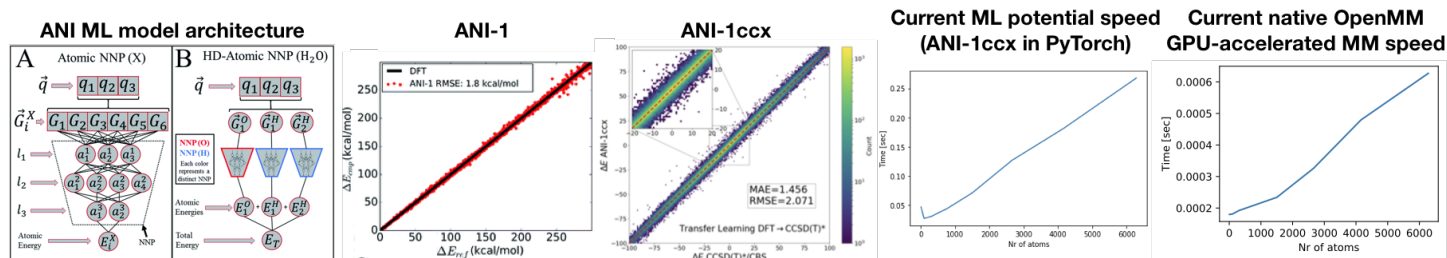
**A new leadership team will lead the transition toward a distributed development and community governance model while driving the next generation of technology innovations to support machine learning for biomolecular simulation.** While Vijay Pande remains involved in the role of collaborator (see Letter) to ensure that OpenMM can drive new research innovations in the emerging field of machine learning for biomolecular modeling, this proposal is led by a new investigator team that has been deeply engaged with OpenMM throughout its lifetime: **PI Markland** (Stanford) has led new developments in path integral molecular dynamics methods within OpenMM; **co-I Chodera** (MSKCC) has been deeply involved in extending OpenMM to treat small molecules and for use in free energy calculations and the use of machine learning potentials in drug discovery; and **co-I de Fabritiis** (Barcelona) is a GPU MD pioneer and leader in machine learning for biomolecular modeling. Senior personnel **Peter Eastman** (lead OpenMM developer) will continue to lead development as OpenMM transitions to a distributed development model, bringing in new developers at MSKCC and Barcelona and cultivating community developer contributions. In addition, the leadership team has formed a new **OpenMM Consortium** consisting of leaders in the field that actively develop with or use OpenMM and will help steer future development and community engagement plans to ensure OpenMM can best meet the needs of the community. Together, the investigators and OpenMM Consortium will ensure OpenMM remains a key tool driving new discoveries and innovations within the biomolecular modeling community.



## INNOVATION

This proposal aims to go much further than simply *sustain* development of OpenMM with the addition of new features and optimization to take advantage of new graphics processor hardware. Here, **we propose to radically extend OpenMM to exploit and drive new innovations within the machine learning community for biomolecular modeling that promises to revolutionize the practice of biomolecular modeling in a few short years.**

**Machine learning is driving a technological revolution in science and engineering.** Advances in a key machine learning (ML) technology—deep learning, in which differentiable universal function approximators composed of simple layers can easily be trained using automatic differentiation techniques—has driven a tidal wave of technology development that is transforming many fields of science [40]. To be successful, ML methods require three main ingredients: (1) a **useful representation** of the system to be modeled that allows a useful relationship to be learned from a tractable amount of data; (2) **large datasets** that can be used to train and assess the models, and (3) **fast hardware-accelerated implementations** for both training and using the models in practical applications.



**Figure 1. Quantum machine learning (QML) potential models like ANI-1 can reproduce quantum chemical potential energies with high accuracy, but current implementations are too slow to be practical for biomolecular simulation.** The ANI-1 model [41] (left) achieves excellent accuracy compared to DFT quantum chemical calculations (middle). Second-generation ANI-1ccx [42] achieves high accuracy compared to high-level CCSD(T) (middle). These methods are currently 200× slower than MM potentials in simulating even alanine dipeptide in a water droplet on GPUs (right; NVIDIA GTX-1080). While ML frameworks like PyTorch and TensorFlow make effective use of the GPU for matrix operations (with hardware improvements focusing on accelerating these operations), current implementations of QML potentials lack GPU-accelerated kernels for atomic features. **Our open source QML GPU kernel library** will enable OpenMM and other codes to use QML potentials at near-MM speeds, enabling existing modeling applications to use QML and QML/MM models.

**Machine learning holds the potential to transform biomolecular modeling and drug discovery.** Several key classes of models have emerged: **Quantum machine learning (QML)** potentials aim to replace the expensive, poorly-scaling machinery of quantum chemistry with fast machine learning (ML) models (often based on deep learning) that can approximate the quantum chemical Born-Oppenheimer potential energy surface to high accuracy at greatly reduced cost [43–45]. The ANI class of models [41, 42, 46] (**Figure 1**), for example, can approximate DFT (ANI-1x [41], ANI-2x [46]) or coupled-cluster (ANI-1ccx [42]) calculations with extremely high accuracy at a fraction of the computational cost. A large number of promising models have been developed in the last few years for learning high-dimensional potential energy surfaces (such as Behler-Parrinello [47], ANI [41, 42, 46], AIMNet [48], SchNet [49, 50]; see **Figure 4**). The speed and flexibility QML models positions them to supplant slow and complex QM/MM methods [51], where subsystems <100 atoms (such as small molecule ligands) are treated with QML and the remainder with MM, as recently demonstrated in [1, 52, 53]; the accuracy of QML methods and ease of training to both quantum chemical and experimental data [2] make them appealing to eventually supplant MM methods entirely in several years as datasets, implementations, hardware progress alongside improvements in the treatment of long-range interactions [54].

Simultaneously, a variety of **translationally- and rotationally-invariant ML architectures** for predicting molecular properties (such as solubilities or binding affinities) from static structures have emerged that promise to make useful predictions from static structures. These models include Tensor Field Networks [55], Clebsch-Gordan nets [56–58] and graph convolutional or message-passing networks [59] such as PotentialNet [60].

At the same time, the ML revolution has driven the development of deep learning models as flexible function approximators for learning complex **collective variables** that describe biomolecular structural changes for use in enhanced sampling algorithms for studying processes like conformational changes and protein-ligand binding [61, 62], or in the extremely popular framework of **metadynamics** [63–65] (which was used in 1,490 papers in 2020 alone [5]). Recent developments include methods such as neural tICA [66, 67], VAMPnets [68], Anncolvar [69], SGOOP [70], and RAVE [71] that encode their collective variables using some differentiable representation of Cartesian coordinates (to allow chain-rule propagation of forces) connected to a differentiable (deep) model. While caution must be exercised in constructing these models in a manner that produces meaningful data, there is no denying the fact that this is a highly active area of research that would benefit from facile tools for moving between efficient molecular simulation engines and high-productivity machine learning frameworks.

**To impact the biological sciences, these models need to be deeply integrated with biomolecular modeling frameworks.** While published accuracies of these ML methods signal their enormous potential, their impact in bioscience has been significantly limited for multiple reasons that we propose to remedy this project:

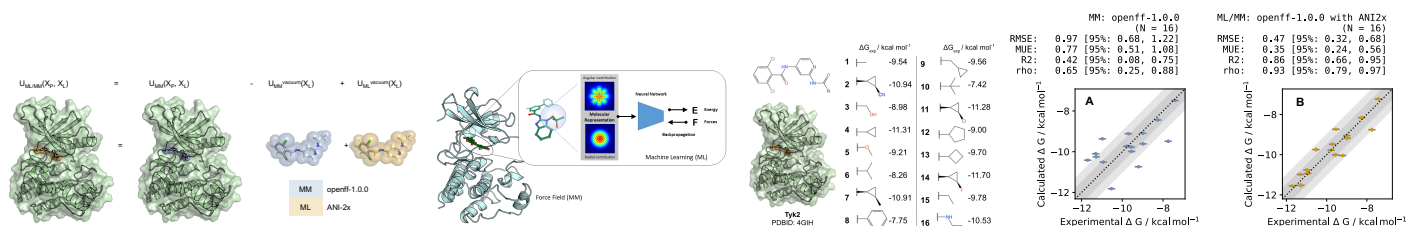
- **Barriers to use:** The appeal of forcefield-agnostic simulation packages like OpenMM is that multiple models can easily be used via the same interface, greatly increasing the productivity of researchers and developers. Currently, using a new QML model requires installing author source code, developing a new interface from scratch, and building all the associated algorithms for molecular simulation (such as integrators) atop it. We will solve this problem by creating plugins that allow ML models to be easily imported into OpenMM from a variety of common ML frameworks (**Aim 2**) and a common repository of popular QML potentials that can be used from within OpenMM as drop-in replacements for existing MM potentials (**Aim 4**).
- **Speed:** While QML potentials like ANI [41, 42, 46] are significantly faster and scale better than quantum chemical calculations, they are still incredibly slow compared to typical MM force fields—current implementations are 200–15,000× slower than typical MM force fields (such as AMBER ff14SB [72]) in simulating solvated droplet or periodic systems (**Figure 1**), rendering them impractically slow for all but the most trivial of applications. Implementations that narrow the speed gap between MM and QML are essential for these models to find practical use. **Initial use cases will focus on hybrid QML/MM potentials**, where a subsystem (such as a small molecule ligand or reactive center) is treated with QML while the remainder is treated with MM [1, 52, 53]—which can provide snear-MM speeds with near-QM accuracy [1]. Later applications will enable QML to treat the entire biomolecular system—once subsequent developments in QML potentials that include accurate long-range deliver improved accuracy and trainable improvements. We will replicate our successful GPU kernel optimization strategies in optimizing the speed bottleneck—computation of atomic feature vectors—to accelerate QML methods (**Aim 2**).
- **Hybrid physical/machine learning models:** Many current QML models include only short-range interactions [41, 42, 47], making it challenging for these systems to accurately model condensed-phase biomolecular systems where long-range interactions are important. While a number of interesting directions are being explored in terms of continuous-convolution [49, 50, 73] or message-passing [48], a promising avenue for development is the use of physical potentials to describe these long-range interactions [54]. In parallel, traditional MM force fields are adopting ML models to describe correlated many-body local effects [74]. **We will make it easy to develop and deploy models that combine short-range QML with long-range physical energy terms**, allowing these to easily be combined with QML models in OpenMM simulations and within popular ML frameworks (TensorFlow, PyTorch, JAX) by developing a library of custom ops that can easily be used in both OpenMM and ML frameworks (**Aim 3**).
- **Data availability:** The construction of new QML potentials requires large, high-level quantum chemical datasets that require access to major supercomputing resources and quantum chemistry expertise to generate, putting them out of reach of most researchers. We will collaborate with the MolSSI QCArchive [75] and Folding@home [76] (which has >1M active CPU cores) to produce large, high-quality quantum chemical datasets for use in developing improved QML potentials for biomolecular modeling (**Aim 4**).
- **Useful abstractions and feature computation kernel libraries:** The development of new QML models also requires researchers re-implement all the featurization schemes necessary prior to the deep learning model(s) used by these methods to fit data. Here, we will build an open source library of optimized GPU kernels that can be used within both molecular simulation and machine learning frameworks as custom ops (PyTorch, TensorFlow, and JAX) to accelerate research and development in this domain (**Aim 3**).
- **Flexibility:** OpenMM will also be able to spur innovation by allowing researchers to flexibly utilize these capabilities in other places arbitrary functions may be of interest, such as using ML models to define collective variables, biasing potentials, or advanced integrators to explore enhanced sampling schemes (**Aim 2**).

**This proposal aims to augment OpenMM with new capabilities that enable it to exploit the machine learning revolution** and power the next generation of research and applications in biomolecular modeling and simulation. **Aim 2** describes how QML potentials, ML models for collective variables, and expressive compute graphs for advanced integrators and samplers can be utilized within OpenMM for immediate use by researchers and developers that use OpenMM in their applications, or in legacy codes that use OpenMM to exploit modern hardware. **Aim 3** will enable the use of OpenMM physical forces within popular ML frameworks, as well as provide large quantum chemical datasets that will help accelerate research into QML potentials—especially hybrid potentials. And **Aim 4** will leverage Folding@home to generate new large datasets of high value for building new QML models with sufficiently expanded coverage and improved accuracy to transform biomolecular modeling.

## New preliminary data

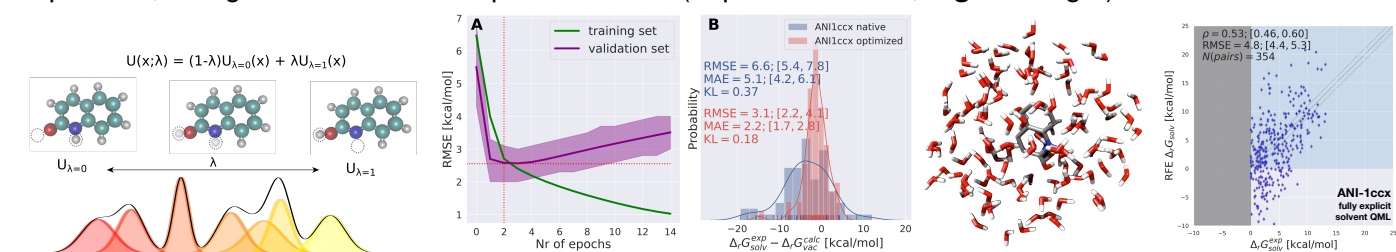
Since initial submission, we have made significant progress in showing the utility of our proposed R&D activities:

**Utility of hybrid QML/MM:** In the short term, hybrid QML/MM systems, in which a <100-atom subsystem is treated with QML (such as a small molecule ligand), will be the primary target for initial development and optimization. While many formulations of hybrid QML/MM potentials are possible, the simplest formulation simply replaces the intramolecular MM ligand interactions—which are often poorly modeled—with QML [1]. As seen in **Figure 2**, even the simplest possible approach can cut the error in protein-ligand alchemical binding free energy calculations **in half**, which would have enormous impact on drug discovery. These QML/MM simulations are currently only on the verge of practicality because of the enormous inefficiencies in current QML / MM interfaces, motivating the proposed work here would allow the use of QML/MM at near-MM speeds.



**Figure 2. Hybrid QML/MM potentials that model the ligand with QML and the environment with MM can reduce the error in protein-ligand alchemical binding free energy calculations from 1.0 kcal/mol to 0.5 kcal/mol [1].** *Left:* While many formulations of hybrid QML/MM potentials are possible, the simplest formulation replaces MM ligand intramolecular interactions with QML, enabling high-accuracy modeling of ligand torsions and intramolecular torsion-torsion and torsion-valence couplings. *Right:* When applied to the Tyk2 benchmark system—a challenging kinase-inhibitor benchmark system from the Schrödinger FEP+ test set (FEP+  $\Delta G$  RMSE  $0.93 \pm 0.12$  kcal/mol [17]; GAFF 1.8  $\Delta G$  RMSE 1.13 kcal/mol), a hybrid QML/MM model using the Open Force Field Initiative [<http://openforcefield.org>] OpenFF 1.0.0 (“Parsley”) small molecule force field [77, 78], AMBER14SB [72], and TIP3P [79] reduces the RMSE from 0.97 [95% CI: 0.68, 1.22] kcal/mol for MM to 0.47 [95% CI: 0.32, 0.68] kcal/mol for QML/MM, which would significantly accelerate drug discovery.

**Alchemical free energy calculations are extremely easy in fully explicit QML simulations:** While it will likely take the field several years—with the right infrastructure support proposed here—to make fully QML simulations practical for fully solvated biomolecular systems by appropriately addressing issues with long-range forces (where work is already underway [54]), these methods hold enormous promise. We have recently demonstrated that we can easily perform alchemical free energy calculations on systems entirely treated with QML (**Figure 3**), and have shown for the first time that both quantum chemical energies and forces and experimental data (here, tautomer free energies) can be used to train these machine learning potentials so as to better generalize from limited data [2]. In addition, alchemical free energy calculations in fully explicit QML are practical and extremely simple to implement, though slow in current implementations (unpublished data; **Figure 3** right).



**Figure 3. Alchemical free energy calculations in pure QML systems are trivial, and can be trained based on experimental free energies to generalize from limited data [2].** *Left:* Alchemical free energy calculations in systems treated entirely by QML are incredibly simple since the absence of singularities in the potential means that linear mixing of potentials in intermediate alchemical states can lead to well-behaved free energy estimates [2]. *Middle:* In addition to quantum chemical energies and gradients, experimental measurements, such as free energies, can be easily used to train QML potentials (A) to reduce the error in predicting properties on other molecules (B), as in this case of tautomer free energies [2]. *Right:* Alchemical free energies can also be computed for fully explicitly solvated QML systems (though at greater cost), as in this case where tautomer ratios are computed with ANI-1ccx (unpublished data). Note that the accuracy of this model (4.8 kcal/mol, compared with 3.1 kcal/mol for ANI-1ccx optimized to experimental data as well) is initially poor because it was parameterized for only gas phase systems; the intent is simply to demonstrate the fact that QML potentials can be integrated into existing MM workflows once they reach maturity.

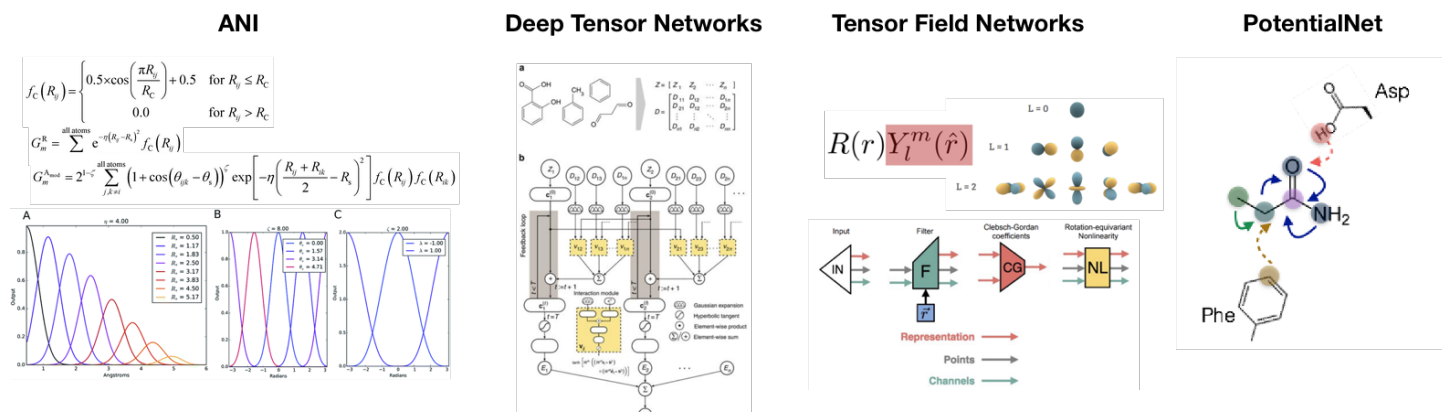
## APPROACH

### AIM 1: Transition OpenMM from a single-laboratory single-developer-led code to a community governance and distributed developer model to ensure long-term sustainability.

**Rationale:** OpenMM was developed in a single laboratory (collaborator **Pande**) with support from the NIH, primarily by a single developer (**Eastman**). To ensure OpenMM can sustainably support its extensive user community, we aim to create an OpenMM Consortium that will oversee its future roadmap and governance, as well as recruit and train two additional paid developers within two key OpenMM-focused laboratories (the **Chodera** lab at MSKCC and the **de Fabritiis** lab in Barcelona). In addition, we will seek to recruit more core developers from other laboratories and the biomolecular simulation community at large that depend on OpenMM.

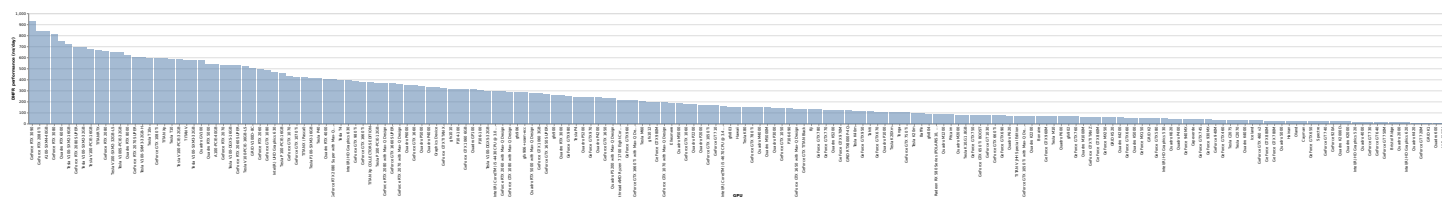
**Approach and Milestones:** We have formed a new core leadership team—consisting of **PI Markland**, **co-I Chodera**, and **co-I de Fabritiis**, together with lead OpenMM developer **Eastman** and collaborator Vijay Pande (original OpenMM PI, machine learning pioneer, and Folding@home founder)—which collectively possess extensive expertise in biomolecular simulation, GPU acceleration, and biomolecular machine learning. This leadership team will be responsible for:

- **Recruit new developers to OpenMM:** We will recruit additional developers at MSKCC and Barcelona that will allow core development tasks to be distributed among the development team, ensuring that development is no longer dependent on a single individual. Additionally, the leadership and developer team will seek to recruit additional contributors from other laboratories that depend on OpenMM (such as Andy Simonett [see Letter] who has made significant contributions already).
- **Develop a public roadmap:** We will develop and maintain a public roadmap for future bugfix, feature, and API-changing OpenMM releases using GitHub Project Boards [80] to ensure the community is well-informed about (and can provide input into) the timing and content of new releases.
- **Accelerate OpenMM on new hardware:** The leadership and developer teams will continue to identify new hardware that would be particularly beneficial for OpenMM to exploit—such as new inexpensive consumer-grade GPUs, new discrete GPUs (such as Intel Xe), integrated hardware accelerators, and tensor processing units (TPUs). Developers—led by Eastman’s extensive experience in hardware acceleration—will tailor kernels to continue to exploit the latest hardware (and software frameworks for using them, such as Apple’s recent Metal and Intel’s oneAPI).
- **Field community-driven feature requests and bug reports:** The development team will clearly and visibly prioritize and address community-provided feature requests and bug reports using GitHub’s issue tracker, issue labeling system, and release milestones, aiming improve on metrics such as open issue duration.
- **Produce enhanced documentation and tutorials:** While OpenMM has extensive user and API documentation and a number of tutorials, we will improve the utility of this documentation and the number of tutorials targeted at new users coming from the machine learning field and users transitioning to OpenMM from other simulation packages could greatly increase the impact and audience of OpenMM. Currently, these resources are already highly utilized, with  $\sim 1500$  unique visitors/month accessing the user guide and API docs and  $\sim 500$  unique visitors/month accessing the tutorials, suggesting effort spent on improving existing



**Figure 4. Popular classes of atomic features for popular atomic machine learning (ML) potential models.** The ANI class of models uses distance- and angle-based features [41, 42]. Continuous convolution networks like Deep Tensor Networks [73] and SchNet [49, 50] use distance-based features. Another class of models, Tensor Field Networks, use spherical harmonics [55–58]. PotentialNet uses a graph convolutional network augmented by distance-dependent edges [60]. To both achieve high performance within OpenMM (**Aim 2**) and enable rapid QML model development in common ML frameworks such as TensorFlow, PyTorch, and JAX (**Aim 3**), we will implement hardware-accelerated compute kernels for these and other common atomic feature computation schemes for use both within OpenMM and as a library of operations ("ops") available within common ML frameworks.





**Figure 5. OpenMM 7.4.2 runs efficiently on an extremely wide variety of GPUs.** Performance (in ns/day) on the joint AMBER-CHARMM (JAC) DHFR benchmark benchmark utilized a 4 fs timestep with hydrogen mass repartitioning (HMR) [86] and the high-quality BAOAB Langevin integrator [87] with bonds to hydrogen constrained [88] and a Monte Carlo barostat [89]. Maximum performance (925 ns/day on the GTX-3090) is nearing a microsecond/day on consumer-grade GPU hardware. Benchmark data courtesy Folding@home.

materials will also be highly valuable. We will focus on improving online documentation and materials, which are broadly accessible to the global biomolecular simulation community, and tutorials that can easily be launched in the browser (such as MyBinder and Google Colab) to further minimize barriers to entry, using web engagement statistics to measure success in increasing engagement.

- **Transition to a community governance model under the OpenMM Consortium:** The leadership team will establish the OpenMM Consortium, consisting of academic researchers, software scientists, and industry scientists with experience in OpenMM and biomolecular modeling to ensure that OpenMM continues to meet the needs of the biomolecular modeling community by setting priorities, integrating best practices, and developing metrics for measuring progress. The leadership team will work closely with the Molecular Sciences Software Institute (MolSSI) [81] (see Letter) to develop a governance model following equitable best practices for a healthy software project responsive to community needs and expectations.

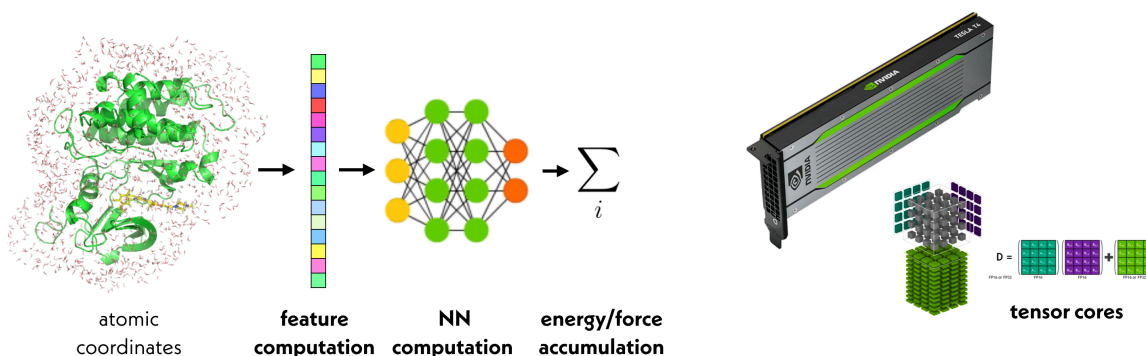
**Preliminary data and feasibility:** We have recently consolidated OpenMM into a central GitHub organization independent of any single laboratory [82]. We have also recruited an initial OpenMM Consortium [83], consisting of members of the biomolecular simulation and modeling field who can advise on how OpenMM can best serve the needs of the community (many of which have contributed Letters): Emilio Gallicchio (Brooklyn College), Justin McCallum (University of Calgary), Vijay Pande (Stanford/a16z), Jean-Philip Piquemal (Sorbonne), Jay Ponder (WashU), Julia Rice (IBM Almaden), Josh Rackers (Sandia), Pengyu Ren (UT Austin), Karmen Condic-Jurkic (MSKCC), Andy Simmonett (NIH), and Bill Swope (IBM). The OpenMM Consortium has already held initial virtual meetings to develop a preliminary roadmap for how OpenMM can best serve both the advanced potential function development community (which led to AMOEBA [84] and HIPPO [85]) as well as the quantum machine learning community. Though the current Consortium currently serves in an advisory role, it will ultimately be tasked with helping develop a community governance model in coordination with MolSSI.

While current benchmarks indicate excellent performance on an extremely wide variety of GPUs (**Figure 5**), new hardware may require the use of new acceleration APIs to fully exploit it. We have recently refactored the OpenMM GPU platforms into a single common API to make it easier to rapidly retarget to new acceleration APIs (such as Apple Metal, AMD RocM, and Intel oneAPI) as needed to achieve state-of-the-art performance on new hardware.

**AIM 2: Develop a next-generation hardware-accelerated molecular simulation platform that will enable direct interoperability with machine learning frameworks.**

**Rationale:** Biomolecular modeling is at the cusp of a revolution, with quantum machine learning (QML) potential energy models now able to achieve quantum chemical accuracy using deep learning models that exploit translationally and rotationally invariant features [43–45, 90] (**Figure 1**). While this field is still young, it is unmistakably clear that there will be an explosion in the exploration and application of QML models and hybrid QML/MM models in biomolecular modeling over the next decade, as researchers seek increased accuracy, the ability to tackle problems previously inaccessible or difficult with QM/MM (such as enzyme catalysis and covalent inhibition), and exploit the ability to easily improve these models with quantum chemical and experimental data.

Current implementations are impractically slow—200–15,000 $\times$  slower than MM simulations within OpenMM—and require writing custom interfaces to glue them to MD codes, which can also result in even more performance loss. To make these methods practically usable by the biomolecular modeling community, we will provide a way to easily import and use ML models developed with popular ML frameworks (TensorFlow [91], PyTorch [92], and JAX [93]) within OpenMM in a flexible manner that allows their use as potentials, collective variables, biasing functions, and integrators. Currently, the major computational bottleneck is that these models must compute fixed-vector-length features from biomolecular systems using ML ops that are not well-suited to dealing with biomolecular data (such as pairlist operations) prior to feeding these feature vectors to highly optimized deep learning models that can exploit hardware that is rapidly improving to accelerate the deep learning stage of the computation (**Figure 6**). We will accelerate these methods by creating custom kernels to eliminate this major computational bottleneck using the same advances that OpenMM has made in accelerating MM potentials.



**Figure 6. Accelerating atomic feature computation through custom GPU-accelerated ops will eliminate the computational bottleneck hindering widespread adoption of QML potentials.** *Left:* The key bottleneck in current QML implementations is the computation of atomic features for biomolecular systems, especially for periodic systems. These operations involve irregular computations like pairlist construction and sums over local atom sets, as well as computation of their derivatives, which are operations that ML frameworks like TensorFlow and PyTorch are not efficient for. By contrast, OpenMM’s GPU-accelerated kernels excel at operations like this. We will implement GPU-accelerated custom ops for TensorFlow/PyTorch that speed up these feature computation operations. *Right:* While the subsequent computational stage—neural network computation—may still be computationally costly on current hardware, both hardware and software are aggressively being optimized by hardware vendors and ML package developers to accelerate these computations, including innovations like adding tensor cores. We will therefore primarily focus our effort on accelerating featurization (**Aim 2**) and enabling developers to distill/compress ML models to tune their performance on common accelerator hardware (**Aim 3**).

### Approach and Milestones:

- Develop hardware-accelerated kernels to accelerate QML potentials within OpenMM:** We have developed prototype OpenMM plugins that allow ML models defined within TensorFlow or PyTorch to compute potential energy components or collective variables within OpenMM. These plugins currently use the corresponding ML framework to compute all features directly from Cartesian atomic coordinates, resulting in poor performance compared with OpenMM’s GPU-accelerated MM energy functions, which use GPU-accelerated kernels that exploit special properties of atomic coordinates (**Figure 1**). We will develop custom GPU-accelerated operations (“ops”) that can be used to replace inefficient atomic feature computation code within these ML frameworks, prioritizing key atomic features used in popular QML models (**Figure 4**): such as Behler-Parrinello [47], ANI [41, 42, 46], AIMNet [48], SchNet [49, 50], Tensor Field Networks [55], Clebsch-Gordan nets [56–58], and graph convolutional or message-passing networks [59] such as PotentialNet [60]. Simultaneously, new generations of hardware will continue to feature optimizations to accelerate the deep learning half of the computation without our help, enabling us to focus on eliminating feature computation as a computational bottleneck. (In **Aim 3**, we discuss how we will also provide model compression tools to help tune models for improved performance.)
- Enable ML models to be flexibly imported and used in a variety of innovative new ways within OpenMM:** In addition to QML potential support, we will make it easy to import models from ML frameworks such as TensorFlow, PyTorch, or JAX for use wherever a function approximator could be useful, such as collective variables, biasing potentials, or integrators (which are simply methods that update positions and velocities and perform computations on them).

**Preliminary data and feasibility:** OpenMM has been at the forefront of offering users and developers convenient ways for specifying differentiable functions to flexibly support new potential functions and collective variables while not compromising on speed. The CustomForce facilities of OpenMM—introduced a decade ago [26]—allow users or developers to specify algebraic expressions or spline functions (without the need to write new code) that are symbolically differentiated and compiled into GPU kernels that run at native speed, even allowing differentiation with respect to force field parameters. For example, implementing a Coulomb potential alchemically modified by the parameter `lambda` is as simple as specifying its algebraic form: `'lambda*C*charge1*charge2/r'`. Extending OpenMM to support differentiable machine learning models through compute graphs via TensorFlow, PyTorch, or JAX is conceptually similar, though differentiation is accomplished through forward- or backward-mode automatic differentiation, rather than via symbolic algebra. In both cases, just-in-time compilation provides the ability to rapidly build optimized kernels targeted to multiple hardware accelerators (such as GPUs).

**PyTorch and TensorFlow plugins for OpenMM:** As a proof of concept, we have developed two new plugins that allow models developed in popular machine learning frameworks to be used within OpenMM as quantum machine learning (QML) potentials or collective variables, both of which can be used alongside normal OpenMM physical potential function forces: (1) the **OpenMM TensorFlow Plugin** [94] allows models exported from TensorFlow [91] to be used within OpenMM, while (2) the **OpenMM PyTorch Plugin** [95] allows models exported from PyTorch [92] to be used within OpenMM.

We have demonstrated that it is possible to import ANI models [41,42,46] defined in TorchANI [96] into OpenMM and simulate a box of water molecules within OpenMM, using our using our prototype PyTorch OpenMM plugin, thereby enabling any applications built on OpenMM to use ANI in place of the wide variety of currently-supported molecular mechanics potentials. While both periodic and non-periodic systems can be simulated, initial performance benchmarks of the unoptimized models show that a box of 215 waters achieves 60 timesteps/second with periodic boundary conditions or 150 timesteps/second without, suggesting the periodic feature computations within the ANI-1 model implemented fully in TorchANI are particularly inefficient.

**Accelerated QML GPU kernel library:** We have started the earliest experimental stages of implementing a library of GPU-accelerated kernels for common QML potentials (available on GitHub at <https://github.com/peastman/NNPOps>). Initial implementations utilize both fast symmetry function evaluation and batched neural network computations, providing an initial 18× speedup over TorchANI for computing energy and forces for a small molecule ligand (46 atoms) with ANI-2x; future optimization will be able to bring this to at least 50×, enabling us to reach nearly-MM speeds for hybrid QML/MM potentials where a <100-atom small molecule ligand is treated with QML and the remainder with MM.

This work will be guided by our machine learning expertise (lead developer **Eastman** and collaborator **Pande** recently co-authored a book on ML for the biosciences [97], and **co-I de Fabritiis** recently achieved highest predictive accuracy in the blinded D3R Grand Challenge 4 [98]), close coordination with the MolSSI Machine Learning team, input from the OpenMM Consortium, and laboratories at the forefront of QML (see Letters).

### **AIM 3: Create a software infrastructure to enable hybrid machine learning / physical model development.**

**Rationale:** Two major software impediments stand in the way of building new QML methods with greatly increased accuracy and expanded coverage for biomolecular systems: (1) Training new ML models on 20M+ molecule datasets [41] is significantly hindered by feature computation—the slow speed of computing complex input features from atomic coordinates subsequently fed into fast deep learning models (which are already hardware accelerated). (2) While QML models excel at modeling short-range interactions, the need to use physical potential models for long-range interactions (as in PhysNet [54], which explicitly includes long-range electrostatics) is a clear direction for future research. Unfortunately, it is very difficult to bring physical models into modern ML frameworks—especially complex electrostatic models that much be implemented from scratch. By enabling OpenMM’s hardware-accelerated feature computation kernels to be used within popular frameworks like TensorFlow and PyTorch, and creating new hardware-accelerated kernels for physical energy terms, we can accelerate research into new ML models not just for fitting quantum chemical energies of biomolecules, but also applications like enzyme design, drug discovery, biotherapeutic optimization, and predicting the functional impact of mutations.

#### **Approach and Milestones:**

- **Produce an open source hardware-accelerated QML kernel library to accelerate QML development, training, and deployment:** To accelerate the development, training, and deployment of new QML models, we will package our custom hardware-accelerated kernels for performing common QML computations (e.g., atomic featurization and message-passing) into a library that can easily be used in both popular ML frameworks (PyTorch, TensorFlow, JAX) and MD engines. Besides making it easier for QML developers to easily build models without coding atomic features and accelerating the training process, this will also make it easier for QML developers to deploy their models within OpenMM (and other MD engines that make use of the kernel library) at full speed, since the same kernels will be used within OpenMM.
- **Facilitate the development of hybrid physical/ML models by creating custom ops for physical interaction kernels:** Hybrid physical/ML models such as PhysNet [54] (which predicts the partial charges at every step and uses a classical electrostatics model) are only beginning to be explored because of the difficulties in implementing optimized physical potential energy functions in ML frameworks. To enable rapid exploration of hybrid models, we will also port OpenMM’s optimize physical energy function kernels (including advanced physical terms such as those in AMOEBA [84, 99–101], HIPPO [85], CHARMM polarizable force fields [102], and Gaussian electrostatics models [103]) over to optimized ML framework ops to use as building blocks for hybrid models, allowing them to run at high speed both within ML frameworks and OpenMM.

- **Unify ML and MM APIs to accelerate science at the interface:** The difficulty of developing hybrid models is exacerbated by the complexity of methods for setting up models for complex biomolecular systems. We will simplify this process by migrating OpenMM's System description classes to a pure Python object model that will allow users to render a parameterized system either to a native OpenMM Context that can run at full speed on the GPU or to a compute graph representation within ML frameworks like TensorFlow, PyTorch, and JAX that allow the model to be differentiated with respect to arbitrary model parameters for use in force field parameterization.
- **Develop a library of tools for optimizing QML model performance within OpenMM (and other MM packages):** While hardware and software developments outside of this project will continue to deliver increased performance for the neural network stages of QML models (**Figure 6**), we can do more to assist QML model developers in delivering increased performance without compromising accuracy. While a number of techniques based on **model distillation** [104], compression [105], and incremental quantization [106] techniques have been proposed for producing new networks that deliver equivalent accuracy but require greatly reduce FLOP counts, these acceleration techniques have not yet seen use in QML potential refinement and deployment. We will introduce a number of promising approaches into our QML model acceleration and export library to make it easier for QML developers to tune the performance of their models prior to deployment in our QML potential model archive.

### Preliminary data and feasibility:

OpenMM's Python API has enabled it to already see significant use within machine learning frameworks. For example, a recent Science paper [107] wrapped OpenMM within a TensorFlow op (without further optimization) to train invertible flows for sampling directly from the Boltzmann distribution. This Aim focuses on finding ways to accelerate this use to make research in this area more practical by eliminating pain points and breaking down the barriers between ML frameworks and physical modeling codes such as OpenMM. We note that others have already demonstrated the feasibility for achieving significant performance optimizations by developing custom hardware-accelerated ops for some physical forces in projects such as `jax-md` [108] and `timemachine` [109], but these attempts demonstrate feasibility and do not contain any of the significant infrastructure required by a fully-featured MD code like OpenMM that is aimed at working with biomolecular systems.

**Implementation of MM potential terms within PyTorch:** We recently released the TorchMD package, which enables physical MM potential terms to be used alongside QML potentials in experimenting with the construction and training of hybrid physical modeling / machine learning potentials that could combine long-range physical interactions with short-range machine learning potentials. The open source TorchMD code is available on GitHub [<https://github.com/torchmd/torchmd>]. This implementation does not yet feature GPU-accelerated MM kernels, which would accelerate PyTorch training by 50×.

### AIM 4: Develop community repositories of large quantum chemical datasets and models to accelerate development and deployment of QML potentials.

**Rationale:** Publicly available quantum chemical datasets of sufficient size for building useful quantum machine learning (QML) potentials covering biomolecules (such as the ANI-1 Dataset [41]) are still highly limited in chemical coverage, conformations, composition, and available levels of quantum theory. To ensure researchers can easily train, validate, and assess new QML potentials—and especially hybrid models that blend ML and physics—we will work in close collaboration with the MolSSI QCArchive project and the Folding@home Consortium to generate large quantum chemical datasets for biomolecules that are useful in both OpenMM and common ML frameworks.

We will additionally make it easy to distribute, modify, and use QML potentials through developing an interchange standard and online repository of QML potentials in close collaboration with the MolSSI QCArchive and Machine Learning team.

### Approach and Milestones:

- **Enable facile and reproducible deployment of QML potentials:** OpenMM currently features a flexible ForceField engine using a force field agnostic force field specification, allowing biomolecular systems to quickly and easily be parameterized with a variety of common force fields by just changing the force field specification string (including nearly all force fields commonly distributed with AMBER [36] and CHARMM [29], as well as small molecule force fields such as GAFF [110] and the Open Force Field Initiative SMIRNOFF force fields [111]). We will extend the OpenMM API to similarly make it easy to construct new QML models in a similar fashion, ensuring that a multitude of popular QML potentials can easily be accessed by users by just



specifying the released model. We will work closely with the MolSSI QCArchive and Machine Learning teams to establish standards for QML model exchange, as well as a common repository that could be shared by multiple packages, and tools for accessing these models programmatically and via the web (inspired by the QCArchive Machine Learning Datasets Repository [112], ModelZoo [113], and the programmatic flexibility of PyTorch Geometric [114]), and work with developers of these models to encourage them to deposit their models at time of publication. Initially, we will work closely with the TorchANI developers (see Letter) to focus on deploying and accelerating the widely popular ANI family of models [41, 42], which have already seen use in modeling protein-ligand binding for drug discovery through a custom TorchANI-NAMD interface [52].

- **Develop a programmatic interface for exporting/deploying QML models to the QML repository:** To make it easier for QML model developers to deploy their models for use within OpenMM (and other packages that ultimately adopt our interchange format and use the MolSSI QML model repository we co-develop), we will create programmatic interfaces for easily exporting and deploying QML models developed using our custom op library.
- **Develop a Folding@home CPU compute core to generate public quantum chemical datasets to accelerate QML potential development:** We will collaborate with MolSSI software scientists and Folding@home developers (see Letters) to develop a Folding@home CPU compute core to perform quantum chemical calculations to populate the MolSSI QCArchive [75]. The compute core will service quantum chemical (QC) computations through the automated MolSSI QCFractal workflow infrastructure to allow MolSSI scientists to utilize the >1M available CPU cores of the planetary-scale Folding@home distributed computing platform to generate extremely large quantum chemical datasets. We will primarily use the open source quantum chemistry code psi4 [115] via the QCEngine [116] quantum chemistry interface, as our collaborators at MolSSI co-develop these codes. These datasets will be made immediately available to the quantum community for use in constructing QML or hybrid QML/MM potentials (or for other purposes, such as constructing new MM force fields or pure ML models). Because QCArchive contains the memory usage and computation time for all submitted molecules, we will be able to develop a simple predictive model to optimize allocation of work to Folding@home clients capable of efficiently handling its resource requirements.
- **Provide immediate open access to generated QC datasets through the MolSSI QCArchive:** The MolSSI QCArchive [75] and machine learning (ML) teams are in the process of collecting all significant extant datasets in the the quantum machine learning (QML) potential space. Collectively, these datasets represent  $\sim 10^8$ – $10^9$  quantum chemical computations. These quantum chemical datasets are made available to the community programmatically through the QCArchive [75], as bulk downloads via the QCArchive Machine Learning Datasets Repository [112], and on Zenodo for archival purposes. These existing datasets are heterogenous, contain different properties as labels, and have been calculated at various levels of theory. As a first step, the QCArchive project has been recomputing these data at standard levels of theory—semiempirical, Hartree Fock, Density Functional Theory (including LDA, GGA, hybrids, and range-separated hybrids) MP2, and coupled cluster—where feasible. This effort, which currently represents only  $\sim 10^7$ – $10^8$  core-hours of computation, would be tremendously bolstered by the resources provided by Folding@Home, which could provide more than  $10^9$  core-hours/year.
- **Select a variety of chemical datasets relevant to constructing high-accuracy QML models with greatly expanded coverage of chemical space:** MolSSI is already engaging the QML community to allow QCArchive users and QML community partners to generate new QC datasets within QCFractal/QCArchive; this systematic, methodical community engagement program that exploits all the tools at MolSSI's disposal—workshops, software fellows, software scientists, visiting scholars, and their position of leadership in the molecular software sciences—will be key to the long-term success of this effort. Initially, however, we plan to generate key datasets of high priority to kickstart exploration in this area, including the PDB Ligand Expo [117]; Enamine REAL building blocks and molecules containing linkers [118]; small biopolymers representative of peptides, nucleic acids; heterocyclic aromatic scaffolds [119]; small biomolecular dimers; and solvated biomolecules. Pipelines developed by the Open Force Field Initiative and MolSSI for fragmenting, capping, and submitting small molecules to QCFractal will be reused for our purposes [111]. Finally, we will need to extend these quantum chemical datasets to cover more of the periodic table. At the moment, this is due largely to a lack of datasets in the literature beyond the p-block (though we are aware of some transition metal complex datasets); we aim for the community to take a leading role in data selection as the field matures. To fill gaps in user-sourced submissions, generative models (such as g-SchNet [120]) will be applied to existing data in the QCArchive to generate new calculation candidates to keep computation pipelines full.

**Preliminary data and feasibility:** Co-I **Chodera**, working in close coordination with MolSSI software scientists (see Letter), has conducted a pilot experiment as part of the Open Force Field Initiative [111] to deploy QCFractal for the purposes of generating new quantum chemical data for use in force field parameterization, with the results deposited in QCArchive (which currently contains >28M quantum chemical calculations for >21M molecules). Tools for the processing of small molecules, their fragmentation into capped small molecules, generation of input geometries, and processing by QCFractal were developed as part of the Initiative [111]. This pilot experiment took advantage of unused CPU computing cycles at multiple academic institutions, with QCFractal computing on ~300 total cores at any one time. Over the course of several months, this pilot experiment carried out ~8M B3LYP-D3(BJ)/DZVP gradient evaluations with Psi4 [115] on a large variety of small molecules, consuming ~1.5M core-hours, and generating a total of 350K geometry optimizations and 3K torsion drives that were immediately deposited into (and are now freely accessible through) QCArchive [75]. As proposed in this Aim, deploying QCFractal on the planetary-scale Folding@home network will provide access to roughly **three orders of magnitude more compute power** than this pilot experiment, enabling the generation of extremely large (~100M–10B geometry optimizations) high-value quantum chemical datasets for immediate public use in building machine learning models with greatly expanded chemical coverage and improved accuracy.

The investigators and collaborators also have extensive experience in developing open standards for (bio)molecular models and their exchange. **Eastman** led the development of a force field agnostic specification for OpenMM's force field descriptions that has enabled multiple families of force fields to be easily encoded within a common XML format. Based on this work, Co-I **Chodera** led the development of the SMIRNOFF specification for biomolecular force fields using hierarchical typing and direct chemical perception that allows a wide class of small molecule and biopolymer force fields to be expressed [111, 121]. MolSSI led the development of the QCSchema specification for quantum chemical computations that allows QCEngine [116, 122] to support multiple back-end quantum chemistry codes via a single standard format.

## Bibliography and References Cited

- [1] Rufa DA, Bruce Macdonald HE, Fass J, Wieder M, Grinaway PB, Roitberg AE, et al. Towards chemical accuracy for alchemical free energy calculations with hybrid physics-based machine learning / molecular mechanics potentials. *bioRxiv*. 2020. Available from: <https://www.biorxiv.org/content/early/2020/07/30/2020.07.29.227959>. doi:10.1101/2020.07.29.227959.
- [2] Wieder M, Fass J, Chodera JD. Fitting quantum machine learning potentials to experimental free energy data: Predicting tautomer ratios in solution. *bioRxiv*. 2020. Available from: <https://www.biorxiv.org/content/early/2020/10/25/2020.10.24.353318>. doi:10.1101/2020.10.24.353318.
- [3] Google Scholar search for "molecular quantum machine learning potential" from 1 Jan 2020–24 Oct 2020;. Available from: [https://scholar.google.com/scholar?as\\_ylo=2020&q=molecular+quantum+machine+learning+potential&hl=en&as\\_sdt=0,5](https://scholar.google.com/scholar?as_ylo=2020&q=molecular+quantum+machine+learning+potential&hl=en&as_sdt=0,5).
- [4] OpenMM Twitter survey, accessed 25 Oct 2020;. Available from: [https://twitter.com/openmm\\_toolkit/status/1319866785236041729](https://twitter.com/openmm_toolkit/status/1319866785236041729).
- [5] Google Scholar from 1 Jan 2020–24 Oct 2020;. Available from: <https://scholar.google.com/scholar>.
- [6] Folding@home Active CPUs and GPUs by OS, accessed 22 Oct 2020;. Available from: <https://stats.foldingathome.org/os>.
- [7] Zimmerman MI, Porter JR, Ward MD, Singh S, Vithani N, Meller A, et al. Citizen Scientists Create an Exascale Computer to Combat COVID-19. *bioRxiv*. 2020. Available from: <https://www.biorxiv.org/content/early/2020/06/30/2020.06.27.175430>. doi:10.1101/2020.06.27.175430.
- [8] NIH RePORTER, accessed 22 Oct 2020;. Available from: <https://projectreporter.nih.gov/reporter.cfm>.
- [9] Frenkel D, Smit B. *Understanding molecular simulation: from algorithms to applications*. vol. 1. Elsevier; 2001.
- [10] Karplus M. *Molecular dynamics simulations of biomolecules*. ACS Publications; 2002.
- [11] Dror RO, Dirks RM, Grossman J, Xu H, Shaw DE. Biomolecular simulation: a computational microscope for molecular biology. *Annual review of biophysics*. 2012;41:429–452.
- [12] Klepeis JL, Lindorff-Larsen K, Dror RO, Shaw DE. Long-timescale molecular dynamics simulations of protein structure and function. *Current opinion in structural biology*. 2009;19(2):120–127.
- [13] Stary A, Kudrncac M, Beyl S, Hohaus A, Timin E, Wolschann P, et al. Molecular dynamics and mutational analysis of a channelopathy mutation in the IIS6 helix of Cav1. 2. *Channels*. 2008;2(3):216–223.
- [14] Wan S, Coveney PV. Molecular dynamics simulation reveals structural and thermodynamic features of kinase activation by cancer mutations within the epidermal growth factor receptor. *Journal of computational chemistry*. 2011;32(13):2843–2852.
- [15] Lemkul JA, Bevan DR. The role of molecular simulations in the development of inhibitors of amyloid  $\beta$ -peptide aggregation for the treatment of Alzheimer's disease. *ACS chemical neuroscience*. 2012;3(11):845–856.
- [16] Bottaro S, Lindorff-Larsen K. Biophysical experiments and biomolecular simulations: A perfect match? *Science*. 2018;361(6400):355–360.
- [17] Wang L, Wu Y, Deng Y, Kim B, Pierce L, Krilov G, et al. Accurate and reliable prediction of relative ligand binding potency in prospective drug discovery by way of a modern free-energy calculation protocol and force field. *Journal of the American Chemical Society*. 2015;137(7):2695–2703.
- [18] Sherborne B, Shanmugasundaram V, Cheng AC, Christ CD, DesJarlais RL, Duca JS, et al. Collaborating to improve the use of free-energy and other quantitative methods in drug discovery. *Journal of computer-aided molecular design*. 2016;30(12):1139–1141.
- [19] Abel R, Wang L, Harder ED, Berne B, Friesner RA. Advancing drug discovery through enhanced free energy calculations. *Accounts of chemical research*. 2017;50(7):1625–1632.
- [20] Gapsys V, Michielssens S, Seeliger D, de Groot BL. Accurate and rigorous prediction of the changes in protein free energies in a large-scale mutation scan. *Angewandte Chemie International Edition*. 2016;55(26):7364–7368.
- [21] Bochicchio A, Jordaan S, Losasso V, Chetty S, Perera RC, Ippoliti E, et al. Designing the sniper: improving targeted human cytolytic fusion proteins for anti-cancer therapy via molecular simulation. *Biomedicines*. 2017;5(1):9.
- [22] Salomon-Ferrer R, Götz AW, Poole D, Le Grand S, Walker RC. Routine microsecond molecular dynamics simulations with AMBER on GPUs. 2. Explicit solvent particle mesh Ewald. *Journal of chemical theory and computation*. 2013;9(9):3878–3888.
- [23] Friedrichs MS, Eastman P, Vaidyanathan V, Houston M, Legrand S, Beberg AL, et al. Accelerating molecular dynamic simulation on graphics processing units. *Journal of computational chemistry*. 2009;30(6):864–872.

- [24] Eastman P, Pande VS. Efficient nonbonded interactions for molecular dynamics on a graphics processing unit. *Journal of computational chemistry*. 2010;31(6):1268–1272.
- [25] Eastman P, Pande VS. Constant constraint matrix approximation: a robust, parallelizable constraint method for molecular simulations. *Journal of chemical theory and computation*. 2010;6(2):434–437.
- [26] Eastman P, Pande V. OpenMM: A hardware-independent framework for molecular simulations. *Computing in science & engineering*. 2010;12(4):34–39.
- [27] Eastman P, Friedrichs MS, Chodera JD, Radmer RJ, Bruns CM, Ku JP, et al. OpenMM 4: a reusable, extensible, hardware independent library for high performance molecular simulation. *Journal of chemical theory and computation*. 2013;9(1):461–469.
- [28] Eastman P, Swails J, Chodera JD, McGibbon RT, Zhao Y, Beauchamp KA, et al. OpenMM 7: Rapid development of high performance algorithms for molecular dynamics. *PLoS computational biology*. 2017;13(7):e1005659.
- [29] Brooks BR, Brooks III CL, Mackerell Jr AD, Nilsson L, Petrella RJ, Roux B, et al. CHARMM: the biomolecular simulation program. *Journal of computational chemistry*. 2009;30(10):1545–1614.
- [30] Harger M, Li D, Wang Z, Dalby K, Lagardère L, Piquemal JP, et al. Tinker-OpenMM: Absolute and relative alchemical free energies using AMOEBA on GPUs. *Journal of computational chemistry*. 2017;38(23):2047–2055.
- [31] Anaconda downloads badge for the OpenMM conda package;. <https://anaconda.org/omnia/openmm/badges>.
- [32] Harvey MJ, Giupponi G, Fabritiis GD. ACEMD: accelerating biomolecular dynamics in the microsecond time scale. *Journal of chemical theory and computation*. 2009;5(6):1632–1639.
- [33] OpenMM API documentation;. <http://docs.openmm.org/latest/api-python/library.html>.
- [34] OpenMMTools documentation;. <https://openmmtools.readthedocs.io/>.
- [35] Anaconda downloads badge for the OpenMMTools conda package;. <https://anaconda.org/omnia/openmmtools/badges>.
- [36] Case D, Ben-Shalom I, Brozell S, Cerutti D, Cheatham III T, Cruzeiro V, et al. AMBER 2018; 2018. University of California, San Francisco.
- [37] Phillips JC, Braun R, Wang W, Gumbart J, Tajkhorshid E, Villa E, et al. Scalable molecular dynamics with NAMD. *Journal of computational chemistry*. 2005;26(16):1781–1802.
- [38] Van Der Spoel D, Lindahl E, Hess B, Groenhof G, Mark AE, Berendsen HJ. GROMACS: fast, flexible, and free. *Journal of computational chemistry*. 2005;26(16):1701–1718.
- [39] The Folding@home Consortium;. <https://foldingathome.org/about/the-foldinghome-consortium/>.
- [40] LeCun Y, Bengio Y, Hinton G. Deep learning. *nature*. 2015;521(7553):436–444.
- [41] Smith JS, Isayev O, Roitberg AE. ANI-1, A data set of 20 million calculated off-equilibrium conformations for organic molecules. *Scientific data*. 2017;4:170193.
- [42] Smith JS, Nebgen BT, Zubatyuk R, Lubbers N, Devereux C, Barros K, et al. Approaching coupled cluster accuracy with a general-purpose neural network potential through transfer learning. *Nature communications*. 2019;10(1):1–8.
- [43] Behler J. Perspective: Machine learning potentials for atomistic simulations. *The Journal of chemical physics*. 2016;145(17):170901.
- [44] Von Lilienfeld OA. Quantum machine learning in chemical compound space. *Angewandte Chemie International Edition*. 2018;57(16):4164–4169.
- [45] Schmitz G, Godtliebsen IH, Christiansen O. Machine learning for potential energy surfaces: An extensive database and assessment of methods. *The Journal of chemical physics*. 2019;150(24):244113.
- [46] Devereux C, Smith JS, Davis KK, Barros K, Zubatyuk R, Isayev O, et al. Extending the applicability of the ANI deep learning molecular potential to Sulfur and Halogens. *Journal of Chemical Theory and Computation*. 2020;16(7):4192–4202.
- [47] Behler J, Parrinello M. Generalized neural-network representation of high-dimensional potential-energy surfaces. *Physical review letters*. 2007;98(14):146401.
- [48] Zubatyuk R, Smith JS, Leszczynski J, Isayev O. Accurate and transferable multitask prediction of chemical properties with an atoms-in-molecules neural network. *Science advances*. 2019;5(8):eaav6490.
- [49] Schütt K, Kindermans PJ, Felix HES, Chmiela S, Tkatchenko A, Müller KR. Schnet: A continuous-filter convolutional neural network for modeling quantum interactions. In: *Advances in neural information processing systems*; 2017. p. 991–1001.
- [50] Schütt KT, Sauceda HE, Kindermans PJ, Tkatchenko A, Müller KR. SchNet—A deep learning architecture for molecules and materials. *The Journal of Chemical Physics*. 2018;148(24):241722.
- [51] Senn HM, Thiel W. QM/MM methods for biomolecular systems. *Angewandte Chemie International Edition*.

- 2009;48(7):1198–1229.
- [52] Lahey SLJ, Rowley CN. Simulating Protein-Ligand Binding with Neural Network Potentials. *Chemical Science*. 2020.
- [53] Bösel L, Thürlmann M, Riniker S. Machine Learning in QM/MM Molecular Dynamics Simulations of Condensed-Phase Systems; 2020.
- [54] Unke OT, Meuwly M. PhysNet: a neural network for predicting energies, forces, dipole moments, and partial charges. *Journal of chemical theory and computation*. 2019;15(6):3678–3693.
- [55] Thomas N, Smidt T, Kearnes S, Yang L, Li L, Kohlhoff K, et al. Tensor field networks: Rotation-and translation-equivariant neural networks for 3d point clouds. *arXiv preprint arXiv:180208219*. 2018.
- [56] Kondor R, Lin Z, Trivedi S. Clebsch–gordan nets: a fully fourier space spherical convolutional neural network. In: *Advances in Neural Information Processing Systems*; 2018. p. 10117–10126.
- [57] Kondor R. N-body networks: a covariant hierarchical neural network architecture for learning atomic potentials. *arXiv preprint arXiv:180301588*. 2018.
- [58] Anderson B, Hy TS, Kondor R. Cormorant: Covariant molecular neural networks. In: *Advances in Neural Information Processing Systems*; 2019. p. 14510–14519.
- [59] Gilmer J, Schoenholz SS, Riley PF, Vinyals O, Dahl GE. Neural message passing for quantum chemistry. In: *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org; 2017. p. 1263–1272.
- [60] Feinberg EN, Sur D, Wu Z, Husic BE, Mai H, Li Y, et al. PotentialNet for molecular property prediction. *ACS central science*. 2018;4(11):1520–1530.
- [61] Sidky H, Chen W, Ferguson AL. Machine learning for collective variable discovery and enhanced sampling in biomolecular simulation. *Molecular Physics*. 2020;118(5):e1737742.
- [62] Wang Y, Ribeiro JML, Tiwary P. Machine learning approaches for analyzing and enhancing molecular dynamics simulations. *Current Opinion in Structural Biology*. 2020;61:139–145.
- [63] Bussi G, Laio A, Parrinello M. Equilibrium free energies from nonequilibrium metadynamics. *Physical review letters*. 2006;96(9):090601.
- [64] Barducci A, Bussi G, Parrinello M. Well-tempered metadynamics: a smoothly converging and tunable free-energy method. *Physical review letters*. 2008;100(2):020603.
- [65] Barducci A, Bonomi M, Parrinello M. Metadynamics. *Wiley Interdisciplinary Reviews: Computational Molecular Science*. 2011;1(5):826–843.
- [66] Wehmeyer C, Noé F. Time-lagged autoencoders: Deep learning of slow collective variables for molecular kinetics. *The Journal of chemical physics*. 2018;148(24):241703.
- [67] Sultan MM, Pande VS. Automated design of collective variables using supervised machine learning. *The Journal of chemical physics*. 2018;149(9):094106.
- [68] Martd A, Pasquali L, Wu H, Noé F. VAMPnets for deep learning of molecular kinetics. *Nature communications*. 2018;9(1):1–11.
- [69] Trapl D, Horvačanin I, Mareška V, Özçelik F, Spiwok V, Unal G. Anncolvar: approximation of complex collective variables by artificial neural networks for analysis and biasing of molecular simulations. *Frontiers in molecular biosciences*. 2019;6:25.
- [70] Pramanik D, Smith Z, Kells A, Tiwary P. Can One Trust Kinetic and Thermodynamic Observables from Biased Metadynamics Simulations?: Detailed Quantitative Benchmarks on Millimolar Drug Fragment Dissociation. *The Journal of Physical Chemistry B*. 2019;123(17):3672–3678.
- [71] Wang Y, Ribeiro JML, Tiwary P. Past–future information bottleneck for sampling molecular reaction coordinate simultaneously with thermodynamics and kinetics. *Nature communications*. 2019;10(1):1–8.
- [72] Maier JA, Martinez C, Kasavajhala K, Wickstrom L, Hauser KE, Simmerling C. ff14SB: improving the accuracy of protein side chain and backbone parameters from ff99SB. *Journal of chemical theory and computation*. 2015;11(8):3696–3713.
- [73] Schütt KT, Arbabzadah F, Chmiela S, Müller KR, Tkatchenko A. Quantum-chemical insights from deep tensor neural networks. *Nature communications*. 2017;8(1):1–8.
- [74] Best RB, Mittal J, Feig M, MacKerell Jr AD. Inclusion of many-body effects in the additive CHARMM protein CMAP potential results in enhanced cooperativity of  $\alpha$ -helix and  $\beta$ -hairpin formation. *Biophysical journal*. 2012;103(5):1045–1051.
- [75] The MolSSI Quantum Chemistry Archive: A central source to compile, aggregate, query, and share quantum chemistry data.; <https://qcarchive.molssi.org/>.
- [76] Folding@home; <https://foldingathome.org/>.
- [77] Qiu Y, Smith D, Boothroyd S, Jang H, Wagner J, Bannan CC, et al. Development and Benchmarking of Open Force Field v1. 0.0, the Parsley Small Molecule Force Field. 2020.

- [78] Qiu Y, Smith DGA, Boothroyd S, Wagner J, Bannan CC, Gokey T, et al.. openforcefield/openforcefields: Version 1.0.0 "Parsley". Zenodo; 2019. doi:10.5281/zenodo.3483227.
- [79] Jorgensen WL, Chandrasekhar J, Madura JD, Impey RW, Klein ML. Comparison of simple potential functions for simulating liquid water. *The Journal of chemical physics*. 1983;79(2):926–935.
- [80] OpenMM GitHub Project Boards, howpublished = <https://github.com/openmm/openmm/projects>;
- [81] The Molecular Sciences Software Institute (MolSSI);. <http://molssi.org>.
- [82] OpenMM GitHub Organization;. <https://github.com/openmm>.
- [83] The OpenMM Consortium;. <http://openmm.org/about.html#consortium>.
- [84] Ponder JW, Wu C, Ren P, Pande VS, Chodera JD, Schnieders MJ, et al. Current status of the AMOEBA polarizable force field. *The journal of physical chemistry B*. 2010;114(8):2549–2564.
- [85] Rackers JA. A Physics-Based Intermolecular Potential for Biomolecular Simulation. Washington University in St. Louis; 2019.
- [86] Hopkins CW, Le Grand S, Walker RC, Roitberg AE. Long-time-step molecular dynamics through hydrogen mass repartitioning. *Journal of chemical theory and computation*. 2015;11(4):1864–1874.
- [87] Leimkuhler B, Matthews C. Robust and efficient configurational molecular sampling via Langevin dynamics. *The Journal of chemical physics*. 2013;138(17):05B601\_1.
- [88] Leimkuhler B, Matthews C. Efficient molecular dynamics using geodesic integration and solvent-solute splitting. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*. 2016;472(2189):20160138.
- [89] Åqvist J, Wennerström P, Nervall M, Bjelic S, Brandsdal BO. Molecular dynamics simulations of water and biomolecules with a Monte Carlo constant pressure algorithm. *Chemical physics letters*. 2004;384(4-6):288–294.
- [90] Huang B, Symonds NO, von Lilienfeld OA. Quantum machine learning in chemistry and materials. *Handbook of Materials Modeling: Methods: Theory and Modeling*. 2020:1883–1909.
- [91] Abadi M, Agarwal A, Barham P, Brevdo E, Chen Z, Citro C, et al.. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems; 2015. Available from: <https://www.tensorflow.org/>.
- [92] Paszke A, Gross S, Massa F, Lerer A, Bradbury J, Chanan G, et al. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In: Wallach H, Larochelle H, Beygelzimer A, d'Alché-Buc F, Fox E, Garnett R, editors. *Advances in Neural Information Processing Systems 32*. Curran Associates, Inc.; 2019. p. 8024–8035. Available from: <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- [93] Bradbury J, Frostig R, Hawkins P, Johnson MJ, Leary C, Maclaurin D, et al.. JAX: composable transformations of Python+NumPy programs; 2018. Available from: <http://github.com/google/jax>.
- [94] The OpenMM TensorFlow plugin;. <http://github.com/openmm/openmm-nn>.
- [95] The OpenMM TensorRT plugin;. <http://github.com/openmm/openmm-tensorrt>.
- [96] TorchANI;. <https://github.com/aiqm/torchani>.
- [97] Ramsundar B, Eastman P, Walters P, Pande V. Deep learning for the life sciences: applying deep learning to genomics, microscopy, drug discovery, and more. " O'Reilly Media, Inc."; 2019.
- [98] Parks CD, Gaieb Z, et al. D3R Grand Challenge 4: blind prediction of protein-ligand poses, affinity rankings, and relative binding free energies. *Journal of computer-aided molecular design*. 2020.
- [99] Piquemal JP, Perera L, Cisneros GA, Ren P, Pedersen LG, Darden TA. Towards accurate solvation dynamics of divalent cations in water using the polarizable amoeba force field: From energetics to structure. *The Journal of chemical physics*. 2006;125(5):054511.
- [100] Shi Y, Xia Z, Zhang J, Best R, Wu C, Ponder JW, et al. Polarizable atomic multipole-based AMOEBA force field for proteins. *Journal of chemical theory and computation*. 2013;9(9):4046–4063.
- [101] Laury ML, Wang LP, Pande VS, Head-Gordon T, Ponder JW. Revised parameters for the AMOEBA polarizable atomic multipole water model. *The Journal of Physical Chemistry B*. 2015;119(29):9423–9437.
- [102] Baker CM, Lopes PE, Zhu X, Roux B, MacKerell Jr AD. Accurate calculation of hydration free energies using pair-specific Lennard-Jones parameters in the CHARMM Drude polarizable force field. *Journal of chemical theory and computation*. 2010;6(4):1181–1198.
- [103] Cisneros GA. Application of Gaussian electrostatic model (GEM) distributed multipoles in the AMOEBA force field. *Journal of chemical theory and computation*. 2012;8(12):5072–5080.
- [104] Hinton G, Vinyals O, Dean J. Distilling the knowledge in a neural network. *arXiv preprint arXiv:150302531*. 2015.
- [105] Cheng Y, Wang D, Zhou P, Zhang T. A survey of model compression and acceleration for deep neural networks. *arXiv preprint arXiv:171009282*. 2017.
- [106] Zhou A, Yao A, Guo Y, Xu L, Chen Y. Incremental network quantization: Towards lossless cnns with

- low-precision weights. arXiv preprint arXiv:170203044. 2017.
- [107] Noé F, Olsson S, Köhler J, Wu H. Boltzmann generators: Sampling equilibrium states of many-body systems with deep learning. *Science*. 2019;365(6457):eaaw1147.
- [108] Schoenholz SS, Cubuk ED. JAX M.D.: End-to-End Differentiable, Hardware Accelerated, Molecular Dynamics in Pure Python; 2019. <https://github.com/google/jax-md>, <https://arxiv.org/abs/1912.04232>.
- [109] TimeMachine: A simplified and differentiable MD engine;. <https://github.com/proteneer/timemachine>.
- [110] Wang J, Wolf RM, Caldwell JW, Kollman PA, Case DA. Development and testing of a general amber force field. *Journal of computational chemistry*. 2004;25(9):1157–1174.
- [111] The Open Force Field Initiative: An open source, open science, and open data approach to better force fields;. <https://openforcefield.org>.
- [112] The MolSSI QCArchive Machine Learning Datasets Repository;. [https://qcarchive.molssi.org/apps/ml\\_datasets/](https://qcarchive.molssi.org/apps/ml_datasets/).
- [113] Model Zoo: Discover open source deep learning code and pretrained models;. <https://modelzoo.co/>.
- [114] Fey M, Lenssen JE. Fast Graph Representation Learning with PyTorch Geometric. In: *ICLR Workshop on Representation Learning on Graphs and Manifolds*; 2019. .
- [115] Parrish RM, Burns LA, Smith DG, Simmonett AC, DePrince III AE, Hohenstein EG, et al. Psi4 1.1: An open-source electronic structure program emphasizing automation, advanced libraries, and interoperability. *Journal of chemical theory and computation*. 2017;13(7):3185–3197.
- [116] QCEngine: A quantum chemistry program executor and IO standardizer (QCSchema) for quantum chemistry;. <https://github.com/MolSSI/QCEngine>.
- [117] Feng Z, Chen L, Maddula H, Akcan O, Oughtred R, Berman HM, et al. Ligand Depot: a data warehouse for ligands bound to macromolecules. *Bioinformatics*. 2004;20(13):2153–2155.
- [118] Enamine REAL;. <https://enamine.net/library-synthesis/real-compounds/real-compound-libraries>.
- [119] Pitt WR, Parry DM, Perry BG, Groom CR. Heteroaromatic rings of the future. *Journal of medicinal chemistry*. 2009;52(9):2952–2963.
- [120] Gebauer N, Gastegger M, Schütt K. Symmetry-adapted generation of 3d point sets for the targeted discovery of molecules. In: Wallach H, Larochelle H, Beygelzimer A, d'Alché-Buc F, Fox E, Garnett R, editors. *Advances in Neural Information Processing Systems 32*. Curran Associates, Inc.; 2019. p. 7564–7576. Available from: <http://papers.nips.cc/paper/8974-symmetry-adapted-generation-of-3d-point-sets-for-the-targeted-discovery-of-molecules.pdf>.
- [121] Mobley DL, Bannan CC, Rizzi A, Bayly CI, Chodera JD, Lim VT, et al. Escaping atom types in force fields using direct chemical perception. *Journal of chemical theory and computation*. 2018;14(11):6076–6092.
- [122] Quantum Chemistry Schema: A JSON Schema for Quantum Chemistry;. <https://molssi-qc-schema.readthedocs.io/>.