# Machine Learning Directed Optimization of Classical Molecular Modeling Force Fields

Bridgette J. Befort,* Ryan S. DeFever,* Garrett M. Tow, Alexander W. Dowling, and Edward J. Maginn†

*Department of Chemical and Biomolecular Engineering,*
*University of Notre Dame, Notre Dame, Indiana 46556, United States*
(Dated: July 15, 2021)

Accurate force fields are necessary for predictive molecular simulations. However, developing force fields that accurately reproduce experimental properties is challenging. Here, we present a machine learning directed, multiobjective optimization workflow for force field parameterization that evaluates millions of prospective force field parameter sets while requiring only a small fraction of them to be tested with molecular simulations. We demonstrate the generality of the approach and identify multiple low-error parameter sets for two distinct test cases: simulations of hydrofluorocarbon (HFC) vapor–liquid equilibrium (VLE) and an ammonium perchlorate (AP) crystal phase. We discuss the challenges and implications of our force field optimization workflow.

## I. INTRODUCTION

Molecular modeling and simulation use computational methods to describe the behavior of matter at the atomistic or molecular level [1]. The veracity and predictive capability of molecular simulations depend critically on the accuracy of the atomic-level interaction energies, and whether the appropriate time- and length-scales are properly sampled. On one hand is a class of techniques broadly termed as *ab initio* or first-principles methods, where atomic interactions are determined from highly accurate quantum chemical methods [2]. Though there are applications that necessitate these methods, *ab initio* energies are computationally expensive to obtain, such that quantum chemical methods are limited to relatively small systems and short timescales. On the other hand, classical molecular simulations represent the atomic interaction energies with an analytical function (a "force field") that can be evaluated much more rapidly than *ab initio* energy, enabling simulations of much larger systems and longer timescales than is possible with *ab initio* techniques. If force fields are highly accurate, classical molecular simulations have been shown to give accurate property predictions in several fields including protein structure refinement [3], drug discovery [4], and energy storage [5].

### A. Developing Accurate Force Fields is Difficult

There are two fundamentally different approaches to developing and improving force fields: bottom-up approaches, wherein parameters are calibrated so the model reproduces the results (e.g., forces, energies, and dipoles) of more expensive and accurate methods (i.e., quantum calculations) [6], and top-down approaches, wherein parameters are calibrated so the model matches experimental results [7]. Emerging bottom-up approaches use

machine learning (ML) to parameterize force fields with black-box potential energy functions[8, 9]. Though these so-called ML force fields[10, 11] have proven successful for an increasing number of systems, the black-box nature of the potential energy function makes the models physically uninterpretable, and hinders model transferability beyond the specific training conditions. Developing accurate and transferable force fields with analytical functional forms is a difficult and laborious endeavor [12]. Significant efforts spanning several decades have resulted in several "off-the-shelf" force fields that describe large swaths of condensed matter chemical space [13–16]. These are most commonly "Class I" force fields that consist of harmonic or sinusoidal intramolecular terms that describe bonded interactions, atomic partial charges that represent electrostatic interactions, and nonbonded repulsion-dispersion terms. Unfortunately, these off-the-shelf force fields can yield poor property predictions, even for relatively common compounds, particularly when they are applied in circumstances beyond the systems and conditions for which they were parameterized [17]. However, since they are well known and the parameter sets are widely distributed, these force fields are used in many molecular simulation studies.

For decades, force field development and optimization has been an active area of research. Several methods and tools have been developed to derive bonded intramolecular parameters and partial charges in a bottom-up fashion from quantum calculations, provided that the desired classical functional form has been selected. Common approaches include gradient-based techniques, evolutionary algorithms, or even analytical solutions [18–25]. These methods work well because the relevant quantities can be computed to a high degree of accuracy with quantum calculations, and evaluating a prospective force field parameter set is computationally trivial. However, optimizing the repulsion-dispersion parameters that are largely responsible for many macroscopic thermodynamic properties (e.g., density, enthalpy of vaporization, vapor pressure, etc.) is more challenging. Since these parameters can be difficult to derive from quantum calculations without special methods [26], top-down parameterization is often necessary. Yet screening thousands of prospective

---

* BJ Befort and RS DeFever contributed equally to this work.
† Corresponding author

parameter sets is computationally expensive due to the need for sufficiently long simulations to accurately compute the relevant experimental properties. Even for relatively simple properties, a single simulation can require hours-to-days of computation time.

It is often desirable to parameterize a force field to reproduce multiple physical properties. A rigorous way to calibrate force fields with multiple properties simultaneously is to use multiobjective optimization[27–31], which can exacerbate the computational burden by an order of magnitude or more. In multiobjective optimization, a solution is Pareto optimal if it is not possible to improve one objective without sacrificing another objective.[32] One approach is to weight each objective and re-solve the optimization problem for many different weights to identify Pareto optimal solutions.[33] Thus computing a set of Pareto optimal solutions is often at least an order of magnitude more computationally expensive than single objective optimization. With much less computational effort, a finite set of candidate solutions can be classified into two groups: the non-dominated set, which comprises the solutions for which no other solution in the set offers improvement in any one objective without degrading performance in another objective, and the dominated set, comprising the solutions for which another solution offers improved performance in one or more objectives without degrading the performance in any other objective. By definition, all points in the Pareto set are non-dominated; the non-dominated set is an easy to compute approximation of the Pareto set.

Given the challenges associated with top-down optimization of the repulsion-dispersion parameters, there are fewer methods and packages available [34, 35] compared to intramolecular parameters and partial charge optimization. Much more frequently, attempts to improve these parameters involve *ad hoc* hand-tuning [36, 37], which is arbitrary and often limited to a few interaction parameters or a scaling thereof, as larger searches quickly become intractable [38]. Instead of performing multiobjective optimization, the more common approach is to use *ad hoc* weights to combine multiple calibration objectives into a single cost function [28, 34, 35]. However, this approach only finds a single Pareto optimal trade-off between the calibration objectives.

## B. Machine Learning Directed Optimization Makes Force Field Calibration More Computationally Tractable

The core challenges of optimizing the repulsion-dispersion parameters can be solved with a computationally inexpensive mapping between the desired physical properties and force field parameters. For certain cases, these mappings can be constructed with statistical mechanics [39, 40], but this approach likely cannot be generalized to arbitrary systems. Alternatively, ML can be used to approximate the relevant mapping. For example, surrogate-assisted optimization (also known as black-box or derivative-free optimization) uses computationally inexpensive surrogate model evaluations to emulate the outputs of a complex computer simulation, e.g., computational fluid dynamics, finite element analysis, or molecular simulations. Several different types of surrogate models have been successfully applied to molecular simulations for uncertainty quantification [41, 42] and force field parameterization [35, 43–45]. Linear regression response surface models were used to predict the optimal combination of scaling factors for the charge and Lennard-Jones (LJ) parameters of General AMBER force field (GAFF) to reproduce four properties of organic liquid electrolytes. While easy to implement and moderately successful at improving the force field's accuracy for some of the properties, this method was limited by the choice of statistically significant parameters in the response surface. [46] For some thermodynamic properties, reweighting methods are an effective tool to test a large number of parameters without performing additional simulations [44, 47, 48], but care must be taken to ensure good phase space overlap between the sampled and reweighted ensembles. [44] Gaussian process regression (GPR) is a popular non-parametric surrogate model that smoothly interpolates between training data. Some applications of GPR in molecular simulations include ML force fields [49–51] and property prediction [52]. In Bayesian optimization, which is a special case of surrogate-assisted optimization, the uncertainty estimates from GPR (or a similar model) are directly used to balance exploration and exploitation. Recent work demonstrates Bayesian optimization can efficiently calibrate force field parameters in coarse-grained models [53–55]. Moreover, computationally inexpensive surrogate models can enable multiobjective optimization algorithms that go beyond *ad hoc* weighting [32] to systematically explore trade-offs when calibrating multiple physical properties.

Here, we demonstrate a new multiobjective surrogate-assisted optimization framework that uses GPRs and support vector machine (SVM) classifiers to improve existing all-atom force fields. The proposed strategy enables extremely accurate property calculations while retaining physically-motivated and interpretable functional forms. We show that the same general approach successfully optimizes force fields for two systems with very different characteristics and property objectives: hydrofluorocarbon (HFC) vapor–liquid equilibrium (VLE) and solid ammonium perchlorate (AP) crystal structure. Our results highlight the versatility of surrogate-assisted optimization approaches for top-down parameterization of all-atom force fields in a wide range of domains. The remainder of the manuscript proceeds as follows: we outline the method and provide technical details in Section II, demonstrate the approach for the two case studies in Section III, discuss the challenges and implications of the method in Section IV, and provide concluding remarks in Section V.

## II. METHODOLOGY

### A. A Machine Learning Directed Force Field Optimization Workflow

An overview of our force field optimization workflow is provided first with a more technical description given in the following subsections. Our strategy in this work is to optimize LJ repulsion-dispersion parameters, which are among the most difficult to calculate from *ab initio* methods [56]. Intramolecular parameters and partial charges, which usually can be reliably and inexpensively determined from bottom-up *ab initio*-based methods, were determined from existing force fields. We stress, however, that this method can be applied to calibrate any force field parameters.

Our force field optimization workflow is shown schematically in Figure 1. First, domain knowledge is used to specify physically reasonable bounds on the search space for the parameters that are being optimized. Next, $\mathcal{O}(10^2)$ initial parameter sets are generated via space-filling Latin hypercube sampling (LHS). Molecular simulations are performed with each parameter set (Figure 1, box 1), and the physical properties of interest are computed from the simulations. These results are used to train surrogate models (box 2, panel d) that predict the simulation results directly from the parameter set, and optionally, the thermodynamic state point, e.g., $T$ and $p$. Additional examples of surrogate model accuracy can be found in SI Figures S1 and S2. The surrogate model is then used to predict the molecular simulation results for a very large number, $\mathcal{O}(10^6)$, of candidate parameter sets, once again generated with LHS (box 3). The $\mathcal{O}(10^2)$ most promising parameter sets are identified via user-selected system-specific metrics including error thresholds, separation in parameter space, and non-dominated status, from the $\mathcal{O}(10^6)$ candidate sets evaluated with the surrogate models (box 4). In multiobjective optimization, the set of non-dominated points includes all parameter sets that are not simultaneously outperformed in every dimension by any other parameter set (Figure 1a) [32]. Finally, the most promising parameter sets are used to initialize the next iteration of molecular simulations (box 1). The process is repeated until parameter sets are generated that provide the desired accuracy for the experimental properties of interest.

The workflow uses a combination of machine learning-based surrogate models and physics-based molecular simulations to quickly optimize force field parameters for a specific system. Physically-motivated potential energy functional forms that have proven successful over decades are retained. Whereas the molecular simulations require hours-to-days to compute experimentally measurable properties arising from a single set of force field parameters, the surrogate models can evaluate millions of parameter sets in minutes-to-hours. This means that once the surrogate models have been trained to predict the results of the molecular simulations, they enable an exhaustive search of large parameter spaces that would require $\mathcal{O}(10^7\text{--}10^9)$ CPU-hours with molecular simulations. We emphasize that although the surrogate models are used to screen millions of candidate parameter sets, all of the promising candidate parameter sets are ultimately tested with physics-based molecular simulations. The role of machine learning is only to act as a surrogate for physics-based simulations, enabling the parameter search through an otherwise intractable space. The iterative procedure allows the surrogate models to improve as additional training data is collected with each iteration. The original molecular simulations are dispersed across the entire parameter space, but subsequent iterations are focused on the smaller regions of parameter space that are predicted to yield good parameter sets, enabling the surrogate models to improve in the most important regions of parameter space. The theory and technical details of each step in Figure 1 are presented in Sections II A 1 to II A 5. Methodological details specific to the HFC and AP examples are reported in Sections II B and II C, respectively.

#### 1. Problem Setup

The interaction potential is taken as a classical molecular mechanics force field, $U(\mathbf{r}) = f(\mathbf{r}, \boldsymbol{\zeta})$, where $U$ is the potential energy, $\mathbf{r} \in \boldsymbol{\Gamma}$ is the vector of position coordinates within configuration space $\boldsymbol{\Gamma}$, $f$ is the functional form for the potential energy, and $\boldsymbol{\zeta} = \zeta_1, \zeta_2, ..., \zeta_N$ are the parameters of $f$ that define the intra- and intermolecular interactions between different types of particles. Molecular simulations can be used to compute $M$ structural, thermodynamic, or dynamic properties, $\mathbf{y}^{\text{sim}} = y_1^{\text{sim}}, y_2^{\text{sim}}, ..., y_M^{\text{sim}}$, from $U(\mathbf{r})$. Depending upon the quality of $U(\mathbf{r})$, $\mathbf{y}^{\text{sim}}$ may or may not be close to the experimental values, $\mathbf{y}^{\text{exp}}$. The goal of this work is to refine $U(\mathbf{r})$ by optimizing $\mathcal{O}(10^1)$ force field parameters, $\boldsymbol{\zeta}' \subseteq \boldsymbol{\zeta}$, such that $\mathbf{y}^{\text{sim}} \approx \mathbf{y}^{\text{exp}}$ for one or more physical properties of interest. In both case studies presented here, the LJ parameters, $\sigma$ and $\varepsilon$, are optimized. Upper and lower bounds for each parameter are selected to span a wide range of physically reasonable values. The initial $\mathcal{O}(10^2)$ parameter sets are randomly selected to be space-filling within these bounds with LHS.

#### 2. Step 1: Perform Molecular Simulations with $\mathcal{O}(10^2)$ Physics-Based Force Fields

Molecular simulations are performed for each parameter set with the molecular dynamics (MD) or Monte Carlo (MC) method. For each parameter set, $\mathbf{y}^{\text{sim}}$ is computed from the simulation output. Simulations may be performed at multiple thermodynamic conditions (e.g., $T$ and $p$) for each parameter set if the experimental data exist. Signac-flow was used to manage the setup and execution of all molecular simulations [57, 58].
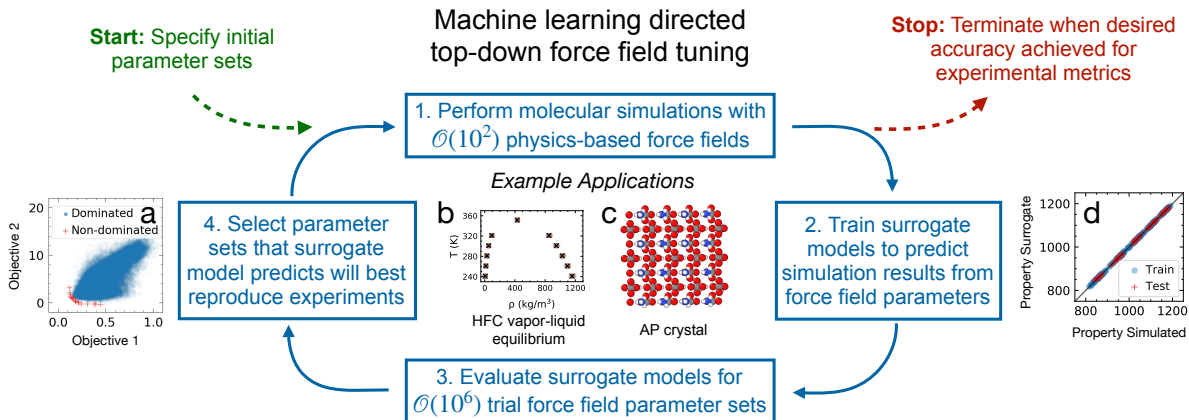
FIG. 1. Overview of the proposed machine learning directed force field optimization procedure. The workflow tests $\mathcal{O}(10^6)$ sets of force field parameters for every $\mathcal{O}(10^2)$ molecular simulations. The four main steps are described in the numbered boxes. Panel (a) shows the difference between dominated and non-dominated solutions for an example where the goal is to minimize two objectives. Panels (b) and (c) highlight the two example applications. Panel (d) shows an example of how the surrogate models accurately predict the outcomes of molecular simulations.

### 3. Step 2: Train Surrogate Models to Predict Simulation Results from Force Field Parameters

Gaussian process (GP) surrogate models are trained to predict $\mathbf{y}^{\text{sim}}$ as a function of the calibrated parameters $\boldsymbol{\zeta}'$. For each property, we train:

$$\hat{y}_i^{\text{sim}} = GP_i(m_i(\boldsymbol{\zeta}'), \text{cov}_i(\boldsymbol{\zeta}', \boldsymbol{\zeta}')) \tag{1}$$

where $\hat{y}_i^{\text{sim}}$ is the surrogate model prediction of $y_i^{\text{sim}}$, $GP_i$ is the GP model for property $i$, $m_i$ is the mean function, and $\text{cov}_i$ is the covariance (kernel) function. All GP models were implemented in GPFlow 2.0.0 [59]. To improve the accuracy of the GP models in regions of parameter space where $\mathbf{y}^{\text{sim}} \approx \mathbf{y}^{\text{exp}}$, we exclude parameter sets that result in extremely poor or unphysical results from the GP training data. We then trained SVM classifiers to predict if a parameter set was unphysical (e.g., simulation fails) so that parameter sets from these regions of parameter space could be excluded when the GP models were used to predict the results of trial parameter sets. All SVM classifiers were implemented in scikit-learn [60] with a radial basis function kernel.

### 4. Step 3: Evaluate Surrogate Models for $\mathcal{O}(10^6)$ Trial Force Field Parameter Sets

After the GP and SVM models are trained, $\mathcal{O}(10^6)$ trial parameter sets are generated with LHS. For each parameter set, the SVM and GP models are used to calculate $\hat{\mathbf{y}}^{\text{sim}}$, the surrogate model estimates of $\mathbf{y}^{\text{sim}}$.

### 5. Step 4: Select Parameter Sets that Surrogate Models Predict Will Best Reproduce Experiments

Parameter sets where the surrogate models predict good agreement with experiment, $\hat{\mathbf{y}}^{\text{sim}} \approx \mathbf{y}^{\text{exp}}$, are selected for the next iteration. In some cases we apply an optional distance-based search algorithm (see SI Methods) to down-select only parameter sets that are far apart in parameter space.

### B. Hydrofluorocarbon Case Study

Force fields were independently developed for two HFCs: difluoromethane (HFC-32) and pentafluoroethane (HFC-125). Two stages of optimization were used for each HFC. The first stage used MD simulations in the $NpT$ ensemble at: 241, 261, 281, 301, and 321 K for HFC-32 and 229, 249, 269, 289, and 309 K for HFC-125. For each temperature, the pressure was set to the experimental[61] saturation pressure. The only property considered during the first stage was the liquid density (LD) ($\mathbf{y} = \{\rho^l\}$). In the second stage of optimization, Gibbs ensemble Monte Carlo (GEMC) was performed. The property objectives were the saturated liquid density, saturated vapor density, vapor pressure, and enthalpy of vaporization, or $\mathbf{y} = \{\rho_{\text{sat}}^l, \rho_{\text{sat}}^v, P_{\text{vap}}, \Delta H_{\text{vap}}\}$. Simulations were performed at the same temperatures used for the first stage. Four iterations of the stage 1 optimization were performed for both HFC-32 and HFC-125. Three and five iterations of stage 2 optimization were performed for HFC-32 and HFC-125, respectively.

### 1. Force Field Parameters

The functional form was taken from GAFF [15]:

$$U(\mathbf{r}) = U^{\text{intra}}(\mathbf{r}) + \sum_i \sum_{j>i} \frac{q_i q_j}{4\pi\epsilon_0 r_{ij}}$$
$$+ \sum_i \sum_{j>i} 4\varepsilon_{ij} \left[ \left(\frac{\sigma_{ij}}{r_{ij}}\right)^{12} - \left(\frac{\sigma_{ij}}{r_{ij}}\right)^6 \right] \quad (2)$$

where $U^{\text{intra}}$ contains all the intramolecular terms, $r_{ij}$ is the distance between atoms $i$ and $j$, $q$ is the atomic charge, $\epsilon_0$ is the permittivity of free space, and $\sigma_{ij}$ and $\varepsilon_{ij}$ parametrize the LJ potential that describes the repulsion-dispersion interactions between atoms $i$ and $j$. The intramolecular interactions are given by:

$$U^{\text{intra}}(\mathbf{r}) = \sum_{\text{bonds}} k_r (r - r_0)^2 + \sum_{\text{angles}} k_\theta (\theta - \theta_0)^2$$
$$+ \sum_{\text{dihedrals}} \nu_n \left[ 1 + \cos(n\phi - \gamma) \right] \quad (3)$$

where $r_0$ and $\theta_0$ are the nominal bond length and angle, respectively, $k_r$, $k_\theta$, and $\nu_n$ are force constants, $n$ is the multiplicity and $\gamma$ is the nominal dihedral angle. The sums are over all bonds, angles, and dihedrals in the system, respectively. The bond, angle, and dihedral parameter for HFC-32 and HFC-125 were taken from GAFF [15]. Partial charges were determined with RESP [18] as implemented in AmberTools 1.4 [62]. The quantum electrostatic potential was computed with Gaussian 09 [63] with B3LYP/6-311++g(d,p) [64, 65]. The intramolecular parameters and partial charges are reported in SI Table S1.

The force field optimization method was used to determine the like-interaction parameters $\sigma_{ii}$ and $\varepsilon_{ii}$ for three atom types (C, F, and H) in HFC-32 and five atom types (C1, C2, F1, F2, and H) in HFC-125. This results in 6 parameters that are optimized for HFC-32 and 10 parameters that are optimized for HFC-125. All unlike interaction parameters were computed with Lorentz–Berthelot mixing rules. For HFC-125, C1 is the carbon bonded to one carbon atom, two fluorine atoms, and one hydrogen atom, while C2 is the carbon bonded to one carbon atom and three fluorine atoms, F1 is bonded with to C1, and F2 is bonded with C2. The lower and upper bounds for each parameter were selected per-element ($\sigma$ in Å, $\varepsilon/k_B$ in K): $3.0 \leq \sigma_C \leq 4.0$, $2.5 \leq \sigma_F \leq 3.5$, $1.7 \leq \sigma_H \leq 2.7$, $20 \leq \varepsilon_C/k_B \leq 60$, $15 \leq \varepsilon_F/k_B \leq 40$, $2 \leq \varepsilon_H/k_B \leq 10$. The parameter bounds for each atom type in HFC-32 and HFC-125 are summarized in SI Tables S2 and S3, respectively.

### 2. Classifier

An SVM classifier was trained to predict parameter sets that yielded spontaneous vaporization ($\rho^l <$ 500 kg/m$^3$) in MD simulations initiated at liquid density from $\boldsymbol{\zeta}'$ and $T$.

### 3. GP Model

The GP models predicted the value of a physical property from $\boldsymbol{\zeta}'$ and $T$. The LD iterations used one GP model that predicted $\rho^l$. Parameter sets with $\rho^l <$ 500 kg/m$^3$ were excluded from the GP training data. The VLE iterations used one GP model for each property: $\{\rho^l_{\text{sat}}, \rho^v_{\text{sat}}, P_{\text{vap}}, \Delta H_{\text{vap}}\}$. All GP models used a radial basis function or Matérn $\nu = 5/2$ kernel and a linear mean function [66].

### 4. Selecting Parameter Sets for the Next Iteration

A new LHS with 1,000,000 (HFC-32) or 500,000 (HFC-125) parameter sets was generated for each iteration. **LD iterations:** Each parameter set was evaluated with the LD SVM classifier at the highest $T$. Each parameter set was evaluated with the LD GP model at each $T$, and the root mean square error (RMSE) between the GP model prediction and experimental liquid density across all five temperatures was calculated for each parameter set. The 100 lowest RMSE parameter sets that the SVM predicted would remain liquid, and the 100 lowest RMSE parameter sets that the SVM predicted would transform to vapor, were selected for the next iteration. The low-RMSE, predicted-vapor parameter sets were included because they reflect disagreement between the SVM and GP models. After four LD iterations, parameter sets for the VLE-1 iteration were selected from the 800 simulated parameter sets. A distance-based search algorithm (see SI Methods) was used to select 25 well-separated parameter sets with RMSE $\leq 10$ kg/m$^3$. **VLE iterations:** Each parameter set from the LHS was evaluated with the LD GP model. Parameter sets predicted to yield LD RMSE $> 25$ kg/m$^3$ were discarded. This step was included to make use of the training data generated during the LD iterations since the LD GP model is very accurate after four LD iterations. The remaining parameter sets were evaluated with the four GP models trained to predict VLE properties ($\rho^l_{\text{sat}}, \rho^v_{\text{sat}}, P_{\text{vap}}, \Delta H_{\text{vap}}$). The RMSE difference between the GP model predictions and experimental values across all five temperatures was calculated for each property and parameter set. All dominated parameter sets were discarded. A parameter set is dominated if one or more parameter sets performs better than it in all of the considered objective dimensions (e.g., physical properties). The 25 parameter sets selected for the next iteration comprised the top performing parameter set for each physical property and 21 parameter sets selected from the remaining non-dominated parameter sets. A distance-based search algorithm identified parameter sets that were well-separated in parameter space.

### 5. MD Simulations

Simulations of 150 HFC molecules were performed in the $NpT$ ensemble at the experimental saturation pressure. Initial configurations were generated at 1000 kg/m$^3$. Following a steepest descent energy minimization, systems were equilibrated for 500 ps with the Bussi thermostat [67] and Berendsen barostat [68] with $\tau_T = 0.1$ ps, $\tau_p = 0.5$ ps. The production simulations were 2.5 ns in length with the Bussi thermostat and Parrinello–Rahman barostat [69] with $\tau_T = 0.5$ ps and $\tau_p = 1.0$ ps. The final 2.0 ns of the production simulations were used to compute the average density.

The equations of motion were integrated with the leapfrog algorithm[70] and a time step of 1.0 fs. LJ interactions and short range electrostatics were cut off at 1.0 nm. The particle mesh Ewald method[71] was used to compute long-range electrostatic interactions. Analytical tail corrections to the LJ potential were applied to energy and pressure. All bonds were constrained with the P-LINCS [72] method with the lincs-order and lincs-iter set to 8 and 4, respectively. Simulations were performed with GROMACS 2020 [73].

### 6. MC Simulations

GEMC simulations were performed with 1000 HFC molecules. The initial liquid box (800 HFC molecules) was generated at the experimental liquid density and pre-equilibrated with a 5000 sweep $NpT$ MC simulation. The initial vapor box (200 HFC molecules) was randomly generated at the vapor density estimated from the ideal gas law. The combined system was simulated with GEMC. The systems were equilibrated for 10,000 MC sweeps followed by a production GEMC simulation was 90,000 MC sweeps.

LJ interactions and short range electrostatics were cut off at 1.2 nm in the liquid box and 2.5 nm in the vapor box. Long-range electrostatics were computed with an Ewald summation with a relative accuracy of $10^{-5}$. Analytical tail corrections to the LJ interactions were applied to energy and pressure. All bonds were fixed at their nominal bond length. Simulations were performed with MoSDeF Cassandra 0.1.1 [74] and Cassandra 1.2.2 [75].

## C. Ammonium Perchlorate Case Study

Simulations of AP were performed at 1 atm and 10, 78, and 298 K. Three properties were considered: (1) the absolute percent error (APE) from the experimental lattice parameters averaged across all three temperatures, i.e. the mean absolute percent error (MAPE), and (2) the mean of the absolute residuals of equilibrium average simulated atomic positions in reference to the experimental unit cell[76] at 10 K, subsequently referred to as unit cell mean distance (UCMD), and (3) hydrogen-bonding symmetry that is present in the experimental crystal structure. Four workflow iterations were performed.

### 1. Force Field Parameters

The Class II force field of Zhu et al.[77] served as a basis for the development of a hand-tuned Class I force field. The partial charges were left unchanged [78]. The Class II intramolecular bonds and angles were recast to the Class I harmonic functional forms; this process was *ad hoc* and involved qualitative matching to the experimental infrared spectrum. The most significant outcome of this procedure was that at 298 K the N—H stretching mode split into two separate peaks for the Class I force field, as opposed to the single peak observed by both experiment and the Class II force field. This is likely due to inherent limitations in the harmonic representation of the vibrational mode; in the context of our work, this trade-off in vibrational behavior for the simplicity and transferability of the Class I AP force field is acceptable. The LJ parameters of the hand-tuned force field were also developed with an *ad hoc* approach, using similar structural metrics as described above. The hand-tuned AP force field parameters are reported in SI Table S4.

The force field optimization workflow was applied to further optimize the $\sigma$ and $\varepsilon$ for the 4 unique atom types in the AP model, giving a total of 8 calibrated parameters. The lower and upper bounds for each parameter were as follows ($\sigma$ in Å, $\varepsilon$ in kcal/mol): $3.5 \leq \sigma_{Cl} \leq 4.5$, $0.5 \leq \sigma_H \leq 2.0$, $2.5 \leq \sigma_N \leq 3.8$, $2.5 \leq \sigma_O \leq 3.8$, $0.1 \leq \varepsilon_{Cl} \leq 0.8$, $0.0 \leq \varepsilon_H \leq 0.02$, $0.01 \leq \varepsilon_N \leq 0.2$, $0.02 \leq \varepsilon_O \leq 0.3$. The parameter bounds are also summarized in SI Table S5. All unlike LJ interactions were calculated with geometric mixing rules.

### 2. Property Calculation Details

In an effort to be more consistent with the refined hydrogen positions described by Choi et al.[76], the hydrogen atoms in the primitive cell were extended along their N—H vectors to match the N—H lengths that they report in Table V. To assess the symmetry that should be present in orthorhombic AP's *Pnma* space group, the differences in the N—H(3)···O(3) mirror symmetric bond lengths and angles were computed. Hydrogen bonds within 0.001 Å and angles within 0.3° were considered symmetric. To determine tolerances for assessing symmetry, the manually tuned force field was utilized and the frequency of saving coordinate data over the 100 ps production run was varied between 100–10,000 fs. When saving the coordinates every 100 fs, the symmetric hydrogen bond lengths were within 0.00003 Å and the angles were within 0.01° of each other. When saving the coordinates every 10,000 fs, the resolution of symmetry decreases to within 0.001 Å for bonds and 0.3° for an-

gles. For data management reasons, the coordinates were saved every 10,000 fs and the corresponding symmetry tolerances were utilized in classifying if a given parameter set was successful in reproducing the experimentally observed symmetry in the hydrogen bonding structure of AP.

### 3. Classifier

Two SVM classifiers were trained. The first classifier predicted whether a parameter set would yield an accurate 10 K unit cell with UCMD < 0.8 Å, and the second classifier predicted whether a parameter set would yield the desired hydrogen bond symmetry, as defined above.

### 4. GP Model

Two GP surrogate models were trained. The first GP model predicted the 10 K UCMD from $\zeta'$. Parameter sets with UCMD $\geq$ 0.8 Å were not included in the training data. The second GP model predicted the APE of the lattice parameters from $\zeta'$ and $T$. Both GP models used a Matérn $\nu = 3/2$ kernel and a linear mean function [66].

### 5. Selecting Parameter Sets for the Next Iteration

1,000,000 new parameter sets were generated using LHS for each iteration. Each parameter set was evaluated with the UCMD and symmetry classifiers. Parameter sets that did not meet the UCMD threshold were discarded. The remaining parameter sets were evaluated with the two GP models. The lattice APE GP model was evaluated at $T = 10$, 78, and 298 K for each parameter set. The mean of the lattice parameter APE at each temperature was calculated and recorded as the lattice MAPE. All parameter sets that did not meet the UCMD and lattice MAPE thresholds listed in the SI Table S6 were discarded. When selecting parameter sets for the fourth iteration, the symmetry SVM was used to remove all parameter sets that did not meet the symmetry threshold (SI Table S6). A total of 250 parameter sets were selected for the next iteration. All non-dominated parameter sets were selected. The remainder of the parameter sets were selected by applying an $L_1$ distance metric in scaled parameter space and the distance-based search to identify well-separated parameter sets.

### 6. MD Simulations

Simulations of orthorhombic AP were performed in the $NpT$ ensemble at 1 atm and 10, 78, and 298 K. The AP structure was taken from the 10 K data of Choi et al. [76] The simulation cell comprised 378 ($6 \times 9 \times 7$)

unit cells. Initial velocities were drawn from a Gaussian distribution with the linear and angular momenta set to zero. A 1.0 fs time step was utilized with the time integration scheme derived by Tuckerman et al. [79] The equations of motions were those of Shinoda et al. [80] Nosé–Hoover style algorithms were utilized for both the thermostat and barostat with relaxation times of 0.1 ps and 1.0 ps, respectively. The $x$-, $y$-, and $z$-dimensions were allowed to fluctuate independently while maintaining an orthorhombic geometry. All simulations utilized 100 ps of equilibration followed by an additional 100 ps for generating production data. Pairwise LJ and Coulombic interactions were computed up to 1.5 nm and long-range electrostatic interactions were computed using the particle–particle particle–mesh method [70] with a relative accuracy of $10^{-5}$. No analytical tail corrections were applied to the repulsion-dispersion interactions. All bonds were fully flexible. Simulations were performed with LAMMPS, version 7 Aug 2019 [81].

## III. RESULTS

### A. Case Study: Hydrofluorocarbon Force Fields

Recent international agreements, including the 2016 Kigali Amendment to the 1987 Montreal Protocol, mandated the phaseout of high global warming potential HFC refrigerants [82]. Accurate HFC force fields that are compatible with typical all-atom functional forms are of interest as part of a broader multi-scale engineering effort to sustainably implement this phaseout. Here, we optimize force fields for HFC-32 and HFC-125, the two components of R-410a, a common household refrigerant, to accurately predict the pure-component VLE properties. While an accurate hand-tuned force field for HFC-32 exists in the literature [37], the existing HFC-125 force fields are either inaccurate [15] or rely on less common functional forms [83–85], which often leads to challenges with force field transferability and simulation software compatibility. For HFC-32, we show that our strategy can develop force fields that outperform expert-created models, while for both HFC-32 and HFC-125, we demonstrate the large improvements that are possible compared against "off-the-shelf" models.

We applied a two-stage approach to improve the HFC force fields. Our workflow was first applied to optimize the force fields to accurately predict the LD at the experimental saturation pressure for five temperatures spanning an 80 K temperature range. Following four iterations (LD-1, LD-2, LD-3, and LD-4), 25 parameter sets with low LD MAPE were used to initiate the second stage of force field optimization. In this stage, force field parameters were optimized to accurately predict VLE properties: saturated liquid density, saturated vapor density, vapor pressure, and enthalpy of vaporization. The two-stage approach has advantages: (1) the MD simulations required to compute LD in the isothermal–isobaric en-
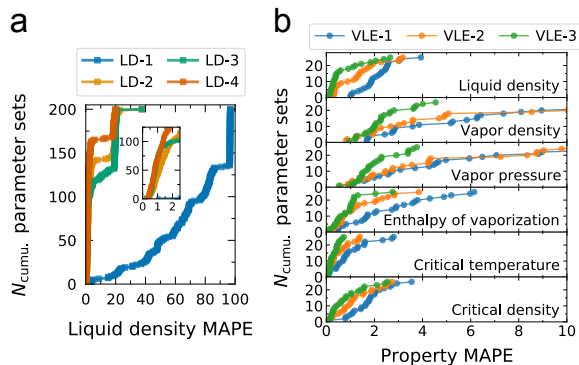
FIG. 2. Cumulative number of HFC-32 parameter sets generated per iteration with less than some MAPE for (a) the liquid density iterations 1–4 (LD-$n$) and (b) vapor–liquid equilibrium iterations 1–3 (VLE-$n$), where $n$ is the iteration number. Inset in panel (a) shows the LD behavior for liquid density MAPE $< 2.5\%$.

semble are computationally less expensive than the MC simulations required to compute VLE properties in the Gibbs ensemble, and (2) the stability of the Gibbs ensemble MC simulations is more sensitive to very poor force field parameters.

Figure 2a shows the cumulative number of parameter sets that yield less than some value of the LD MAPE for each HFC-32 LD iteration. Analogous results for HFC-125 are reported in SI Figure S3. The strength of the surrogate model approach is highlighted by the improvement from the initial liquid density iteration, LD-1, which evaluated 250 parameter sets generated directly from LHS, to the second liquid density iteration, LD-2, which evaluated parameter sets predicted by the surrogate models to yield low LD MAPE. In LD-1 fewer than 5 parameter sets had an LD MAPE below 10%, but LD-2 yielded more than 100 parameter sets with LD MAPE below 2.5%. Limited additional improvements are observed in LD-3 and LD-4, but additional parameter sets with low LD MAPE are nonetheless generated. Figure 2b shows the same information for three VLE workflow iterations (VLE-1, VLE-2, and VLE-3). Consistent improvements in the saturated liquid density, saturated vapor density, vapor pressure, and enthalpy of vaporization are observed from VLE-1 to VLE-3. The results for the critical temperature and critical density also show improvement even though these properties were not explicitly included in the parameter optimization workflow. Note that the saturated liquid density in VLE-1, which evaluated 25 parameter sets generated during the LD stage, performs slightly worse than the results from LD-4 for two reasons: (1) the model vapor pressure is not precisely equal to the experimental vapor pressure, and (2) a smaller system size and shorter interaction cutoff were used to minimize the computational overhead of the LD iterations. Despite the approximation errors introduced by smaller system sizes and cutoffs, the success of our two-

stage optimization strategy shows that initial iterations can be performed with less computationally expensive simulations.
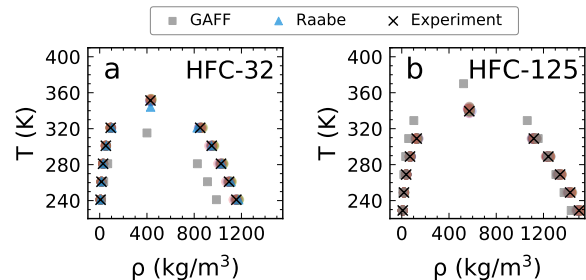


FIG. 3. Vapor–liquid equilibrium envelopes for (a) HFC-32 and (b) HFC-125. The 26 (HFC-32) and 45 (HFC-125) non-dominated parameter sets identified in this work are reported as the transparent colored circles and are compared with literature [15, 37] and experiment [61]. All the non-dominated parameter sets for both HFCs well reproduce the experimental values and are thus highly overlapped.

After completing the three(five) VLE iterations, our force field parameterization workflow yielded 26 HFC-32(45 HFC-125) non-dominated parameter sets. Figure 3 compares vapor–liquid coexistence curves predicted by our non-dominated parameter sets with experiments [61] and force fields for HFC-32 and HFC-125 found in the literature. Results for vapor pressure and enthalpy of vaporization are shown in SI Figure S4. The optimized HFC-32 and HFC-125 force fields are notably better than GAFF, and multiple optimized HFC-32 force fields give improved accuracy in all properties compared to the Raabe force field[37]. We chose an error threshold metric to select a subset of top-performing parameter sets from the non-dominated sets. This yielded four HFC-32 top parameter sets with MAPE of less than 1.5% and four HFC-125 top parameter sets with MAPE of less than 2.5% for the four properties included in the optimization workflow and the critical temperature and critical density. Comparisons of critical temperature and critical density values between experiment, the top four optimized force fields, and literature force fields for both HFCs are shown in SI Tables S7 and S8.

### B. Case Study: Ammonium Perchlorate Force Field

AP is a key ingredient in some solid rocket propellants. Experimental data for physical properties of AP are readily available and a Class II force field parameterized by Zhu et al. [77] has been used to predict [78] pure AP properties at temperatures up to 298 K. The Class II functional form supplements the harmonic diagonal constants found in the more common Class I force fields through the inclusion of cross terms, namely, the stretch–stretch and stretch–bend interactions. The cross
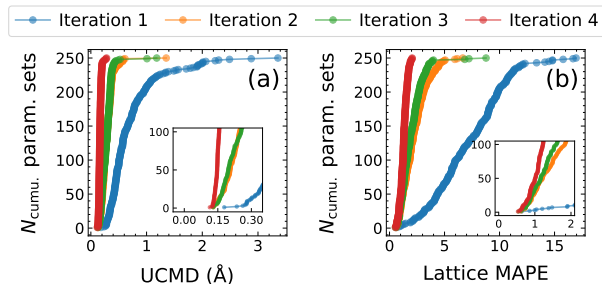
FIG. 4. Cumulative number of AP parameter sets per iteration with less than some value of (a) the 10 K unit cell mean distance (UCMD) and (b) the lattice MAPE. Insets have the same axis titles and focus on the improvement from iteration 3 to iteration 4. Less strict UCMD and lattice MAPE criteria were applied when selecting parameter sets for iterations 2 and 3, and stricter criteria were applied when selecting parameter sets for iteration 4. Threshold values for selecting next iteration points are shown in SI Table S6.

terms couple internal coordinates in an effort to better reproduce the molecular energetics as well as the dynamics of a system by accounting for anharmonic and coupling interactions. However, it is of interest to develop a Class I force field for AP to use in conjunction with existing Class I force fields for the other components of conventional solid propellant, aluminum oxide [86] and the polymeric binder [87]. Here, we parameterize an AP force field with our force field optimization workflow; we previously had utilized hand-tuning methods to develop a Class I AP force field. We present a comparison between the conventional hand-tuning approach and our workflow. In addition to the motivation provided above, we selected solid AP as our second case study because it represents a very different system than the HFC VLE investigated in the first case study.

The properties to which we calibrated our Class I force field were: (1) UCMD at 10 K, defined as the mean of the absolute residuals of equilibrium average simulated atomic positions in reference to the experimentally observed unit cell atomic positions (low values indicate the simulation maintains the experimental AP crystal structure); (2) unit cell lattice parameter mean absolute percent error at the three temperatures of interest (10, 78, and 298 K); and (3) correct hydrogen bond symmetry.

Four iterations of the force field optimization workflow were performed. The cumulative error plots are shown in Figure 4. Once again, we observe substantial improvement between the first and second workflow iteration. Here, the cumulative error plots also show that the criteria for selecting parameter sets for the next iteration can significantly affect the improvement in objective performance between iterations. Less strict UCMD and lattice MAPE criteria were applied when selecting parameter sets for iterations 2 and 3, and stricter criteria were applied when selecting parameter sets for iteration 4; iteration 4 showed much greater improvement over itera-
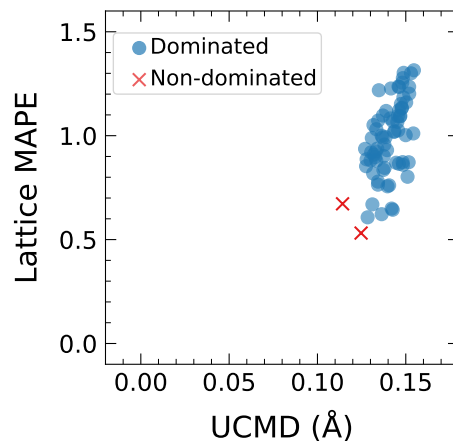


FIG. 5. 70 AP parameter sets that yield lower UCMD and lattice parameter errors than the hand-tuned values while maintaining the correct hydrogen bonding symmetry. The red points are non-dominated and indicate our top two AP parameter sets. The blue points are dominated.

TABLE I. The crystal structure results for the top two AP parameter sets, "Top A and Top B", identified via the workflow presented in this study, the hand-tuned parameter set (HT), and the Class II parameter set of Zhu et al [77]. Lattice parameter results are reported in terms of percent error relative to experimental results [76]. The UCMD results are given in Å.

| Property | T (K) | Top A | Top B | HT | Class II |
|---|---|---|---|---|---|
| | 298 | -0.77 | -0.40 | -2.09 | -0.21 |
| Lat. $a$ | 78 | -0.88 | -0.48 | -1.87 | -2.79 |
| | 10 | -0.24 | 0.26 | -1.38 | -3.10 |
| | 298 | 1.13 | 0.61 | 1.96 | 7.00 |
| Lat. $b$ | 78 | 1.11 | 0.89 | 1.68 | 8.19 |
| | 10 | 0.63 | 0.26 | 1.16 | 8.22 |
| | 298 | -0.18 | -0.74 | -1.04 | 1.64 |
| Lat. $c$ | 78 | -0.71 | -1.10 | -1.31 | 0.46 |
| | 10 | 0.39 | 0.04 | -0.30 | 0.32 |
| MAPE | — | 0.67 | 0.53 | 1.42 | 3.55 |
| UCMD | 10 | 0.1142 | 0.1247 | 0.1560 | 0.3485 |

tion 3 whereas iterations 2 and 3 are very similar. Our workflow generated 70 parameter sets over the four iterations which gave lower UCMD and lattice parameter errors than the hand-tuned values while maintaining the correct hydrogen bonding symmetry. We found two non-dominated parameter sets, as shown in Figure 5. These two non-dominated parameter sets will subsequently be referred to as our top two AP parameter sets. Table I compares the AP results for these top parameter sets with the hand-tuned and Class II force field results.

## IV. DISCUSSION

### A. Many Distinct Parameter Sets Yield Equally Accurate Results

The conventional wisdom in molecular modeling often seems to be that there is a single "correct" or "best" set of force field parameters, but this may be a misleading way to think about force field optimization. No force field is a perfect representation of the physical world. Therefore, model limitations will result in trade-offs between different objectives, and, depending on the property priorities for a specific application, lead to different optimal parameter sets [88]. However, our results clearly show that multiple parameter sets can reproduce several experimental properties with very low error. For the HFCs, our procedure yielded 26 (HFC-32) and 45 (HFC-125) non-dominated parameter sets, which are distinctly different parameterizations, all of which display good performance on our optimization objectives and the critical temperature and density. A visual representation of the non-dominated parameter sets and their performance for the optimization objectives is shown in Figure 6. For HFC-32, where there are 6 optimized force field parameters, the non-dominated parameter sets show variation of up to ~0.3 Å in the carbon and fluorine $\sigma$ values and up to ~10 K/$k_B$ in the carbon and fluorine $\varepsilon$ values. For HFC-125, there is even larger variation in the $\sigma$ and $\varepsilon$ values among the non-dominated parameter sets. We suspect this is because there are a larger number of parameters for HFC-125 (10) than for HFC-32 (6), allowing for compensating behavior between different parameters. For example, consider $\sigma_{F1}$ and $\sigma_{F2}$. There is a clear compensating effect: when $\sigma_{F1}$ is larger, $\sigma_{F2}$ is smaller, and vice-versa. On the other hand, $\sigma_{F1}$ and $\sigma_{F2}$ do appear to be different, as some parameterizations of $\sigma_{F1}$ are 0.3 Å larger than any of the parameterizations of $\sigma_{F2}$.

The visualizations in Figure 6 suggest that the 26 (HFC-32) and 45 (HFC-125) non-dominated parameter sets are indeed distinct parameterizations, rather than closely related parameterizations with small variations along a continuous manifold of good parameters. To further investigate this question, the $L_1$ distance between the best-performing parameter set in each property and every other non-dominated parameter set was calculated and plotted against the property error (SI Figure S5). No correlation is observed between the similarity of a parameter set to the top-performing parameter set in a given property and the property error for that parameter set. This strongly suggests that our non-dominated parameter sets are indeed distinct parameterizations. In part, this can be attributed to our procedure for advancing parameter sets to the next iteration, where we intentionally selected points that were well-separated in parameter space (Section II B 4).

Similar behavior is observed in the AP system, where we identified 70 parameter sets that outperform the hand-tuned Class I and existing Class II force fields[77].
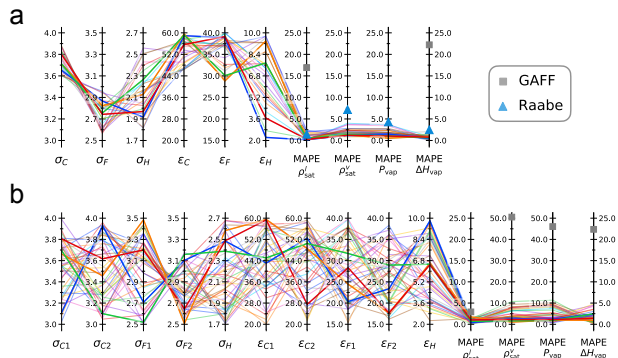


FIG. 6. Repulsion-dispersion parameters for (a) 26 HFC-32 and (b) 45 HFC-125 high quality parameter sets. $\sigma$ is reported in units of Å and $\varepsilon$ is reported in units of K/$k_B$. Each parameter set is connected by a different color line. Thick lines indicate the top 4 parameter sets for each molecule. The y-axes are scaled to show the full range investigated for each parameter. The final four y-axes show the performance for the training objectives. The gray squares and cyan triangles show the performance of GAFF [15] and the force field of Raabe [37], respectively. For HFC-32 the GAFF MAPE for $\rho_{vap}$ and $P_{vap}$ are not shown as they are 133 and 104, respectively.
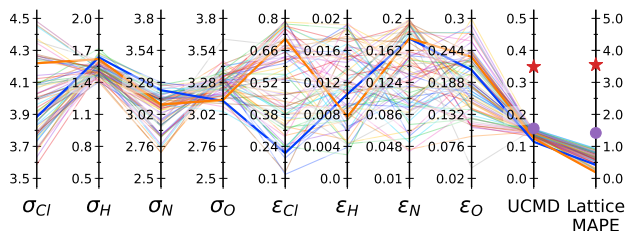


FIG. 7. Repulsion-dispersion parameters for the final 70 AP parameter sets. $\sigma$ is reported in units of Å and $\varepsilon$ is reported in units of kcal/mol. Each parameter set is connected by a different color line. The thick lines show the top 2 AP parameter sets. The y-axes are scaled to show the full range investigated for each parameter. The final two y-axes show the training objectives. The red stars and purple circles show the performance of the Class II force field of Zhu et al. [77] and the hand-tuned Class I force field, respectively.

Figure 7 shows the variation in the optimized AP force field parameters. Once again, a number of distinct parameterizations yield similar accuracy for the optimization objectives. The $\sigma$ values vary by ~0.3 Å for the hydrogen and oxygen atom types that are more exposed to intermolecular interactions, and up to as much as nearly 1.0 Å for the buried Cl atom type. The $\varepsilon$ values vary by as much as ~0.6 kcal/mol, with the largest variation once again observed for the Cl atom type. Although there is a large variation in the individual parameter values between different parameter sets, it is the entire parameter set, taken together, that provides good performance. The results presented here do not suggest that a parameter

can take any value within the ranges shown in Figure 7, e.g., any value of $\sigma_{Cl}$ between 3.5 and 4.5 Å, and yield good performance if all other other parameter values are held constant. Rather, correlations between the different parameters enable a number of distinct yet highly accurate force field parameterizations.

Finding many distinct well-performing non-dominated parameter sets suggests the model may be overparameterized. To investigate this, we performed a local identifiability analysis by inspecting the eigenvalues of the Fisher information matrix (FIM) for the top four parameter sets for both the HFC-32 and HFC-125 models. As detailed in the SI Discussion, we find the FIM has one and five near-zero eigenvalues for HFC-32 and HFC-125, respectively, when considering only the liquid density data. This means we can only identify five (HFC-32: 6 total parameters minus 1 near-zero eigenvalue equals 5 identifiable directions, HFC-125: 10 minus 5 equals 5) parameters using only experimental liquid density data. The corresponding eigenvectors for these near zero eigenvalues reveal the direction in parameter space in which the regression objective is flat (near zero curvature). Unfortunately, these eigenvectors do not point in the direction of a single parameter, which complicates their interpretation. More importantly, the FIM is full rank when simultaneously regressing both liquid density and VLE experimental datasets, which implies both models are locally fully identifiable. Thus, this analysis resolves one aspect of overparameterization by mathematically quantifying the importance of including multiple types of experimental data in the model calibration process. Moreover, our results suggest all of the top parameter sets are near locally optimal solutions (all with positive curvature, thus locally identifiable).

Another aspect of overparameterization is that we find a large number of high-quality solutions. These results are not surprising, given that many inverse problems based on engineering models have numerous locally optimal parameter sets that lead to accurate in-sample predictions.[89] In this case, we hypothesize that parameterizing each molecule individually leads to many locally optimal parameter sets. Extending our method to simultaneously optimize force field parameters for an entire class of molecules (e.g., all hydrofluorocarbons) with a number of shared atom types will likely reduce the overparameterization. While we leave the development of an HFCs force field for future work, here, we explore the effects of using shared atom types for HFC-32 and HFC-125 on the number of high-quality model parameterizations. We consider four atom-typing schemes (AT-1, AT-2, AT-3, and AT-4), shown in Figure 8b. AT-1 is the scheme we have used thus far; there are eight total atom types, three for HFC-32 and five for HFC-125. In AT-2, we use a total of three atom types across both molecules, C, F, and H. AT-3 and AT-4 both use five atom types, but differ in how these atom types are distributed. In AT-3, we maintain the original scheme for HFC-125, but then re-use the C1, F1, and H1 types for

HFC-32. In AT-4, the C and H types are shared as they are either small or buried, while each fluorine is a different atom type. The surrogate models trained during this work were used to evaluate the performance of the different atom typing schemes. LHS was used to generate 500,000 parameter sets. First, the liquid density GP surrogate model was used to eliminate any parameter sets with RMSE greater than 100 kg/m$^3$. For each of the remaining parameter sets, the VLE GP surrogate models were used to predict the MAPE for each VLE property (saturated liquid and vapor densities, vapor pressure, and enthalpy of vaporization). Figure 8a reports the percentage of the original 500,000 parameter sets that yield less than a given MAPE threshold for all four VLE properties, simultaneously. The atom-typing schemes with a reduced number of atom types have a much smaller percentage of parameter space containing low-error parameter sets. In fact, AT-2, with only 3 atom types, does not result in any parameterizations that are predicted to have below 46% MAPE for all four VLE properties. AT-3 and AT-4 show that even with the same number of atom types, one atom-typing scheme may result in superior performance. This naturally raises another question: given different atom-typing schemes, which should be used? Recent work [90] demonstrates the promise of using Bayes factors to compare models with different levels of complexity (e.g., different atom-typing schemes) and make a justified selection.

Since the prior analysis was performed entirely with the predictions of the GP surrogate models, we performed molecular simulations with two top-performing parameter sets for each of the shared atom-typing schemes (AT-2, AT-3, and AT-4) in order to compute the simulated MAPE values and compare them with the surrogate model predictions. The results are reported in SI Table S9. Overall, the surrogate model predictions were excellent, often showing less than 0.5% MAPE deviation from the simulated MAPE. GEMC simulations for AT-2 were unstable at the highest temperature, confirming the surrogate models' prediction that AT-2 would not yield any good parameter sets. We also explored HFC-125-only force fields with a reduced number of atom types (SI Table S10), and found that we were able to identify parameter sets with less than 3% MAPE using only 3 atom types (C, F, and H). However, as noted above, when we attempted to use three atom types (C, F, and H) for both HFC-32 and HFC-125, no good force fields were identified. This finding is strong evidence that the fluorine atom types in HFC-32 and HFC-125 should be different (e.g., AT-4), and shows how developing parameterizations for an entire class of molecules will reduce the number of viable parameter sets.

Adding additional objective properties is a complementary strategy to reduce the number of viable parameter sets. In that case, it is important that the additional properties are orthogonal in the sense that good performance for one property is not highly correlated with good performance for another property. If prop-
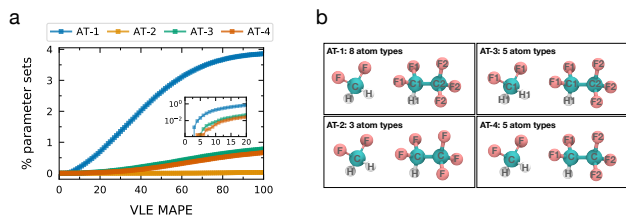
FIG. 8. (a) Cumulative percent of parameter sets from a large ($\mathcal{O}(10^5)$) Latin hypercube that yield less than each value of MAPE for all four VLE properties. For a given MAPE, a higher percentage indicates that more parameterizations achieve at least that threshold level of accuracy. Results are shown for four different atom-typing schemes (AT-1, AT-2, AT-3, and AT-4). The inset focuses on the low-MAPE region and reports the data on a log-scale. (b) Schematic of AT-1, AT-2, AT-3, and AT-4. AT-1 is the original atom-typing scheme where no atom types were shared between HFC-32 and HFC-125.

erty performance is highly correlated, then adding additional properties to the optimization workflow may not substantially reduce the number of viable parameter sets. The apparent overparameterization observed in this work emphasizes why tuning force fields for specific systems and using a few objective properties via relatively simple methods such as epsilon-scaling, manipulating mixing rules, or varying a single parameter value are often quite successful. However, our findings suggest that the force fields developed via these methods are most likely only one of a large number of possible parameterizations that would yield at least equal accuracy.

A further question involves how final parameter sets should ultimately be selected, given that many high-quality parameter sets are available. Our workflow is explicitly not designed to identify a single optimal set of force field parameters. Instead, it searches for and identifies high quality parameter sets with respect to all of the optimization objectives, e.g., points in the non-dominated set. Selecting a single specific parameter set from the optimized parameter sets identified by the workflow requires additional *post hoc* criteria that are application specific. Here, we chose non-dominated status and error thresholds for all properties. Alternative strategies include creating a weighted sum of errors in the properties based upon the desired application and domain knowledge, ranking force fields by their error in the various properties studied via statistical tests [54], evaluating the force field's performance for properties not included in the optimization procedure, or selecting parameter sets based upon a measure of compatibility with the force fields being used for other components of a system. One could also consider chemical intuition when selecting the final parameter sets, e.g., for HFC-125, perhaps a parameter set with more similar values for both fluorine atoms would be preferred. Though our preference is to minimize the number of *ad hoc* choices, ultimately, selecting the final force field for a given application will be system

and application dependent and rely heavily on domain expertise.

## B. Maintaining a Physically-Motivated Analytical Functional Form Aids Transferability to Properties Not Included as Optimization Objectives

One important question is whether the force field parameters developed with this workflow will yield accurate property predictions for properties not included in the optimization workflow. We have already shown that the HFC force fields developed during the VLE tuning stage result in accurate critical temperature and density even though these properties were not optimization objectives. However, these critical properties are largely determined by accurately capturing the temperature dependence of the saturated liquid and saturated vapor density, both of which were optimization objectives. To further investigate the transferability of force field parameters developed with our workflow to properties not included as optimization objectives, we examine the performance of the 25 parameter sets used during the VLE-1 iteration. These parameter sets were used for VLE-1 because they were identified as good at predicting the temperature dependence of the liquid density during the LD iterations. Figure 2 shows that when applied for VLE-1, many perform quite well for VLE properties. In fact, three of the HFC-32 parameter sets used for the VLE-1 iteration had less than 2% MAPE in all six properties. Furthermore, when compared with GAFF, all 25 parameter sets selected from the LD stage yield better performance for all six properties. This is strong evidence that our force field optimization workflow can, with the correct optimization objectives, yield force fields that accurately predict properties beyond the optimization objectives.

The transferability of the LD-optimized parameters to VLE gives credence to our overall force field optimization philosophy, which maintains traditional analytical functional forms and uses machine learning as a guide to identify optimal parameters. However, *a priori*, it is unclear that there should be such a strong correlation between the liquid density and VLE properties. For many systems, accurately predicting the liquid density is a necessary, but often quite insufficient, condition for an accurate force field. We hypothesize there is a key factor that contributes to the transferability of the parameters developed during the LD iterations to VLE: the LD simulations were performed at the saturated vapor pressure across an 80 K temperature range, up to within 30 K of the experimental critical temperature. Accurately capturing the liquid density at saturation across a relatively large temperature range and avoiding spontaneous vaporization, especially at conditions closer to the critical point, requires capturing a careful balance of the cohesive energy and molecular size, which are closely related to the LJ repulsion-dispersion parameters that were calibrated. If the correlation between LD-optimized pa-

rameters and VLE properties proves applicable to other classes of molecules, it may offer a rapid method for developing force fields with accurate VLE properties.

### C. Selecting Good Properties for Force Field Optimization is Challenging

When optimizing force fields for the HFC case study, we were interested in developing force fields that accurately predict HFC VLE behavior. As such, we chose to calibrate parameters to the saturated liquid and vapor densities, vapor pressure, and enthalpy of vaporization. However, these properties are expensive to compute in molecular simulations, making it difficult to evaluate a large parameter space. Therefore we used less computationally expensive LD iterations to generate good parameter sets for VLE and narrow the parameter search space. Furthermore, we continued to use the highly accurate LD GP surrogate models to screen out poor parameter sets during the VLE iterations. The success of this approach demonstrates that a cheaper "screening" property can be used to narrow the parameter search space drastically when good parameter sets for the screening property are a superset of the good parameter sets for the final properties of interest.

The AP case study had different challenges. The MD simulations required to predict the AP properties were computationally inexpensive, so there was no need to first use a screening property. However, it was not immediately clear what experimental properties we should target. Our first implementation attempted to reproduce the temperature dependence of the crystal lattice parameters alone; this proved ineffective, and naive in hindsight, as we generated many force fields that yielded the correct crystal lattice parameters but incorrect crystal structures. To overcome this issue, we added the 10 K UCMD as an objective because it is a measure of how accurately the force field reproduces the experimental crystal structure at 10 K. The lattice MAPE was still included to capture the temperature dependence of the crystal dimensions since the experimental unit cell coordinates are only reported at 10 K.

The UCMD surrogate model has a notable difference from the others; whereas the other surrogate models predict a property (e.g., lattice $a$ or $p_{vap}$), the UCMD is itself an objective function. The UCMD surrogate model predicts the mean distance of all of the unit cell atoms from their respective coordinates in the experimental unit cell. By definition, this distance is zero if the simulated structure perfectly matches experiment. There are benefits to using physical experimentally measured properties compared to an objective function within the optimization workflow, including providing a clear mapping between a surrogate model and the objective metric. However, using surrogate models to predict the value of an objective function provides the opportunity to combine multiple pieces of information into a single quantity, as is the case
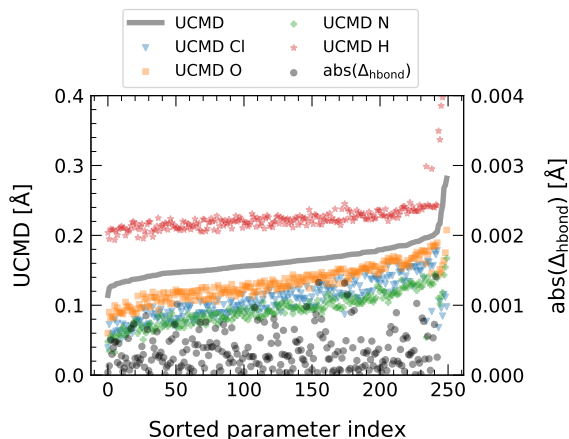


FIG. 9. Overall (gray line) unit cell mean distance (UCMD) compared with UCMD of the four different atom types (points) for the parameter sets tested during iteration 4 of the ammonium perchlorate force field optimization. The hydrogen bond symmetry is reported as abs($\Delta_{hbond}$), where $\Delta_{hbond}$ is the difference in the symmetric hydrogen bond lengths.

with UCMD, which combines the distance of 40 atoms from their positions in the experimental unit cell into a single value. This strategy can drastically reduce the number of required surrogate models. In general, our experience with the AP case study emphasizes that careful thought must be given as to which experimental properties are best to target and how these should be accounted for within the workflow. Roughly 75% of our effort for the AP case study was dedicated to identifying the appropriate experimental properties to target.

### D. Systematic Parameter Search Provides Insights into Model Limitations

The exhaustive search of parameter space enabled by our workflow provides opportunities to distinguish between inaccurate results from poor parameter sets and physical limits from our choice in force field functional form and unoptimized parameters. For example, although our workflow finds high-quality AP parameter sets, we encountered limitations that likely arise from parameters that were not calibrated, and possibly even the force field functional form that we selected. No parameter set predicted an overall UCMD of less than 0.1 Å. Given the exhaustive search enabled by our force field optimization workflow, this suggests that there are no parameter sets capable of yielding a crystal structure with UCMD below 0.1 Å, given the selected functional form, intramolecular parameters, and partial charges. Figure 9 shows the per-element UCMD distances after iteration 4. Although the UCMD for the chlorine, oxygen, and nitrogen atoms fall between 0.1 Å and 0.15 Å for many parameter sets, the hydrogen UCMD rarely falls below 0.2 Å.

Further investigation suggests that this effect is because the N—H bond stretching is insufficiently susceptible to the three unique local hydrogen-bonding chemical environments; experiments report[76] that the N—H bond lengths range between 1.028–1.058 Å whereas in simulations the N—H bond lengths typically cover a much smaller range — between 1.025–1.033 Å — for parameter sets that well reproduce the experimental physical properties. The N—H stretching force constant was not included in our parameterization process. However, even if it was, it is not clear that it would be possible to capture the correct bond stretching behavior and match the vibrational spectra and the N—H bond lengths with a Class I functional form. The exhaustive search provides confidence that the limitations of the model arise from the functional form and unoptimized parameters, rather than the selected parameterization.

## V. CONCLUSIONS

We have presented a machine learning directed workflow for top-down optimization of force field parameters. By harnessing surrogate-assisted optimization, our workflow drastically reduces the number of simulations necessary to find optimal force field parameters by replacing them with computationally tractable surrogate model evaluations. We synthesize GPR and SVM surrogate models and multiobjective optimization into a generic approach to optimize all-atom force fields for realistic systems. We have applied our workflow to optimize HFC force fields to VLE properties and an AP force field to the experimental crystal structure. These case studies show that our workflow can be used for systematic exhaustive screening of parameter space and that surrogate models are highly effective at predicting both simulated physical properties and objective metrics, enabling us to find multiple low-error force fields. The approach presented here could be further combined with gradient-based methods or other approaches such as trust region surrogate-based optimization [91] to further refine the final force fields.

Based upon the success of our approach for the two disparate case studies presented here, we believe that this workflow can be applied to most molecular systems and optimization objectives, provided sufficient reference data. Surrogate models could be used to predict difficult-to-compute thermodynamic properties such as solubilities and binding energies, and transport properties such as self-diffusivity and thermal conductivity. While we have focused on calibrating repulsion-dispersion parameters in this work, this workflow could be used to calibrate any parameters within the force field in a fully top-down approach or as part of a bottom-up force field development workflow, by including *ab initio* data in the fitting procedure [45]. Additionally, we discussed the reasons for successes and limitations of the workflow, the potential challenges of applying this workflow to a particular system (i.e. choosing optimization objectives), and the questions about molecular modeling these results present. We highlight that this workflow is built on a foundation of domain knowledge in selecting the parameters to calibrate, the parameter bounds, and the experimental properties to ensure results are reasonable.

Finally, while we believe that our workflow will enable more efficient force field development and optimization in the future, reducing the need for laborious hand-tuning practices, quantifying the workflow's efficiency was beyond the scope of this work. We can, however, anecdotally note for the AP case study that the hand-tuning approach utilized ∼15,000 simulations and only found 1 optimal parameter set. This is in contrast to our presented workflow, which evaluated ∼3,000,000 parameter sets using surrogate models, $O(10^3)$ times as many as the hand-tuning method, but only required 3,000 simulations, to find 70 parameter sets with lower error in the metrics of interest than the hand-tuned parameter set. We anticipate further refining the proposed workflow, e.g., incorporating adaptive sampling via Bayesian optimization, can dramatically reduce the number of molecular simulations required to identify parameter sets that accurately predict several physical properties.

TABLE II. Description of abbreviations utilized.

| Abbreviation | Expansion |
| --- | --- |
| AP | Ammonium perchlorate |
| APE | Absolute percent error |
| AT | Atom-typing scheme |
| FIM | Fischer information matrix |
| GAFF | General AMBER Force Field |
| GEMC | Gibbs ensemble Monte Carlo |
| GP | Gaussian process |
| GPR | Gaussian process regression |
| HFC | Hydrofluorocarbon |
| HFC-125 | Pentafluoroethane |
| HFC-32 | Difluoromethane |
| LD | Liquid density |
| LHS | Latin hypercube sampling |
| LJ | Lennard-Jones |
| MAPE | Mean absolute percent error |
| MC | Monte Carlo |
| MD | Molecular dynamics |
| RMSE | Root mean square error |
| SVM | Support vector machine |
| UCMD | Unit cell mean distance |
| VLE | Vapor–liquid equilibrium |

**DATA AND SOFTWARE AVAILABILITY**

Codes used to perform the HFC case study and all generated parameters sets are available at: https://github.com/dowlinglab/hfcs-ffit. Codes used to perform the AP case study and all generated parameter sets are available at: https://github.com/dowlinglab/ap-ffit.

[1] E. J. Maginn, AIChE J. **55**, 1304 (2009).

[2] R. Iftimie, P. Minary, and M. E. Tuckerman, Proc. Natl. Acad. Sci. U.S.A. **102**, 6654 (2005).

[3] L. Heo and M. Feig, Proc. Natl. Acad. Sci. U.S.A. **115**, 13276 (2018).

[4] S. A. Hollingsworth and R. O. Dror, Neuron **99**, 1129 (2018).

[5] A. A. Franco, A. Rucci, D. Brandell, C. Frayret, M. Gaberscek, P. Jankowski, and P. Johansson, Chem. Rev. **119**, 4569 (2019).

[6] T. G. A. Youngs, M. G. Del Pópolo, and J. Kohanoff, J. Phys. Chem. B **110**, 5697 (2006).

[7] O. Lobanova, A. Mejía, G. Jackson, and E. A. Müller, J. Chem. Thermodyn. **93**, 320 (2016).

[8] F. Noé, A. Tkatchenko, K.-R. Müller, and C. Clementi, Ann. Rev. Phys. Chem. **71**, 361 (2020).

[9] O. T. Unke, S. Chmiela, H. E. Sauceda, M. Gastegger, I. Poltavsky, K. T. Schütt, A. Tkatchenko, and K.-R. Müller, Chem. Rev. (2021), 10.1021/acs.chemrev.0c01111.

[10] J. Behler and M. Parrinello, Phys. Rev. Lett. **98**, 146401 (2007).

[11] K. T. Schütt, F. Arbabzadah, S. Chmiela, K. R. Müller, and A. Tkatchenko, Nature Commun. **8**, 1 (2017).

[12] J. A. Harrison, J. D. Schall, S. Maskey, P. T. Mikulski, M. T. Knippenberg, and B. H. Morrow, Appl. Phys. Rev. **5**, 031104 (2018).

[13] A. K. Rappe, C. J. Casewit, K. S. Colwell, W. A. Goddard, and W. M. Skiff, J. Am. Chem. Soc. **114**, 10024 (1992).

[14] W. L. Jorgensen, D. S. Maxwell, and J. Tirado-Rives, J. Am. Chem. Soc. **118**, 11225 (1996).

[15] J. Wang, R. M. Wolf, J. W. Caldwell, P. A. Kollman, and D. A. Case, J. Comput. Chem. **25**, 1157 (2004).

[16] K. Vanommeslaeghe, E. Hatcher, C. Acharya, S. Kundu, S. Zhong, J. Shim, E. Darian, O. Guvench, P. Lopes, I. Vorobyov, and A. D. Mackerell Jr., J. Comput. Chem. **31**, 671 (2010).

[17] M. G. Martin, Fluid Phase Equilibr. **248**, 50 (2006).

[18] C. I. Bayly, P. Cieplak, W. Cornell, and P. A. Kollman, J. Phys. Chem. **97**, 10269 (1993).

[19] J. Wang and P. A. Kollman, J. Comput. Chem. **22**, 1219 (2001).

[20] V. Barone, I. Cacelli, N. D. Mitri, D. Licari, S. Monti, and G. Prampolini, Phys. Chem. Chem. Phys. **15**, 3736 (2013).

[21] R. M. Betz and R. C. Walker, J. Comput. Chem. **36**, 79 (2015).

[22] F. Zahariev, N. D. Silva, M. S. Gordon, T. L. Windus, and M. Dick-Perez, J. Chem. Inf. Model **57**, 391 (2017).

[23] L.-P. Wang, K. A. McKiernan, J. Gomes, K. A. Beauchamp, T. Head-Gordon, J. E. Rice, W. C. Swope, T. J. Martínez, and V. S. Pande, J. Phys. Chem. B **121**, 4023 (2017).

[24] M. V. Ivanov, M. R. Talipov, and Q. K. Timerghazin, J. Phys. Chem. A **119**, 1422 (2015).

[25] R. Wang, M. Ozhgibesov, and H. Hirao, J. Comput. Chem. **39**, 307 (2018).

[26] S. Grimme, WIREs Comp. Mol. Sci. **1**, 211 (2011).

[27] S. Mostaghim, M. Hoffmann, P. H. König, T. Frauenheim, and J. Teich, Proceedings of the 2004 Congress on Evolutionary Computation (IEEE Cat. No.04TH8753) **1**, 212 (2004).

[28] A. Jaramillo-Botero, S. Naserifar, and W. A. Goddard, J. Chem. Theory Comput. **10**, 1426 (2014).

[29] K. Stöbener, P. Klein, M. Horsch, K. Küfer, and H. Hasse, Fluid Phase Equilibr. **411**, 33 (2016).

[30] A. Krishnamoorthy, A. Mishra, D. Kamal, S. Hong, K.-i. Nomura, S. Tiwari, A. Nakano, R. Kalia, R. Ramprasad, and P. Vashishta, SoftwareX **13**, 100663 (2021).

[31] A. Krishnamoorthy, A. Mishra, N. Grabar, N. Baradwaj, R. K. Kalia, A. Nakano, and P. Vashishta, Comput. Phys. Commun. **254**, 107337 (2020).

[32] K. Miettinen, *Nonlinear Multiobjective Optimization* (Springer Science & Business Media, New York, NY, 1998).

[33] A. W. Dowling, G. Ruiz-Mercado, and V. M. Zavala, Comput. Chem. Eng. **90**, 136 (2016).

[34] L.-P. Wang, T. J. Martinez, and V. S. Pande, J. Phys. Chem. Lett. **5**, 1885 (2014).

[35] M. Hülsmann and D. Reith, Entropy **15**, 3640 (2013).

[36] K. Murzyn, M. Bratek, and M. Pasenkiewicz-Gierula, J. Phys. Chem. B **117**, 16388 (2013).

[37] G. Raabe, J. Chem. Eng. Data **58**, 1867 (2013).

[38] Y. Zhang, Y. Zhang, M. J. McCready, and E. J. Maginn, J. Chem. Eng. Data **63**, 3488 (2018).

[39] T. Lafitte, A. Apostolakou, C. Avendano, A. Galindo, C. S. Adjiman, E. A. Müller, and G. Jackson, J. Chem. Phys. **139**, 154504 (2013).

[40] E. A. Müller and G. Jackson, Annu. Rev. Chem. Biomol. Eng. **5**, 405 (2014).

[41] F. Rizzi, H. N. Najm, B. J. Debusschere, K. Sargsyan, M. Salloum, H. Adalsteinsson, and O. M. Knio, Multiscale Model. Sim. **10**, 1460 (2012).

[42] M. J. Zimoń, R. Sawko, D. R. Emerson, and C. Thompson, Fluids **2**, 12 (2017).

[43] S. Wu, P. Angelikopoulos, G. Tauriello, C. Papadimitriou, and P. Koumoutsakos, J. Chem. Phys. **145**, 244112 (2016).

[44] R. A. Messerly, S. M. Razavi, and M. R. Shirts, J. Chem. Theory Comput. **14**, 3144 (2018).

[45] H. Liu, Z. Fu, Y. Li, N. F. A. Sabri, and M. Bauchy, MRS Commun. **9**, 593 (2019).

[46] Y. Zhang, *Computational Methods to Assist in Material Discovery: Membranes and Lithium-Ion Battery Electrolytes*, Ph.D. thesis, University of Notre Dame, Notre Dame, IN (2019).

[47] M. Pechlaner, M. M. Reif, and C. Oostenbrink, Mol. Phys. **115**, 1144 (2017).

[48] M. Diem and C. Oostenbrink, J. Chem. Inf. Model. **60**, 279 (2020).

[49] R. Tamura, J. Lin, and T. Miyazaki, J. Phys. Soc. Japan **88**, 044601 (2019).

[50] M. J. Burn and P. L. Popelier, J. Chem. Phys. **153**, 054111 (2020).

[51] J. Vandermause, S. B. Torrisi, S. Batzner, Y. Xie, L. Sun, A. M. Kolpak, and B. Kozinsky, npj Computat. Mater. **6**, 1 (2020).

[52] D. Stephenson, J. R. Kermode, and D. A. Lockerby, Microfluid. Nanofluid. **22**, 1 (2018).

[53] J. L. McDonagh, A. Shkurti, D. J. Bray, R. L. Anderson, and E. O. Pyzer-Knapp, J. Chem. Inf. Model. **59**, 4278 (2019).

[54] J. M. Sestito, M. L. Thatcher, L. Shu, T. A. Harris, and Y. Wang, J. Phys. Chem. A **124**, 5042 (2020).

[55] M. Razi, A. Narayan, R. M. Kirby, and D. Bedrov, Comput. Mater. Sci. **176**, 109518 (2020).

[56] E. Boulanger, L. Huang, C. Rupakheti, A. D. MacKerell, Jr., and B. Roux, J. Chem. Theory Comput. **14**, 3121 (2018).

[57] C. S. Adorf, P. M. Dodd, V. Ramasubramani, and S. C. Glotzer, Comput. Mater. Sci. **146**, 220 (2018).

[58] V. Ramasubramani, C. Adorf, P. Dodd, B. Dice, and S. Glotzer, in *Proceedings of the Python in Science Conference* (2018).

[59] A. G. Matthews, M. van der Wilk, T. Nickson, K. Fujii, A. Boukouvalas, P. León-Villagrá, Z. Ghahramani, and J. Hensman, J. Mach. Learn. Res. **18**, 1 (2017).

[60] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, J. Mach. Learn. Res. **12**, 2825 (2011).

[61] E. Lemmon, I. H. Bell, M. Huber, and M. McLinden, "NIST Standard Reference Database 23: Reference Fluid Thermodynamic and Transport Properties-REFPROP, Version 10.0, National Institute of Standards and Technology," (2018), Standard Reference Data Program, Gaithersburg.

[62] D. Case, T. Darden, I. T.E. Cheatham, C. Simmerling, J. Wang, R. Duke, R. Luo, R. Walker, W. Zhang, K. Merz, B. Roberts, B. Wang, S. Hayik, A. Roitberg, G. Seabra, I. Kolossváry, K. Wong, F. Paesani, J. Vanicek, J. Liu, X. Wu, S. Brozell, T. Steinbrecher, H. Gohlke, Q. Cai, X. Ye, J. Wang, M.-J. Hsieh, G. Cui, D. Roe, D. Mathews, M. Seetin, C. Sagui, V. Babin, T. Luchko, S. Gusarov, A. Kovalenko, and P. Kollman, "AMBER 11," (2010), university of California, San Francisco.

[63] M. J. Frisch, G. W. Trucks, H. B. Schlegel, G. E. Scuseria, M. A. Robb, J. R. Cheeseman, G. Scalmani, V. Barone, G. A. Petersson, H. Nakatsuji, X. Li, M. Caricato, A. Marenich, J. Bloino, B. G. Janesko, R. Gomperts, B. Mennucci, H. P. Hratchian, J. V. Ortiz, A. F. Izmaylov, J. L. Sonnenberg, D. Williams-Young, F. Ding, F. Lipparini, F. Egidi, J. Goings, B. Peng, A. Petrone, T. Henderson, D. Ranasinghe, V. G. Zakrzewski, J. Gao, N. Rega, G. Zheng, W. Liang, M. Hada, M. Ehara, K. Toyota, R. Fukuda, J. Hasegawa, M. Ishida, T. Nakajima, Y. Honda, O. Kitao, H. Nakai, T. Vreven, K. Throssell, J. A. Montgomery Jr., J. E. Peralta, F. Ogliaro, M. Bearpark, J. J. Heyd, E. Brothers, K. N. Kudin, V. N. Staroverov, T. Keith, R. Kobayashi, J. Normand, K. Raghavachari, A. Rendell, J. C. Burant, S. S. Iyengar, J. Tomasi, M. Cossi, J. M. Millam, M. Klene, C. Adamo, R. Cammi, J. W. Ochterski, R. L. Martin, K. Morokuma, O. Farkas, J. B. Foresman, and D. J. Fox, "Gaussian 09, Revision E.01," (2016), gaussian, Inc. Wallingford CT.

[64] A. D. Becke, J. Chem. Phys. **98**, 1372 (1993).

[65] P. J. Stephens, F. J. Devlin, C. F. Chabalowski, and M. J. Frisch, J. Phys. Chem. **98**, 11623 (1994).

[66] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning* (The MIT Press, Cambridge, MA, 2006).

[67] G. Bussi, D. Donadio, and M. Parrinello, J. Chem. Phys. **126**, 014101 (2007).

[68] H. J. Berendsen, J. Postma, W. F. Gunsteren, A. DiNola, and J. R. Haak, J. Chem. Phys. **81**, 3684 (1984).

[69] M. Parrinello and A. Rahman, J. Appl. Phys. **52**, 7182 (1981).

[70] R. W. Hockney and J. W. Eastwood, *Computer Simulation Using Particles* (McGraw-Hill, New York, 1981).

[71] U. Essmann, L. Perera, M. L. Berkowitz, T. Darden, H. Lee, and L. G. Pedersen, J. Chem. Phys. **103**, 8577 (1995).

[72] B. Hess, J. Chem. Theory Comput. **4**, 116 (2008).

[73] M. J. Abraham, T. Murtola, R. Schulz, S. Páll, J. C. Smith, B. Hess, and E. Lindahl, SoftwareX **1**, 19 (2015).

[74] R. S. DeFever, R. A. Matsumoto, A. W. Dowling, P. T. Cummings, and E. J. Maginn, J. Comput. Chem. **42**, 1321 (2021).

[75] J. K. Shah, E. Marin-Rimoldi, R. G. Mullen, B. P. Keene, S. Khan, A. S. Paluch, N. Rai, L. L. Romanielo, T. W. Rosch, B. Yoo, and E. J. Maginn, J. Comput. Chem. **38**, 1727 (2017).

[76] C. S. Choi, H. J. Prask, and E. Prince, J. Chem. Phys. **61**, 3523 (1974).

[77] W. Zhu, X. Wang, J. Xiao, W. Zhu, H. Sun, and H. Xiao, J. Hazard. Mater. **167**, 810 (2009).

[78] G. M. Tow and E. J. Maginn, J. Chem. Phys. **149**, 244502 (2018).

[79] M. E. Tuckerman, J. Alejandre, R. López-Rendón, A. L. Jochim, and G. J. Martyna, J. Phys. A: Math. Gen. **39**, 5629 (2006).

[80] W. Shinoda, M. Shiga, and M. Mikami, Phys. Rev. B **69**, 134103 (2004).

[81] S. Plimpton, J. Comput. Phys. **117**, 1 (1995).

[82] *Handbook for the Montreal Protocol on Substances that Deplete the Ozone Layer*, Tech. Rep. (2006).

[83] M. Fermeglia, M. Ferrone, and S. Pricl, Fluid Phase Equilibr. **210**, 105 (2003).

[84] M. Lísal and V. Vacek, Fluid Phase Equilibr. **127**, 83 (1997).

[85] J. Stoll, J. Vrabec, and H. Hasse, J. Chem. Phys. **119**, 11396 (2003).

[86] R. T. Cygan, J.-J. Liang, and A. G. Kalinichev, J. Phys. Chem. B **108**, 1255 (2004).

[87] G. M. Tow and E. J. Maginn, Macromolecules **53**, 2594 (2020).

[88] K. Stöbener, P. Klein, S. Reiser, M. Horsch, K.-H. Küfer, and H. Hasse, Fluid Phase Equilibr. **373**, 100 (2014).

[89] L. Tenorio, *An Introduction to Data Analysis and Uncertainty Quantification for Inverse Problems* (SIAM, 2017).

[90] O. C. Madin, S. Boothroyd, R. A. Messerly, J. D. Chodera, J. Fass, and M. R. Shirts, arXiv preprint arXiv:2105.07863 (2021).

[91] A. R. Conn, N. I. Gould, and P. L. Toint, *Trust Region Methods* (SIAM, 2000).

# Supporting Information for: Machine Learning Directed Optimization of Classical Molecular Modeling Force Fields

Bridgette J. Befort,* Ryan S. DeFever,* Garrett M. Tow, Alexander W. Dowling, and Edward J. Maginn†

*Department of Chemical and Biomolecular Engineering,*

*University of Notre Dame, Notre Dame, Indiana 46556, United States*

(Dated: July 15, 2021)

---

* BJ Befort and RS DeFever contributed equally to this work.

† Corresponding author

## S1.   METHODS

### Identifying parameter sets that are well-separated in parameter space

The distance between each parameter set is taken as the $L_1$ norm in scaled parameter space. Scaled parameter space is defined such that the lower bound of a parameter is equal to 0.0 and the upper bound is equal to 1.0. The following algorithm was used to select well-separated points: (1) define a distance threshold, (2) select one parameter set at random and add it to the list of those for the next iteration (3) discard all parameter sets within the distance threshold of the parameters sets selected for the next iteration, (4) return to (2) and continue iterating until no parameter sets remain, (5) check the final number of parameter sets identified for the next iteration, and if more than desired, start over and return to (1) with a larger distance threshold.

## S2.   DISCUSSION

### HFC Identifiability Analysis

Local identifiability analysis was performed via eigenvalue-eigenvector decomposition of the Fischer information matrix (FIM) which describes sensitivity of the fitted parameters to the experimental data. An FIM with eigenvalues of or near zero is singular and indicates that (some) of the parameters are insensitive to the data (see Chapter 10 of Ref. 1 and Section 3.4 of Ref. 2). To build the FIM, first a Jacobian matrix $J^T$ was approximated with the central finite difference formula. Properties resulting from the perturbed parameter sets and used in the gradient calculation were obtained from the GP surrogate models or simulations. The product of the Jacobian and its transpose, $J^T \cdot J$, approximates the FIM. Eigenvalue-eigenvector decomposition was performed on the FIM. The number of non-zero eigenvalues of the FIM indicate the number of directions which are identifiable and the eigenvectors corresponding to near-zero eigenvalues indicate the directions of unidentifiability. Each component in an eigenvector corresponds with a single parameter. If one of these components is of a much larger order of magnitude than the other components in an eigenvector (or the other components were zero), the parameter corresponding to that component would be unidentifiable. However, the components of the eigenvector could all be non-zero and of

similar order of magnitude, indicating that the unidentifiability is in a direction that is the linear combination of all of the parameters.

We applied this analysis to the top four force fields for both HFCs. For the HFC identifiability analyses which used only liquid density, we found there was a single direction of unidentifiability that was a linear combination of parameters for the HFC-32 case and five directions of unidentifiability that were linear combinations of parameters for the HFC-125 case. Upon adding the VLE data into the sensitivity analysis, the models for both HFCs became fully identifiable. We performed the identifiability analysis two different ways: in the first case, we used the GP models to build the Jacobian matrix, and in the second case we performed additional molecular simulations to build the Jacobian matrix. In both cases, we obtained the same conclusions, indicating once again that the GP models are very good at predicting the results from molecular simulations.

Eigenvalue and eigenvector results using GP and simulation predictions for each HFC for liquid density and VLE data are included in the Supporting Information spreadsheets in the zip files 'HFC32-Identifiability.zip' and 'HFC125-Identifiability.zip'.

## S3. FIGURES AND TABLES



FIG. S1. The simulation result compared with the GP surrogate model prediction for the surrogate models trained during the VLE-2 iteration for HFC-32. Comparisons are shown for $\rho_{\mathrm{sat}}^{l}$ (a, b), $\rho_{\mathrm{sat}}^{v}$ (c, d), $P_{\mathrm{vap}}$ (e, f), $\Delta H_{\mathrm{vap}}$ (g, h). Comparisons for the training data are shown in the left column (a, c, e, g) and the comparisons for the test data are shown in the right columns (b, d, f, h). Radial basis function (RBF), Matérn $\nu = 3/2$ (Matern32), and Matérn $\nu = 5/2$ (Matern52) refer to the kernel for the GP surrogate models [3]. These results are representative of the GP model accuracy for LD and VLE iterations of the HFC force field optimization.

FIG. S2. Examples of the GP surrogate model means (lines) and variances (shaded regions) for one parameter set from the VLE-2 iteration for HFC-32. Radial basis function (RBF), Matérn $\nu = 3/2$ (Matern32), and Matérn $\nu = 5/2$ (Matern52) refer to the kernel for the GP surrogate models [3]. Points shown in black were included in the training data for the GP models, whereas points in red were excluded. GP surrogate models shown for $\rho_{\text{sat}}^l$ (a), $\rho_{\text{sat}}^v$ (b), $P_{\text{vap}}$ (c), $\Delta H_{\text{vap}}$ (d). These results are representative of the GP model accuracy for LD and VLE iterations of the HFC force field optimization.

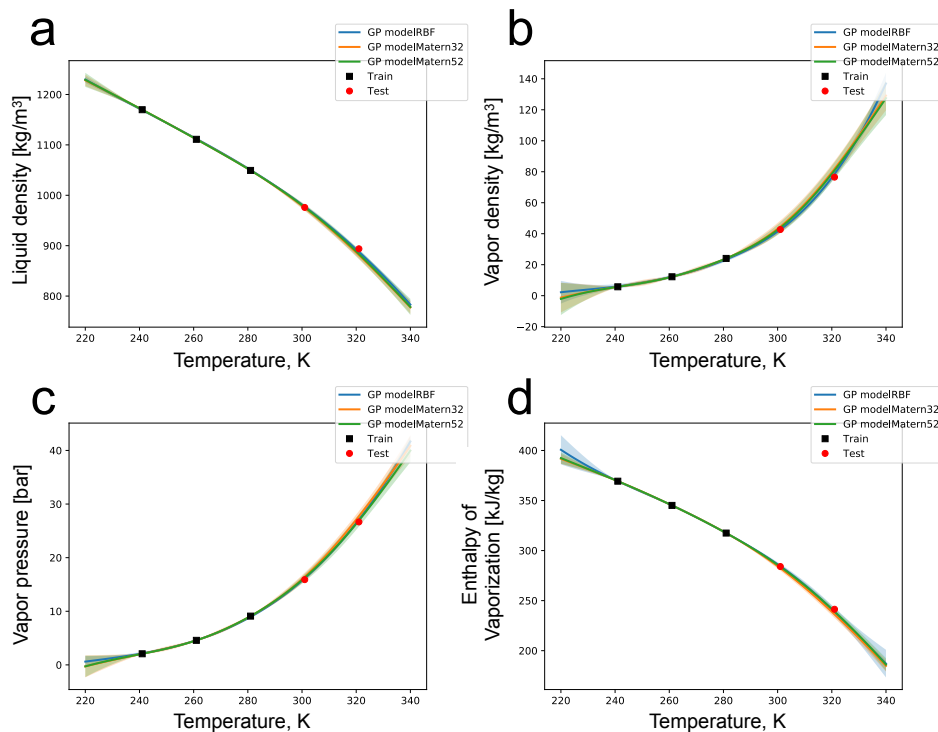FIG. S3. Cumulative number of HFC-125 parameter sets per iteration with less than some MAPE for (a) the liquid density iterations 1–4 (LD-$n$) and (b) vapor–liquid equilibrium iterations 1–5 (VLE-$n$), where $n$ is the iteration number. Inset in panel (a) shows the LD behavior for liquid density MAPE < 2.5%.

FIG. S4. Vapor pressure and enthalpy of vaporization for HFC-32 and HFC-125 force fields compared with literature [4, 5] and experiment [6]. The 26 (HFC-32) and 45 (HFC-125) non-dominated parameter sets are shown as lightly shaded colored circles. All the non-dominated parameter sets for both HFCs well reproduce the experimental values and are thus highly overlapped.

FIG. S5. Distance between the best parameter set for each property ($L_1$ norm with normalized parameter values) and all other parameter sets versus the property error for HFC-32 VLE iterations. The point with an $L_1$ norm of 0.0 shows the performance of the best parameter set for a given property. One point is shown for each parameter set tested during the VLE iterations. The lack of correlation between the $L_1$ distance from the top performing parameter set and the property error emphasizes that high quality parameter sets are distributed throughout parameter space.

TABLE S1. Partial charges and intramolecular parameters for HFC-32 and HFC-125

| | Partial Charges | |
|---|---|---|
| Type | GAFF Type | $q$ ($e$) |
| C | c3 | 0.405467 |
| F | f | -0.250783 |
| H | h2 | 0.0480495 |
| C1 | c3 | 0.224067 |
| C2 | c3 | 0.500886 |
| F1 | f | -0.167131 |
| F2 | f | -0.170758 |
| H1 | h2 | 0.121583 |

| | Bonds | |
|---|---|---|
| GAFF Type | $k_r$ (kcal mol$^{-1}$ Å$^{-2}$) | $r_0$ (Å) |
| c3-f | 356.9 | 1.3497 |
| c3-h2 | 331.7 | 1.0961 |
| c3-c3 | 300.9 | 1.5375 |

| | Angles | |
|---|---|---|
| GAFF Type | $k_\theta$ (kcal mol$^{-1}$ rad$^{-2}$) | $\theta_0$ (deg) |
| f-c3-f | 70.9 | 107.36 |
| f-c3-h2 | 51.1 | 108.79 |
| c3-c3-f | 66.1 | 109.24 |
| c3-c3-h2 | 46.2 | 110.22 |
| h2-c3-h2 | 39.0 | 110.20 |

| | Dihedrals | | |
|---|---|---|---|
| GAFF Type | $\nu_n$ (kcal mol$^{-1}$) | $n$ | $\gamma$ (deg) |
| f-c3-c3-f | 1.20 | 1 | 180.0 |
| f-c3-c3-h2 | 0.1556 | 3 | 0.0 |

TABLE S2. HFC-32 force field tuning parameters

| Intermolecular parameters | | |
|---|---|---|
| Type | $\sigma$ Bounds (Å) | $\varepsilon$ Bounds (K/$k_B$) |
| C | 3.0–4.0 | 20.0–60.0 |
| F | 2.5–3.5 | 15.0–40.0 |
| H | 1.7–2.7 | 2.0–10.0 |

TABLE S3. HFC-125 force field tuning parameters

| Intermolecular parameters | | |
|---|---|---|
| Type | $\sigma$ Bounds (Å) | $\varepsilon$ Bounds (K/$k_B$) |
| C1 | 3.0–4.0 | 20.0–60.0 |
| C2 | 3.0–4.0 | 20.0–60.0 |
| F1 | 2.5–3.5 | 15.0–40.0 |
| F2 | 2.5–3.5 | 15.0–40.0 |
| H | 1.7–2.7 | 2.0–10.0 |

TABLE S4. Hand-tuned AP force field parameters

| Intermolecular parameters | | | |
|---|---|---|---|
| Type | $q$ ($e$) | $\sigma$ (Å) | $\varepsilon$ (kcal/mol) |
| Cl | 1.5456 | 3.9140 | 0.5018 |
| H | 0.387625 | 1.7361 | 0.0027 |
| N | -0.5505 | 3.3078 | 0.0406 |
| O | -0.6364 | 3.3107 | 0.0954 |

| Bonds | | |
|---|---|---|
| Type | $k_r$ (kcal mol$^{-1}$ Å$^{-2}$) | $r_0$ (Å) |
| Cl-O | 426.42 | 1.4523 |
| H-N | 413.55 | 1.0300 |

| Angles | | |
|---|---|---|
| Type | $k_\theta$ (kcal mol$^{-1}$ rad$^{-2}$) | $\theta_0$ (deg) |
| H-N-H | 33.45 | 109.5 |
| O-Cl-O | 107.60 | 109.5 |

TABLE S5. AP force field tuning parameters

| Intermolecular parameters | | |
|---|---|---|
| Type | $\sigma$ Bounds (Å) | $\varepsilon$ Bounds (kcal/mol) |
| Cl | 3.5–4.5 | 0.1–0.8 |
| H | 0.5–2.0 | 0.0–0.02 |
| N | 2.5–3.8 | 0.01–0.2 |
| O | 2.5–3.8 | 0.02–0.3 |

TABLE S6. Screening criteria for AP iterations

| Iteration | Structure Classifier Threshold (Å) | Symmetry Classifier Threshold (Å) | UCMD Threshold (Å) | Lattice MAPE Threshold |
|---|---|---|---|---|
| 1-2 | 0.8 | - | 0.35 | 2.5 |
| 2-3 | 0.8 | - | 0.35 | 2.5 |
| 3-4 | 0.8 | 0.001 | 0.2 | 1.5 |

TABLE S7. Critical temperatures ($T_c$) and densities ($\rho_c$) predicted by GAFF [4], the force field of Raabe [5], and the top four HFC-32 parameter sets compared to experiment [6]

| Force Field | $T_c$ (K) | $\rho_c$ (kg/m$^3$) |
|---|---|---|
| GAFF | 315.3 | 400.1 |
| Raabe | 344.1 | 430.9 |
| Top A | 351.1 | 431.0 |
| Top B | 352.8 | 430.5 |
| Top C | 351.9 | 431.8 |
| Top D | 352.9 | 430.9 |
| Experiment | 351.4 | 429.8 |

TABLE S8. Critical temperatures ($T_c$) and densities ($\rho_c$) predicted by GAFF [4] and top four HFC-125 parameter sets compared to experiment [6]

| Force Field | $T_c$ (K) | $\rho_c$ (kg/m$^3$) |
|---|---|---|
| GAFF | 370.0 | 523.4 |
| Top A | 342.5 | 570.9 |
| Top B | 341.5 | 562.9 |
| Top C | 341.8 | 567.5 |
| Top D | 343.1 | 576.6 |
| Experiment | 339.4 | 571.9 |

TABLE S9. Performance of HFC-32 and HFC-125 force fields with shared atom types. Results reported for the simulated (sim.) and surrogate model (sur.) predictions. The simulated results for HFC-32 with AT-2 are not reported as the highest temperature GEMC simulation was unstable.

| | HFC-32 MAPE | | | | HFC-125 MAPE | | | |
|---|---|---|---|---|---|---|---|---|
| | $\rho_{\text{sat}}^l$ | $\rho_{\text{sat}}^v$ | $P_{\text{vap}}$ | $\Delta H_{\text{vap}}$ | $\rho_{\text{sat}}^l$ | $\rho_{\text{sat}}^v$ | $P_{\text{vap}}$ | $\Delta H_{\text{vap}}$ |
| AT-2 (sim.) | - | - | - | - | 1.5 | 31.6 | 27.6 | 15.4 |
| AT-2 (sur.) | 2.3 | 43.6 | 34.8 | 1.3 | 2.0 | 46.7 | 38.2 | 16.2 |
| AT-3 (sim.) | 0.8 | 2.4 | 1.8 | 1.8 | 2.8 | 4.5 | 2.4 | 4.8 |
| AT-3 (sur.) | 0.8 | 2.3 | 2.2 | 2.0 | 2.7 | 4.0 | 3.0 | 3.8 |
| AT-4 (sim.) | 1.5 | 2.2 | 1.8 | 1.5 | 0.4 | 0.5 | 2.2 | 1.4 |
| AT-4 (sur.) | 1.4 | 2.6 | 1.9 | 1.5 | 0.3 | 3.2 | 3.8 | 1.5 |

TABLE S10. MAPE of HFC-125 force fields with reduced number of atom types. Results reported for the simulated (sim.) and surrogate model (sur.) predictions.

| | HFC-125 MAPE | | | |
|---|---|---|---|---|
| Atom types | $\rho_{\text{sat}}^l$ | $\rho_{\text{sat}}^v$ | $P_{\text{vap}}$ | $\Delta H_{\text{vap}}$ |
| C1, C2, F, H (sim.) | 0.7 | 3.4 | 3.3 | 2.2 |
| C1, C2, F, H (sur.) | 0.5 | 2.5 | 0.5 | 2.4 |
| C1, C2, F, H (sim.) | 0.8 | 4.1 | 4.8 | 1.7 |
| C1, C2, F, H (sur.) | 0.9 | 0.9 | 2.0 | 2.0 |
| C, F, H (sim.) | 0.4 | 2.8 | 1.2 | 1.7 |
| C, F, H (sur.) | 0.5 | 2.5 | 2.4 | 1.1 |
| C, F, H (sim.) | 0.5 | 3.0 | 1.3 | 1.9 |
| C, F, H (sur.) | 0.5 | 1.1 | 1.9 | 2.1 |

[1] Y. Bard, *Nonlinear Parameter Estimation* (Academic Press, Inc., 1974).

[2] G. Seber and C. Wild, *Nonlinear Regression* (John Wiley & Sons, 1989).

[3] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning* (The MIT Press, Cambridge, MA, 2006).

[4] J. Wang, R. M. Wolf, J. W. Caldwell, P. A. Kollman, and D. A. Case, J. Comput. Chem. **25**, 1157 (2004).

[5] G. Raabe, J. Chem. Eng. Data **58**, 1867 (2013).

[6] E. Lemmon, I. H. Bell, M. Huber, and M. McLinden, "NIST Standard Reference Database 23: Reference Fluid Thermodynamic and Transport Properties-REFPROP, Version 10.0, National Institute of Standards and Technology," (2018), standard Reference Data Program, Gaithersburg.