

Topology Automated Force-Field Interactions (TAFFI): A Framework for Developing Transferable Force-Fields

Bumjoon Seo¹, Zih-Yu Lin¹, Qiyuan Zhao¹, Michael A. Webb², and Brett
M. Savoie¹

¹*Davidson School of Chemical Engineering, Purdue University, West
Lafayette, Indiana 47906, United States*

²*Department of Chemical and Biological Engineering, Princeton University,
Princeton, New Jersey 08540, United States*

May 2, 2021

Abstract

Force-field development has undergone a revolution in the past decade with the proliferation of quantum chemistry based parameterizations and the introduction of machine learning approximations of the atomistic potential energy surface. Nevertheless, transferable force-fields with broad coverage of organic chemical space remain necessary for applications in materials and chemical discovery where throughput, consistency, and computational cost are paramount. Here we introduce a force-field development framework called Topology Automated Force-Field Interactions (TAFFI) for developing transferable force-fields of varying complexity against an extensible database

of quantum chemistry calculations. TAFFI formalizes the concept of atom typing and makes it the basis for generating systematic training data that maintains a one-to-one correspondence with force-field terms. This feature makes TAFFI arbitrarily extensible to new chemistries while maintaining internal consistency and transferability. As a demonstration of TAFFI, we have developed a fixed-charge force-field, TAFFI-gen, from scratch that includes coverage for common organic functional groups that is comparable to established transferable force-fields. The performance of TAFFI-gen was benchmarked against OPLS and GAFF for reproducing several experimental properties of 87 organic liquids. The consistent performance of these force-fields, despite their distinct origins, validates the TAFFI framework while also providing evidence of the representability limitations of fixed-charge force-fields.

1 Introduction

Molecular dynamics (MD) simulations are a ubiquitous tool in contemporary materials and chemical characterization. The development of approximations to the atomistic potential energy surface (PES) has been central to extending MD simulations to address large systems, condensed phases, and long timescales.¹⁻⁷ Over the past several decades, many PES approximations (i.e., force-fields) have been implemented, spanning the gamut from relatively simple non-reactive, fixed-charged, and harmonic forms⁸⁻¹⁵ to more recent and complex machine-learning based approximations.¹⁶⁻²⁵ Along this continuum there is an intrinsic trade-off between accuracy and complexity, with fixed-charge force-fields being the most economical description but also exhibiting the most limited representability with respect to approximating the PES. Nevertheless, for specific force-field forms it is still unclear in the extent to which representability limitations versus limited training data cause errors in the properties simulated by MD. This distinction is crucial because representability limitations are fundamental to the form of the force-field,²⁶⁻²⁹ whereas errors associated with training data or parameterization protocols can be redressed without increasing the computational cost or

16 complexity of the force-field.³⁰⁻³⁷ It would thus be desirable to develop a framework capa-
17 ble of parameterizing force-fields of varying complexity against common training data such
18 that representability limitations could be established. In the current work, we demonstrate
19 the implementation of such a framework to benchmark a new fixed-charged force-field from
20 scratch, with the long-term goal of flexibly matching force-field complexity to the required
21 accuracy of an MD simulation.

22 Apart from the specific form of the PES approximation, force-fields are also distinguished
23 by whether they are transferable across chemical species or only applicable to specific sys-
24 tems. The latter strategy is in principle more accurate and easier to implement, as transfer-
25 ability imposes additional requirements on the force-field that may lead to accuracy trade-offs
26 and also necessarily more training data. In a typical system-specific workflow, a user supplies
27 one or more molecules that they want to simulate, a set of quantum chemistry calculations
28 are performed to generate training data, and a one-off approximate force-field is parameter-
29 ized to the training data.³⁸⁻⁴⁰ However, there are many applications, including molecular
30 discovery and reactive systems, where transferable force-fields with general applicability are
31 clearly desirable due to the cost of parameterizing a force-field from scratch every time a new
32 molecule or material is encountered. Nevertheless, the on-the-fly parameterization concept
33 is potentially still applicable to extending transferable force-fields if the associated quan-
34 tum chemistry data is stored and parameterizations of new molecules are performed in a
35 backwards-compatible fashion. This is the approach adopted in the force-field framework
36 developed here.

37 The most mature transferable force-fields are based on the concept of atom types, where
38 the local bonding environment about each atom is used as the basis for transferring force-
39 field terms across recurring bonding motifs. Atom typing reduces the number of parameters
40 required to simulate new molecules, and the concept has precedence in thermodynamic incre-
41 ment theories going back to Pauling. However, even in modern machine learning force-fields,
42 atom types are often latent variables that are learned during training.^{20,23} The challenge

43 for transferable force-fields has always been with extending them to include coverage for
44 new chemistries.⁴¹⁻⁴⁶ Among the specific challenges are generating training data for new
45 chemistries that are consistent with the existing training corpus and performing new param-
46 eterizations with backwards compatibility with the rest of the force-field. For these reasons,
47 the most popular transferable force-fields with the largest chemical coverage are built on
48 top of legacy force-fields with decades of development (GAFF,^{47,48} CGenFF,^{49,50} and OPLS-
49 AA^{15,42,51,52}). Nevertheless, expanding the coverage of these force-fields still typically involves
50 retraining the whole force-field. Although not yet fully realized, machine learning force-fields
51 present a parallel approach to achieving transferability by simply expanding training data
52 to the point that de facto transferability is achieved. Among the ideas presented here, is
53 that these two approaches are not as incompatible as they seem. Specifically, the data gen-
54 eration problem for machine learning force-fields is largely shared with the data generation
55 problem for simpler force-fields, and a framework that systematically expands a corpus of
56 training data on the basis of new atom types would be advantageous regardless of the specific
57 functional form used for the force-field.

58 The current work addresses the challenges of producing arbitrarily extensible and trans-
59 ferable force-fields based on quantum chemistry training data. The presented framework,
60 topology automated force-field interactions⁵³⁻⁵⁵ (TAFFI), accomplishes this by formalizing
61 the concept of atom types using molecular graphs and defining a one-to-one correspondence
62 between force-field parameters and the model compounds used to generate training data.
63 These features are compatible with on-the-fly parameterization of new force-field parameters
64 while maintaining self-consistency and backwards compatibility. The result is an extensible
65 force-field supported by a continuously growing body of training data that can be fit to
66 flexible force-field forms. In the current work, TAFFI is used to derive a fixed-charge force-
67 field (TAFFI-gen) for 87 organic molecules as a case study to illustrate the methodology
68 and benchmark its performance. Additionally, over 2000 distinct force-field terms involving
69 270 unique atom types for TAFFI-gen are distributed with this work, including coverage

70 for many common organic moieties. Condensed-phase simulation results using TAFFI-gen
71 are compared with the GAFF and OPLS-AA force-fields for the reproduction of a range
72 of experimental liquid properties. The consistent performance of these force-fields, despite
73 their distinct origins, validates the TAFFI framework while also providing evidence of the
74 representability limitations of fixed-charge force-fields.

75 **2 Methods**

76 **2.1 Methodology Overview**

77 An overview of the three stages of data generation and force-field parameterization within
78 the TAFFI framework is provided here using diethyl carbonate as an example to guide the
79 reader (Fig. 1). A detailed description of each step is provided in the subsequent sections
80 (2.2-2.4).

81 In Stage 1 (Fig. 1a-c), the atom types and modes associated with the user-supplied
82 molecule(s) are determined (Fig. 1a, Sec. 2.2.1) and the model compounds necessary to pa-
83 rameterize any missing terms are generated (Fig. 1b, Sec. 2.2.2). Rules based on chemical
84 topology are used for both of these steps to yield a unique dependency graph that can be
85 sorted (Fig. 1c, Sec. 2.2.3) to schedule the parameterization calculations. Assuming no pre-
86 vious parameters exist, parameterization (i.e., Stages 2 and 3) begins with simple molecules
87 like ethane and methanol which are at the base of the sorted dependency graph (Fig. 1c,
88 group 1) followed sequentially by larger molecules like ethanol, methoxyethane, and dimethyl
89 carbonate. There is a one-to-one mapping between force-field terms and model compounds,
90 such that each term is derived exclusively from the quantum chemistry training data of a
91 single model compound, which ensures the extensibility and backwards compatibility of the
92 force-field. At higher levels of the dependency graph, force-field parameters inherited from
93 model compounds at lower levels are held fixed during parameterization.

94 In Stage 2 (Fig. 1e), the data generation and force-field parameterizations associated

95 with intramolecular modes are performed. Each model compound is first initialized in a
96 canonical conformation (Sec. 2.3.1) then optimized at the target quantum chemistry level
97 of theory. The optimized geometry is then used as an input for constrained mode scans
98 of unique bonds, angles (Sec. 2.3.2), and dihedrals (Sec. 2.3.3). The intramolecular force-
99 field modes associated with the model compounds are then parameterized to the quantum
100 chemistry mode scans self-consistently with all other intramolecular parameters.

101 In Stage 3 (Fig. 1f), the data generation and force-field parameterizations associated
102 with intermolecular interactions are performed. Condensed-phase molecular dynamics is
103 used to sample molecular and pairwise configurations of each model compound (Sec. 1.1 in
104 the S.I.). Quantum chemistry calculations of electrostatic potentials and interaction energies
105 are performed on the molecular and pairwise configurations, respectively, and serve as the
106 reference data for parameterizing the intermolecular force-field terms (Sec. 2.4.1-2.4.2).
107 Finally, the intramolecular modes associated with the model compounds are refit to ensure
108 self-consistency with the final intermolecular terms (e.g., partial charges and Lennard-Jones
109 interactions).

110 Model compounds that are in the same group of the dependency graph are parameterized
111 in parallel during Stages 2 and 3. In the current example, the intramolecular and intermolec-
112 ular terms for methanol and ethane would be derived first, followed by the compounds in
113 group two (Fig. 1c), and so forth, until all parameters are obtained that are necessary to
114 simulate diethyl carbonate. The TAFFI database is updated at each step of the process to
115 avoid redundant calculations when parameterizing new molecules. For example, the force-
116 field terms associated with ethanol and ethane are at the base of the dependency graphs
117 of many potential organic species, but they are only evaluated once and then stored for all
118 future parameterizations.

2.2 Stage 1 - Organization of Calculations

Stage 1 of TAFFI consists of identifying the requisite force-field parameters (Fig. 1a, Sec. 2.2.1), generating model compounds for those parameterizations (Fig. 1b, Sec. 2.2.2), and ordering the parameterizations to ensure internal consistency (Fig. 1c, Sec. 2.2.3). Chemical topology (i.e., the molecular graph) plays a central role in Stage 1 for automating the assignment and parameterization of the force-field. The chemical topology can be expressed in a computationally useful form as an adjacency matrix, \mathbf{A} , with dimensions equal to the number of atoms in the molecule, and elements defined by

$$A_{ij} = \begin{cases} 1 & \text{if a bond exists between atom } i \text{ and atom } j \\ 0 & \text{if a bond does not exist between atom } i \text{ and atom } j. \end{cases} \quad (1)$$

Chemical topology is used in Stage 1 in three ways: (i) the definition of atom types, (ii) the definition of the model compounds, and (iii) for determining the molecular dependencies and order of calculations.

2.2.1 Definition of Atom Types

In TAFFI, the concept of an atom type is formalized based on the local molecular subgraph about each atom out to a specified number of bonded neighbors, d . In turn, bonds, angles, and dihedrals are uniquely defined based on the atom types involved in each mode. For the current work, a bond-depth $d = 2$ has been uniformly used for defining atom types. This choice enables TAFFI-gen to support a greater degree of chemical specificity than is present in other transferable force-fields (e.g., a mixture of $d = 1$ and $d = 2$ types are common depending on the available experimental parameterization data) while still being usefully transferable.

Atom typing in TAFFI occurs via breadth-first searches of the molecular graph out to d -bonds from the atom being typed. This procedure is seeded by querying the row of the

141 adjacency matrix (Eq. 1) corresponding to the atom being typed and identifying the atoms
142 bonded to it. This process is recursively applied $d-1$ additional times by reseeding the search
143 with the bonded atoms and excluding the atom seed from the previous generation to avoid
144 backtracking. The subgraphs obtained in this way uniquely define the atom types in the
145 molecule. TAFFI utilizes a string syntax for canonicalizing these subgraphs and expressing
146 them in a machine-readable format. In this syntax, all numbers refer to atomic numbers
147 (i.e., 1 corresponds to hydrogen, and 6 to carbon), open brackets (“[”) designate bonds, and
148 closed brackets (“]”) designate the end of bonded groups (i.e., either the point at which d
149 bonds is reached or at which a branch terminates). A bond is indicated between the atom
150 directly following the open bracket, “[”, and the first atom preceding the bracket that is not
151 enclosed by a “]”. The atom being typed is always designated first. For example, the atom
152 type of the central carbon atom in ethanol is encoded as `[6[6[1][1][1]][8[1]][1][1]]`, where the
153 first 6 refers to the central carbon atom itself, the `[6[1][1][1]]` refers to the bonded methyl, the
154 `[8[1]]` refers to the bonded alcohol, and the final `[1][1]` specifies the two hydrogens directly
155 bonded to the central carbon. To resolve the ambiguity associated with graph isomorphism,
156 the ordering of branches within each atom type is determined by the mass of the bonded
157 atoms, followed by the mass and number of next-nearest bonded atoms (similar to Cahn-
158 Ingold-Prelog priority rules). Labels for unique bond, angle, and dihedral types are defined
159 based on the atom types involved in each mode (e.g., `[6[6[1][1][1]][8[1]][1][1]] [1[6[8][6][1][1]]]`
160 is the bond type associated with the C-H bond about the central carbon atom in ethanol).

161 2.2.2 Definition of Model Compounds.

162 In TAFFI, all force-field parameters are derived from a set of algorithmically generated model
163 compounds for which reference quantum chemistry data can be generated. For a given force-
164 field term (e.g., a partial charge, bond type, angle type, etc.), the model compound is defined
165 as the smallest acyclic molecule that both exhibits the required force-field term and conserves
166 the Lewis structure of the associated atom types. For example, as shown in Fig. 1b for $d = 2$,

167 the model compound used to parameterize the partial charges of the terminal alkyl hydrogen,
168 [1[6[6][1][1]]], is ethane, because ethane is the smallest molecule containing that atom type.

169 Starting with the target compound supplied by the user, these model compounds are
170 generated in two steps. First, all atoms more than d bonds away from the targeted term
171 are removed to form a preliminary compound. For atom types, bond types, and angles, this
172 means truncating all atoms more than $d + 1$ bonds away from any atom involved in the
173 targeted mode. For dihedrals, this means truncating all atoms more than $d + 1$ bonds away
174 from the atoms defining the rotatable bond (i.e., the 2-3 atoms of the dihedral). Second,
175 any undercoordinated atoms that result from this truncation are hydrogenated to a level
176 that is consistent with the hybridization of the subgraph and necessary to form a valid
177 Lewis structure. We emphasize that the resulting model compounds are independent of the
178 specific user-supplied structure that initiated their generation. That is, each force-field term
179 is parameterized using a unique model compound, and the user-supplied structures only play
180 a role in identifying force-field terms in need of parameterization.

181 This definition of model compounds has two shortcomings that we note here but leave
182 to future work to address. First, this definition leads to ambiguity in cases involving double
183 bonds between nearest and next-nearest neighbors of the atoms associated with the force-field
184 term (e.g., keto-enol tautomers). In these cases, double bonds with the highest bond energy
185 are preferentially formed in the model compound.⁵⁵⁻⁵⁷ For example, the model compound
186 for the atom type [6[6[1][1][1]][6[8][6]][1][1]] is 2-butanone rather than 1-buten-2-ol (i.e., the
187 ketone as opposed to the alcohol, consistent with the Erlenmeyer rule). This ambiguity
188 could be addressed in the future by introducing bond-orders into the atom labels (e.g.,
189 using distinct symbols for double and triple bonds instead of specifying bonds generically
190 with '[' and ']') such that distinct tautomers would be parameterized to distinct model
191 compounds. Second, this definition leads to force-field terms associated with cyclic structures
192 being derived from data for acyclic model compounds. We note that rings, and similarly
193 conjugated groups, have intrinsically non-local contributions to their configurational energy

194 that represents a challenge to the locality assumption of any force-field based on atom types.
195 This could be addressed in the future by using model compounds for rings and conjugated
196 subunits that preserve these components, but this is outside of the scope of the current study.

197 It may happen that the model compounds exhibit new force-field terms that are distinct
198 from the parent molecule. Thus, model compound generation is recursively performed for
199 these new force-field terms until all model compounds have been generated for all unknown
200 terms. Because each model compound is smaller than its parent, this recursion will eventually
201 terminate with small model compounds containing approximately d non-hydrogen atoms.
202 This procedure yields model compounds that are generally small and amenable to high-level
203 quantum chemistry calculations. For example, 90% of model compounds generated in this
204 study had six or fewer heavy atoms (the mode is four), and no model compound had greater
205 than eight heavy atoms (Fig. S1).

206 **2.2.3 Definition of the Dependency Graph.**

207 The recursive generation of model compounds creates dependencies based on shared force-
208 field terms. To account for these dependencies, it is necessary to order data generation
209 and parameterizations (Subsections 2.2-2.3) such that all force-field terms, besides those
210 associated with a given model compound, have been obtained prior to performing each
211 parameterization. These dependencies are enumerated during model compound generation
212 and stored in a dependency graph. The dependency graph has nodes for all model compounds
213 and directed connections between all dependent compounds (e.g., ethanol depends on ethane
214 for the partial charges of atom type [1[6[6][1][1]]], but ethane does not depend on ethanol,
215 Fig. 1c). Prior to performing force-field parameterizations, a topological sort is applied to
216 the dependency graph such that no dependencies exist within the same level of the sorted
217 graph. Data generation and parameterization (Stages 2 and 3) are then performed beginning
218 with model compounds in the bottom level of this graph and working to the top (i.e., level
219 1 to level 4 in Fig. 1c). This addresses the issue of force-field terms potentially being

220 missing during parameterization because the terms at each level can be directly determined
 221 when all of the dependent terms in the lower levels of the dependency graph are known.
 222 The algorithm for model compound generation (Sec. 2.2.2) in TAFFI has the important
 223 property that dependent model compounds are always identical to or smaller than their
 224 parent molecule. Consequently, the dependency graph for any molecule is directed and
 225 acyclic, and it is always possible to order calculations such that all dependencies exist at the
 226 time of parameterization.

227 2.2.4 Force-field Expression

While the particular force-field expression used for fitting the data in the TAFFI database is flexible, this choice does guide which calculations are performed on the model compounds in the subsequent stages. For the current study, we employ the following fixed-charge functional form:

$$V_{\text{FF}} = \sum_{\text{bonds}} k_r (r - r_0)^2 + \sum_{\text{angles}} k_\theta (\theta - \theta_0)^2 + \sum_{\text{dihedrals}} \sum_{i=1}^4 \frac{1}{2} V_i (1 + (-1)^{i+1} \cos(i\phi)) + \sum_{i>j} \left\{ \frac{q_i q_j e^2}{4\pi \epsilon_0 r_{ij}} + 4\epsilon_{ij} \left[\left(\frac{\sigma_{ij}}{r_{ij}} \right)^{12} - \left(\frac{\sigma_{ij}}{r_{ij}} \right)^6 \right] \right\}, \quad (2)$$

228 where k_r and r_0 are a bond-specific force constant and equilibrium displacement, respec-
 229 tively; k_θ and θ_0 are an angle-specific force constant and equilibrium angle, respectively; the
 230 V_i terms are dihedral-specific Fourier coefficients, r_{ij} are the interatomic separations, q_i are
 231 the atomic partial charges, e is the elementary charge, and ϵ_{ij} and σ_{ij} are the Lennard-Jones
 232 (LJ) parameters for each pairwise interaction. The summation for the Lennard-Jones and
 233 Coulomb potentials runs over all intermolecular atomic pairs and all intramolecular atomic
 234 pairs separated by more than three bonds (i.e., 1-4 intramolecular interactions are excluded).
 235 All dihedrals that rotate about double-bonds are modeled as invertible harmonic modes by
 236 only using the $i = 2$ term in the dihedral expression. These functional forms are largely
 237 standardized in general force-fields and are broadly implemented in existing MD packages,

238 which makes them an obvious starting point for this initial benchmark study.

239 **2.3 Stage 2 - Intramolecular Parameterizations**

240 Stage 2 consists of generating reference quantum chemistry data and performing force-field
241 parameterizations related to the intramolecular force-field parameters (Sec. 2.3.1-2.3.3).
242 Parameterizing intramolecular modes is a prerequisite for generating reference data for in-
243 termolecular force-field terms in Stage 3. Thus, Stage 2 occurs first to yield a provisional
244 force-field, with the final intramolecular force-field terms refit after the Stage 3 intermolecular
245 terms.

246 **2.3.1 Conformer Generation**

247 The first step in generating reference quantum chemistry data for fitting intramolecular
248 force-field terms is generating an optimized geometry for the model compounds. Here, all
249 model compounds are initialized as the conformer with *trans* relationships for all backbone
250 dihedrals (i.e., the all-*trans* conformer). The all-*trans* conformer is generated by (i) iden-
251 tifying the atoms belonging to the longest connected path in the molecular graph (i.e., the
252 molecular backbone), (ii) aligning the backbone dihedrals in all-*trans* geometries, and (iii)
253 repeating with the remaining branches of the molecular graph until all non-terminal dihe-
254 drals exhibit *trans* relationships. Since the all-*trans* designation leaves the conformation of
255 terminal dihedrals ambiguous (e.g., the dihedral involving chlorine in 1-chlorobutane), the
256 conformation of end groups is determined by explicitly generating and optimizing all end
257 group conformers by steepest descent using the Universal Force Field (UFF),⁵⁸ then using
258 the lowest energy conformer as the input structure for quantum chemical geometry optimiza-
259 tion. This procedure yields a deterministic conformer and initial geometry for each model
260 compound.

261 2.3.2 Parameterization of Harmonic Modes

262 Bonds, angles, and dihedrals about double bonds are modeled here with harmonic forms
 263 (Eq. 2). In cases where a model compound has multiple resonance structures, if a dihedral
 264 has a double bond in any resonance structure then it is modeled as a harmonic mode (e.g.,
 265 all dihedrals in benzene are considered harmonic in TAFFI-gen). All harmonic modes are
 266 self-consistently fit to constrained quantum chemistry mode scans. Bond mode scans consist
 267 of compression and extension by 0.1 pm about the optimized bond length in steps of 0.02 pm.
 268 Angle mode scans consist of compression and expansion by 0.5° about the optimized angle
 269 in steps of 0.1°. Harmonic dihedral scans consist of compression and expansion by 0.5°
 270 about the optimized dihedral angle in steps of 0.1°. At each scan configuration, geometry
 271 optimizations are performed with the mode being parameterized constrained to a fixed value
 272 while optimizing all remaining degrees of freedom.

273 The harmonic modes associated with the model compounds are parameterized to mini-
 274 mize the following objective function:

$$\chi_{\text{harm}}^2 = \sum_i \left(E_{\text{QC},i} - \sum_{\nu_j \in \text{local}} V_{\text{FF}}(\nu_j) \right)^2, \quad (3)$$

275 where the index i runs over all scanned configurations, $E_{\text{QC},i}$ is the single-point energy
 276 of configuration i relative to the minimum-energy configuration, the index j runs over all
 277 bonds, angles, and harmonic dihedrals that share an atom with the scanned mode (i.e.,
 278 “local” modes), and $V_{\text{FF}}(\nu_j)$ is the force-field energy of mode ν_j in configuration i . The self-
 279 consistent fit over all local modes is performed because the force-field terms are generally
 280 not linearly independent. All fits are performed using the limited-memory Broyden-Fletcher-
 281 Goldfarb-Shanno algorithm with bound constraints (L-BFGS-B) to limit the fit variables to
 282 positive values. Initial guesses for the force-constant and equilibrium displacement for each
 283 scanned mode are obtained by a linear least-squares fit to the quantum chemistry single-
 284 point energies with respect to the mode being fit. This procedure is repeated until reaching

285 self-consistency among all intramolecular modes. During these fits, only the force-field terms
 286 associated with the model compound are parameterized, and any terms inherited from model
 287 compounds lower in the dependency graph are held fixed.

288 2.3.3 Parameterization of Flexible Dihedral Potentials

289 Dihedrals that rotate about single and triple bonds are modeled by TAFFI-gen with a
 290 truncated Fourier series. All flexible dihedrals are self-consistently fit to constrained quantum
 291 chemistry scans from $[0, 2\pi)$ and $[0, -2\pi)$, in 5° steps, about each rotatable bond. Two scans
 292 are performed to mitigate the path-dependence of the scan (e.g., this can be important
 293 for sterically crowded dihedrals) and the lowest energy union of the two scans is used as
 294 reference data for the parameterization. During each scan, the dihedral being parameterized
 295 is constrained to a fixed value while optimizing all remaining degrees of freedom. In the case
 296 where multiple dihedrals exist about the same bond, only the dihedral involving the heaviest
 297 atoms—or secondarily, the longest chain—is explicitly constrained during the scan.

298 The Fourier coefficients are fit to minimize the residual between the quantum chemistry
 299 and force-field potentials for the constrained dihedral rotation according to the following
 300 objective function:

$$\chi_{\text{Fourier}}^2 = \sum_i \left(E_{\text{QC},i} - \sum_{\nu_j \notin \text{fit}} E_{\text{FF},i}(\nu_j) - \sum_{\nu_j \in \text{fit}} \sum_{k=1}^4 \frac{1}{2} V_{j,k} (1 + (-1)^{k+1} \cos(k\phi_{i,j})) \right)^2 + \omega_{\text{L2}} N_{\text{fit}}^{-1} \sum_{i,j \in \text{fit}} V_{i,j}^2, \quad (4)$$

301 where the index i runs over all scan configurations, $E_{\text{QC},i}$ is the single-point energy of the
 302 configuration, the second summation runs over all force-field terms that are not being fit
 303 (i.e., bonds, angles, unscanned dihedrals, electrostatics, and Lennard-Jones terms), the third
 304 summation runs over all dihedrals that share the scanned bond (i.e., $\nu_j \in \text{fit}$), $V_{j,k}$ are the
 305 dihedral-specific force constants, and $\phi_{i,j}$ is the angle of dihedral j in configuration i . The last
 306 summation is an L2 regularization of the average magnitude of the dihedral fit coefficients
 307 that reduces overfitting to noisy data. ω_{L2} is set to 0.1 percent of the range of the fit

308 values (i.e. the difference between $E_{\text{QC},i}$ and the second summation in Eq. 4). All fits are
 309 performed using the L-BFGS-B algorithm with bound constraints limiting the magnitude of
 310 the dihedral fit coefficients to two hundred percent of the range of fit potential.

311 During Stage 2, the Lennard-Jones parameters and partial charges are not yet determined,
 312 so UFF parameters and approximate partial charges fit to the optimized geometry of the
 313 model compound (Sec. 2.4.1) are used as an approximation. After Stage 3, all intramolecular
 314 parameters are refit with updated partial charges and Lennard-Jones parameters using the
 315 same procedure.

316 2.4 Stage 3 - Intermolecular Parameterizations

317 Stage 3 consists of generating reference quantum chemistry data and performing force-field
 318 parameterizations related to the intermolecular force-field parameters (Sec. 2.4.1-2.4.2).
 319 Configurational sampling is critical for generating reference data for intermolecular terms,
 320 which requires Stage 3 to occur after a preliminary force-field is obtained from Stage 2.
 321 After configurational sampling (Sec. 1.1 in the S.I.) , quantum chemistry calculations on
 322 molecular and pairwise configurations are used to parameterize the partial-charges (Sec.
 323 2.4.1) and Lennard-Jones parameters (Sec. 2.4.2), respectively.

324 2.4.1 Parameterization of Partial Charges

325 The electric potential calculated on a grid about each molecule in each sampled configuration
 326 (see Sec. 1.1 in the S.I.) is used as reference data for the partial charge parameterization.
 327 The partial charges are fit to minimize the following objective function:

$$\chi_q^2 = \sum_s^{N_{\text{samples}}} \left(\omega_{\text{pot}} N_{\text{pot}}^{-1} \sum_i^{N_{\text{pot}}} (V_{\text{QC},i} - V_{\text{FF},i})^2 + \omega_{\text{D}} \sum_i^3 (D_{\text{QC},i} - D_{\text{FF},i})^2 + \omega_{\text{T}} \left(\sum_i^{N_{\text{atoms}}} q_i - q_{\text{T}} \right)^2 \right), \quad (5)$$

328 where the first summation (s) is over the sampled configurations, the second summation is
 329 over the squared deviations of the force-field description ($V_{\text{FF},i}$) from the reference electric

330 potential ($V_{\text{QC},i}$) as calculated on the N_{pot} grid points, the third summation corresponds to
331 the element-wise deviations of the force-field description ($D_{\text{FF},i}$) from the reference molecular
332 dipole (D_{QC}), and the fourth summation corresponds to deviations from the total molecular
333 charge (q_{T}). ω_{pot} , ω_{D} , ω_{T} are weighting coefficients for penalizing the electric potential,
334 dipole, and total charge deviations, respectively. The s index is implied in all terms, but
335 dropped for clarity. Partial charges (q_i) are fit using $\omega_{\text{pot}} = 1.0$, $\omega_{\text{D}} = 0.1$, $\omega_{\text{T}} = 1.0$, specified
336 in inverse atomic units. As implemented in ORCA v.4.1.2, the electric potential is calculated
337 on a cubic grid with a grid spacing of 0.3 Å, and any grid points further than 2.8 Å from
338 any atom or within the COSMO radius of any atom are discarded.

339 The partial charges are fit in two steps. First, Eq. 5 is minimized while constraining
340 polar atoms of identical TAFFI atom type to have the same partial charge. Polar atoms are
341 considered to be any non-hydrogen atoms besides carbon, and hydrogen atoms that are not
342 bonded to carbon. A second fit is then performed by minimizing Eq. 5 while holding the
343 partial charges for the polar atom types constant and constraining all non-polar atoms of
344 the same type to have the same partial charge. This two step procedure is similar to the
345 RESP algorithm⁵⁹ and is meant to improve the accuracy of the electric potential near the
346 polar atoms. This procedure differs from the RESP algorithm in (i) the form of the objective
347 function and (ii) the use of 200 configurations rather than a single configuration. We note
348 that fitting to multiple configurations tends to reduce the magnitude of the partial charges,
349 alleviating the need for the heuristic hyperbolic restraint used in RESP. Partial charge fits
350 are performed using the BFGS algorithm.

351 2.4.2 Parameterization of Pairwise Interactions

352 Counter-poise corrected interaction energies (IE) of the sampled pairwise configurations (see
353 Sec. 1.1 in the S.I.) are used as reference data for the Lennard-Jones parameterization. The

354 Lennard-Jones parameters are fit to minimize the following objective function:

$$\chi_{\text{LJ}}^2 = \omega_{\text{IE}} N_{\text{IE}}^{-1} \sum_i^{N_{\text{IE}}} (IE_{\text{QC},i} - IE_{\text{FF},i})^2 + \omega_{\epsilon} N_{\epsilon}^{-1} \sum_i^{N_{\epsilon}} (\epsilon_{\text{UFF},i} - \epsilon_{\text{FF},i})^2 + \omega_{\sigma} N_{\sigma}^{-1} \sum_i^{N_{\sigma}} (\sigma_{\text{UFF},i} - \sigma_{\text{FF},i})^2, \quad (6)$$

355 where the first summation corresponds to squared deviations of the force-field interaction en-
 356 ergy (IE_{FF}) from the counter-poise corrected interaction energy (IE_{QC}) over all N_{IE} pairwise
 357 samples, the second summation corresponds to the L2 regularization of the Lennard-Jones
 358 energy parameters ($\epsilon_{\text{FF},i}$) with respect to the UFF reference values ($\epsilon_{\text{UFF},i}$), and the third
 359 summation corresponds to the L2 regularization of the Lennard-Jones atomic radii ($\sigma_{\text{FF},i}$)
 360 with respect to the UFF reference values ($\sigma_{\text{UFF},i}$). The latter terms in the objective function
 361 are included to avoid extreme values in ϵ and σ that can occur when using only a least-
 362 squares objective function. The Lennard-Jones parameters are fit using $\omega_{\text{IE}} = 1.0$ mol/kcal
 363 $\omega_{\epsilon} = 1.0$ mol/kcal, and $\omega_{\sigma} = 0.1 \text{ \AA}^{-1}$. A comparison of the interaction energies calculated at
 364 the UFF level and with the regularized and unregularized TAFFI-gen parameters (Fig. S2)
 365 confirms that the regularization terms have only a small effect on the reproduction of the
 366 interaction energies. The interaction energies are calculated in the force-field description as
 367 the sum of intermolecular Lennard-Jones and electrostatic terms between the molecules in
 368 each configuration. The partial charges are held fixed during the fitting of the Lennard-Jones
 369 parameters. Any configurations with unstable interaction energies (i.e., $IE_{\text{QC}} > 0$ kcal/mol)
 370 are excluded from the fit. Lennard-Jones fits are performed using the L-BFGS-B algorithm.

371 2.5 Dataset Description

372 LAMMPS⁶⁰ and ORCA v.4.1.2⁶¹ were used to perform the molecular dynamics simulations
 373 and quantum chemistry calculations, respectively, associated with reference data generation.
 374 All quantum chemistry calculations were performed at the $\omega\text{B97X-D3}^{62}/\text{def2-TZVP}^{63,64}$ level
 375 of theory for training the version of TAFFI-gen reported here.

376 To assess the performance of TAFFI-gen, we present a benchmark on the dataset of

377 small organic molecules introduced by Caleman et al. for GAFF and OPLS-AA.⁶⁵ The
378 original MD-based predictions of liquid properties by Caleman included 147 molecules in
379 their benchmark set. In the current study, we have excluded ring, nitro, and phosphate
380 containing compounds, as they require a more sophisticated treatment of atom types and
381 model compounds that is beyond the scope of the current work. After these exclusions, a
382 total of 87 molecules at 146 distinct state points (i.e., multiple temperatures per molecule
383 where included by Caleman et al.) are in the presented benchmark. A list of all bench-
384 mark compounds and individual property predictions are distributed in the supplementary
385 information of this work.

386 Six properties were calculated from the MD trajectories: the density, enthalpy of va-
387 porization, static dielectric constant, volumetric thermal expansion coefficient, isothermal
388 compressibility, and quantum-corrected heat capacity at constant volume. Following the
389 reference benchmark by Caleman, three types of MD simulations were performed to extract
390 these properties. Gas phase simulations were run to obtain the expected potential energy per
391 molecule in the gas phase for the enthalpy of vaporization calculation. Relatively long liquid
392 phase simulations (i.e., 10 ns) in the NPT ensemble were run to compute all properties other
393 than the heat capacity. Short liquid phase simulations (i.e., 100 ps) were run in the NVT
394 ensemble with high sampling frequency to calculate the constant volume heat-capacity using
395 the two-phase method.^{66,67} Details of the simulation and analysis methods are described in
396 the SI. We note that the dielectric constants of methanoic acid have been omitted in analysis
397 due to lack of convergence, which will be revisited in the discussion. Besides this case, all
398 available experimental data in Ref. 65 for the benchmark molecules has been included for
399 comparison.

400 Finally, four error measures are reported for comparing the results for TAFFI-gen against
401 experimental data and the other force-fields (Eq. 7-10). The mean absolute difference (MAD)
402 is calculated as

$$\text{MAD} = \frac{1}{N} \sum_i^N |x_{i,\text{sim}} - x_{i,\text{ref}}| \quad (7)$$

403 where N is the total number of data points, $x_{i,\text{sim}}$ is the simulated value for each data point
404 and $x_{i,\text{ref}}$ is the corresponding reference value (DFT calculated value or experimental value).
405 The mean signed difference (MSD) is calculated as

$$\text{MSD} = \frac{1}{N} \sum_i^N x_{i,\text{sim}} - x_{i,\text{ref}} \quad (8)$$

406 with positive values indicating an average overestimation of the value by simulations. The
407 mean absolute percent difference (MAPD) is calculated as

$$\text{MAPD} = \frac{100}{N} \sum_i^N \frac{|x_{i,\text{sim}} - x_{i,\text{ref}}|}{x_{i,\text{ref}}} \quad (9)$$

408 The mean signed percent difference (MSPD) is calculated as

$$\text{MSPD} = \frac{100}{N} \sum_i^N \frac{x_{i,\text{sim}} - x_{i,\text{ref}}}{x_{i,\text{ref}}} \quad (10)$$

409 We note that MAD and MSD are more sensitive to the large magnitude samples in the
410 dataset, whose deviations tend to be correspondingly larger than the small magnitude sam-
411 ples. MAPD and MSPD are more sensitive to the small magnitude samples in the dataset.

412 **3 Results and Discussion**

413 TAFFI-gen is parameterized to DFT reference data for small model compounds. Thus,
414 the errors in TAFFI-gen predictions can be decomposed into errors associated with the un-
415 derlying DFT parameterization data and representability errors associated with the limited
416 functional form of the force-field. Regarding the first source of error, the dispersion-corrected
417 range-separated functional used here is among the highest performing in benchmarks of con-
418 formational energetics and cluster interactions for organic species.^{62,68-72} Nevertheless, even
419 modern functionals have documented deficiencies for aqueous solutions and reaction barriers

420 that would require higher fidelity training data for models of water or reactive force-fields,
421 which are beyond the present scope. Thus, for the current study, we acknowledge this poten-
422 tial source of error but consider it negligible in comparison with the representability errors
423 associated with the simple functional form of the force-field.

424 To quantify the magnitude of errors associated with the functional form of the force-field,
425 we have compared the TAFFI-gen predictions for normal modes and optimized geometries
426 against DFT results for the benchmark compounds (Fig. 2). Comparing the normal mode
427 frequencies provides a measure of the accuracy of forces in the force-field representation
428 (Fig. 2a). We observe a MAD of 52 cm^{-1} and MAPD of 6%, which is comparable to non-
429 transferable quantum chemistry derived force-fields using more complex forms.^{39,40} This
430 suggests that in general TAFFI-gen exhibits accurate force-behavior near equilibrium struc-
431 tures. Notably, the largest percent deviations are associated with low frequency modes
432 ($< 1000\text{ cm}^{-1}$), which is expected given the lack of explicit coupling between dihedral terms
433 and the exclusion of improper modes in the current force-field.

434 The predicted equilibrium structures of the benchmark compounds provides a second
435 point of comparison between TAFFI-gen and the reference DFT level of theory (Fig. 2b).
436 These comparisons are performed by optimizing the compounds at the DFT and force-
437 field levels starting from the same all-trans conformer, then aligning the structures via the
438 Kabsh algorithm. First, we observe that the deviations of atomic positions (MAD= 0.1Å),
439 bonds lengths (MAD= 0.002Å), and bending angles (MAD= 0.7°) are all extremely small on
440 a per molecule basis, which confirms the generally excellent agreement between TAFFI-gen
441 and DFT for local structural features. Larger deviations are observed for proper dihedrals
442 (MAD= 7°) and improper dihedrals (MAD= 6°). From the distribution of proper dihedral
443 deviations, it is evident that these errors are driven by a small number of outliers that
444 adopt distinct conformers at the TAFFI-gen level upon geometry optimization. In particu-
445 lar, terminal methyl groups proximate to esters and amides tend to twist relative to DFT
446 predictions (Fig. S4), which occurs for methyl acetate (dihedral MAD= 33°), methyl formate

447 (36°), acetyl acetate (37°), N,N-dimethylacetamide (34°), N-methylformamide (36°), and N-
448 methylacetamide (45°). In contrast, the errors in improper dihedrals appear to be systematic,
449 with a relatively large standard deviation in MAD across the reference structures (5.99°).
450 This is a consequence of not explicitly including improper modes in the force-field form.
451 The errors in improper angles are intuitively largest for planar conjugated units. For example, the
452 largest error is exhibited by the improper defined about the carbonyl in 2,6-dimethylheptan-
453 4-one, where TAFFI-gen exhibits an improper angle of 32° in contrast to 0° predicted by
454 DFT. The optimized geometries for DFT and TAFFI-gen for the molecules with large MADs
455 are compared in Fig. S4. Although we have focused on the largest error cases to illustrate
456 the limitations of the common force-field form employed here, the overall mean performance
457 is nevertheless very accurate (Table 1). Namely, the overwhelming majority of structural fea-
458 tures are quantitatively reproduced by TAFFI-gen and the cases where incorrect conformers
459 are stabilized are rare and isolated to the periphery of the molecules.

460 We note that the above comparison has been performed for the benchmark molecules
461 and not for the model compounds actually used for TAFFI-gen parameterization. Fig. S3
462 presents the analogous comparisons with DFT results for model compounds only, which
463 show very similar deviations compared with the benchmark structures. The similar errors
464 observed between these two cases provides evidence that the $d = 2$ atom typing of TAFFI-
465 gen leads to excellent transferability between model compounds and larger molecules for
466 structural features, while the limited representability of the force-field is the main source of
467 error with respect to the DFT reference data.

468 MD Simulations of six liquid properties were performed to establish the performance
469 of TAFFI-gen relative to OPLS-AA and GAFF in predicting experimental liquid properties
470 (Fig. 3). These properties include the density (ρ), heat of vaporization (ΔH_{vap}), static dielec-
471 tric constant (ϵ), volumetric thermal expansion coefficient (α_P), isothermal compressibility
472 (κ_T), and quantum-corrected heat capacity at constant volume (c_v) for the 87 molecules in
473 the current benchmark. We note that among the liquid properties, ρ , ΔH_{vap} , and ϵ have

474 historically been utilized as part of the OPLS-AA and GAFF parameterizations, whereas
475 for TAFFI-gen this data is not utilized in any way and represents a test for the force-field.
476 Summary statistics across the benchmark are presented in Table 2, and the TAFFI-gen
477 predictions for individual simulation conditions are presented in Table S1.

478 The summary error statistics calculated across all systems for each force-field illustrates
479 the similar accuracy (and inaccuracy) of the three force-fields for the various properties.
480 Although some specific differences occur, which are discussed below, it is perhaps surprising
481 that the mean performance is so consistent despite the distinct parameterization protocols
482 and training data for the three force-fields. For instance, all of the force-fields exhibit rela-
483 tively small errors for ρ and c_v , large systematic errors for ε (e.g., $R^2 < 30\%$ in all cases),
484 and high correlation but large variances for ΔH_{vap} , α_{P} , and κ_{T} . These trends can be ratio-
485 nalized by the common functional form of these force-fields. For instance, the Lennard-Jones
486 potential is capable of recapitulating the molecular volume, which is the leading order con-
487 tribution to density, but is an approximate description of van der Waals interactions which
488 significantly contribute to ΔH_{vap} . Similarly, a fixed point-charge model is an aggressive
489 simplification of electrostatic interactions, which explains the poor dielectric results in all
490 cases, and also contributes to the high variances of the other fluctuation-based condensed
491 phase properties. The heat capacity is well reproduced in all cases, which is also consis-
492 tent with the generally accurate reproduction of local configurational energetics (i.e., bond,
493 angle, and to a lesser degree dihedral terms) in these force-fields. Thus, despite their in-
494 dependent reference data, the force-fields exhibit similar average prediction behaviors that
495 reflect the representability limitations of the functional form of the force-field. The approx-
496 imate treatment of intermolecular interactions, in particular, leads to shared trade-offs in
497 reproducing thermodynamic properties. This is further evidenced by the fact that interac-
498 tion energy errors in TAFFI-gen exhibit the largest variance of all training properties (Fig.
499 S2). Specifically, while TAFFI-gen exhibits excellent reproduction of the mean interaction
500 energies (MSE of -0.09 kcal/mol for the model compound training data), the error residuals

501 exhibit very long tails (kurtosis=20.25) which is clear evidence of representability limitations
502 associated with the pairwise fixed-charge form of the force-field.

503 Although our interpretation of the similar mean performance of the three force-fields is
504 that representability limitations dominate the general behaviors, this does not exclude some
505 specific cases being the result of inaccurate parameterizations. For instance, the efforts of
506 the Open Force-Field Consortium have highlighted many cases where additional accuracy
507 can be squeezed from fixed-charge force-fields by refining specific parameters.^{15,73-77} Like-
508 wise, the fact that OPLS generally outperforms the other force-fields illustrates that the
509 specific force-field terms for TAFFI-gen might be improved by tuning the parameterization
510 hyperparameters or supplementing the training data.

511 To facilitate a more fine-grained comparison between the force-fields, the MAPD with
512 respect to each liquid property is presented on a per functional group basis in Figure 4.
513 Molecules were included in a category if they exhibit the specified functional group; thus,
514 some molecules are included in multiple categories. We have also combined similar functional
515 groups in some cases to avoid scarce or empty categories. We note that experimental data
516 is not available for all compounds for all properties, thus the number of compounds in each
517 category varies across properties, and bars have been omitted for cases where less than three
518 datapoints were available. We note that a large number of distinct outliers are observable
519 for GAFF that have previously been discussed by Coleman et al., and are thus not further
520 remarked on here.

521 There are several informative outliers observed for all of the force-fields that shed further
522 light on representability limitations. For example, all of the force-fields exhibit underesti-
523 mated dielectric constants for the amides, which suggests the need for polarizable terms to
524 accurately account for the large molecular dipoles and strong hydrogen bonding associated
525 with this functional group.⁷⁸ Another noticeable trend is large overestimations of the volu-
526 metric thermal expansion coefficient and isothermal compressibility for the halides, which are
527 mainly driven by small molecules with multiple halogens such as chloroform (>48/74% devia-

528 tions, respectively), dichloro(fluoro)methane ($>49\%/n.a.$), 1,1-dichloroethene ($>43\%/n.a.$),
529 and 1,1,2,2-tetrachloroethane ($>14/34\%$). There is also a trend for the heat capacity of
530 halides to be underestimated (on average by $>23\%$). It is known that halogens often exhibit
531 anisotropic charge distributions with a positive electrostatic potential on the outermost part
532 of the halogen, which cannot be accurately described using fixed-charge models.^{79,80} Based
533 on this understanding, various models have been developed for halides that include a virtual
534 site with positive charge,^{15,81-85} multipole electrostatics,^{86,87} polarizability,⁸⁷⁻⁸⁹ and angular-
535 dependent LJ terms.⁹⁰

536 A distinct outlier for TAFFI-gen is diethyl carbonate, which exhibits a large density
537 underestimation in comparison with experiment (MSPD = -18% ; this is the outlier visible
538 in Fig. 3a at $\rho_{\text{exp}} \sim 0.9$). This is the only carbonate in the benchmark, and carbonates
539 are unique in that they are the only benchmarked functional group that extends beyond
540 the $d = 2$ graph specificity explored here for TAFFI-gen. In particular, the $d = 2$ model
541 compound for the backbone oxygens (atomtype [8[6[6][1][1]][6[8][8]]) is ethoxyformic acid,
542 which fails to preserve the carbonate structure. The other benchmarked properties of diethyl
543 carbonate are also relatively poorly reproduced ($\Delta H_{\text{vap}}:-21\%$, $\varepsilon:-19\%$, $\alpha_{\text{P}}:86\%$, $\kappa_{\text{T}}:200\%$, $c_{\text{v}}:-$
544 10%), which we attribute to the poor congruence between the model compounds and the
545 target carbonate. This is further confirmed by an experiment where we reparameterized the
546 diethyl carbonate LJ force-field terms for the ether oxygens and the carbonate carbon with
547 ethyl methyl carbonate, which preserves the carbonate. In this case, the errors in comparison
548 with experiment are much smaller ($\rho:-1\%$, $\Delta H_{\text{vap}}:4\%$, $\varepsilon:-2\%$, $\alpha_{\text{P}}:34\%$, $\kappa_{\text{T}}:25\%$, $c_{\text{v}}:-4\%$). This
549 is a revealing example of how a fixed graph specificity (i.e., $d = 2$ in the current study) can
550 lead to non-systematic errors when applied to large functional groups.

551 Methanoic acid is also a distinct outlier for all of the force-fields. This system exhibits
552 long correlation times for the system dipoles, which have been previously established to
553 originate from strong dimer interactions.^{65,91} For TAFFI-gen, the dipole correlation decay
554 could not be converged even with longer 50 ns trajectories (not shown). Additionally, the

555 overestimation of the heat of vaporization for the ketone, aldehyde, and carboxylic acid
556 group is disproportionately affected by methanoic acid (>110% deviation), where the other
557 outliers are relatively minor [1-methoxy-2-(2-methoxyethoxy)ethane (>30%), and pentane-
558 2,4-dione (>35%)]. Excluding methanoic acid from the group for heat of vaporization results
559 in the MAPD values similar to other oxygen-containing functional groups (GAFF:20.24%,
560 OPLS-AA:13.94% and TAFFI:27.58%). This is an illustrative case of how fixing the force-
561 field complexity does not lead to systematic errors across distinct chemistries. To achieve
562 a target accuracy for a given set of properties, it is possible to simplify the force-field in
563 some cases, while it is necessary to add complexity in others. The development of more
564 sophisticated models for hydrogen-bonding in methanoic acid indirectly substantiates this
565 point.⁹¹⁻⁹⁵

566 As noted by Caleman et al., there are also cases where the simulation conditions may
567 exacerbate prediction errors in comparison with experiment. For example, the benchmarks
568 for some alcohols and amines, including propane-1,2,3-triol and (2-hydroxyethoxy)ethan-2-
569 ol, and ethane-1,2-diamine, are performed near their melting point. This results in highly
570 viscous liquids at the simulation temperatures (Table S1) and likely exacerbates errors in the
571 fluctuation-derived properties that are not representative of simulations further away from
572 the phase transition.

573 4 Conclusions

574 It would be useful to have a force-field framework that could bridge simple fixed-charge
575 force-fields on the one hand and complex machine learning force-fields on the other. The
576 present work takes the first step in this direction by establishing a parameterization frame-
577 work (TAFFI) based on an extensible quantum chemistry dataset that can be used to fit
578 transferable force-fields of varying complexity. With the TAFFI framework we have formal-
579 ized the concept of atom typing and made it the basis for generating systematic training

580 data that maintains a one-to-one correspondence with force-field terms. This feature makes
581 TAFFI arbitrarily extensible to new chemistries while maintaining internal consistency and
582 transferability. As a demonstration of TAFFI, we have developed a fixed-charge force-field,
583 TAFFI-gen, from scratch that includes coverage for many common organic moieties. The
584 performance of TAFFI-gen was benchmarked against OPLS-AA and GAFF for reproducing
585 several experimental properties of 87 organic liquids. The comparable accuracy between
586 TAFFI-gen and existing force-fields in this benchmark is quite encouraging in light of the
587 decades of optimization the existing force-fields have undergone and their use of experimental
588 data. Nevertheless, a major conclusion from this case-study is that the similar qualitative be-
589 haviors of these force-fields reflects the representability limitations of their simple functional
590 form in approximating the atomistic PES. In particular, similar trade-offs and inaccuracies
591 are observed in all of the force-fields which motivates a more sophisticated treatment of
592 intermolecular interactions.

593 We have been careful to document the shortcomings of TAFFI-gen, since our long-term
594 goal is not to simply make the best fixed-charge force-field, but to develop a data-driven
595 means of matching force-field complexity to simulation targets. For instance, amide and
596 halogen containing molecules exhibited among the largest deviations in TAFFI-gen for var-
597 ious liquid properties. Although it would be possible to introduce *ad hoc* corrections to the
598 LJ parameters and partial charges associated with these functional groups, it would come at
599 the expense of increasing errors in reproducing the interaction energies in the training data,
600 and thus would likely lead to uncontrolled errors in other liquid properties. Such *ad hoc*
601 corrections are what we want to avoid with TAFFI. From our perspective, a better pathway
602 forward is to systematically increase the complexity of specific force-field terms based on
603 well-defined error metrics. For example, selectively adding lone-pair sites or Drude particles
604 to specific functional groups could foreseeably be done in a data-driven manner to improve
605 the accuracy of a specific property without introducing *ad hoc* corrections. Likewise, we ob-
606 served that carbonates require larger model compounds than other functional groups, which

607 motivates potentially treating distinct functional groups at variable levels of graph speci-
608 ficity (i.e., in contrast to the fixed $d = 2$ specificity used here for all benchmarks). Within
609 the context of the TAFFI framework, such comparative retraining against shared training
610 data is possible while retaining transferability and on-the-fly extensibility. Additionally, the
611 systematic expansion of training data based on the occurrence of new atom types is also a
612 promising basis for training transferable ML force-fields for organic chemistry.

613 The current study is limited to liquid simulations of neutral non-cyclic organic species,
614 but several extensions to other classes of molecules and force-field forms are obvious and
615 underway. Because TAFFI is based solely on quantum chemistry data, it can be extended
616 to ionic and radical species that have limited coverage in existing experimentally based
617 force-fields. The extension to ions and radicals will require a more general treatment of
618 formal charges in the atom types and model compounds than has been presented here. We
619 have also noted that cyclic molecules and large conjugated groups fundamentally challenge
620 the locality assumption implicit in the use of atom types. A workable near-term solution
621 is to parameterize such systems whole and later use the data generated in this way to
622 establish general ring and conjugation corrections. With respect to extending TAFFI to
623 support the parameterization of more complex force-fields, it will be necessary to augment
624 the calculations currently performed on model compounds to include properties like atomic
625 polarizability, heat of formation, and bond-dissociation energies that would justify more
626 complex parameterizations. The small model compounds used by TAFFI for generating
627 reference data is an advantage in this respect, as higher levels of theory and more extensive
628 characterizations can be afforded while pursuing broad coverage of organic chemical space.

629 **Acknowledgement**

630 Acknowledgment is made to the Donors of the American Chemical Society Petroleum Re-
631 search Fund for support of the work by B.S and Z-Y. L. The work of Q. Z. was made

632 possible through support of the Purdue Process Safety and Assurance Center. M.A.W. ac-
633 knowledges support from Princeton University. The work performed by B.M.S. was made
634 possible through the Dreyfus Program for Machine Learning in the Chemical Sciences and
635 Engineering. This work used the Extreme Science and Engineering Discovery Environment
636 (XSEDE), which is supported by National Science Foundation grant no. ACI-1548562. Sim-
637 ulations were performed on the Comet supercomputer at the University of California, San
638 Diego, under the Allocation no. TG-CHE190014

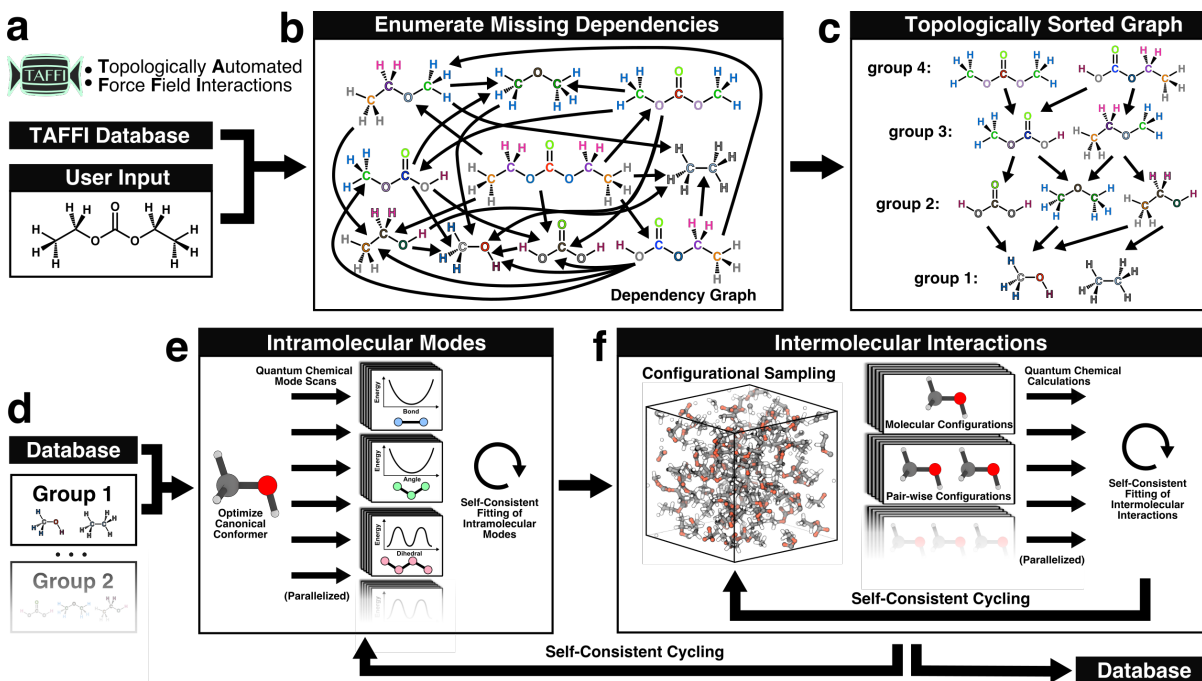


Figure 1: Structure to simulation overview of the TAFFI methodology using diethyl carbonate as an example. (a) Topological criteria are used to determine the necessary parameters for the simulation and identify the missing parameters in the database. (b) An unsorted graph of the molecular dependencies for simulating diethyl carbonate. For simplicity only the dependencies associated with atom types (i.e., not bonds, angles, etc.) are shown, arrows point toward dependencies, and unique atom types at a bond depth of two are distinctly colored. (c) TAFFI model compound rules produce directed acyclic dependency graphs that can always be linearized to sequentially organize calculations. (d) Hierarchical organization ensures that all dependencies exist prior to attempting the parameterization. (e) Intramolecular modes are parameterized using constrained mode scans from quantum chemistry. (f) Intermolecular interactions are parameterized using quantum chemical calculations on molecular configurations sampled from molecular dynamics. The TAFFI database is updated each cycle and all quantum chemistry data is retained for future refitting and force-field extension.

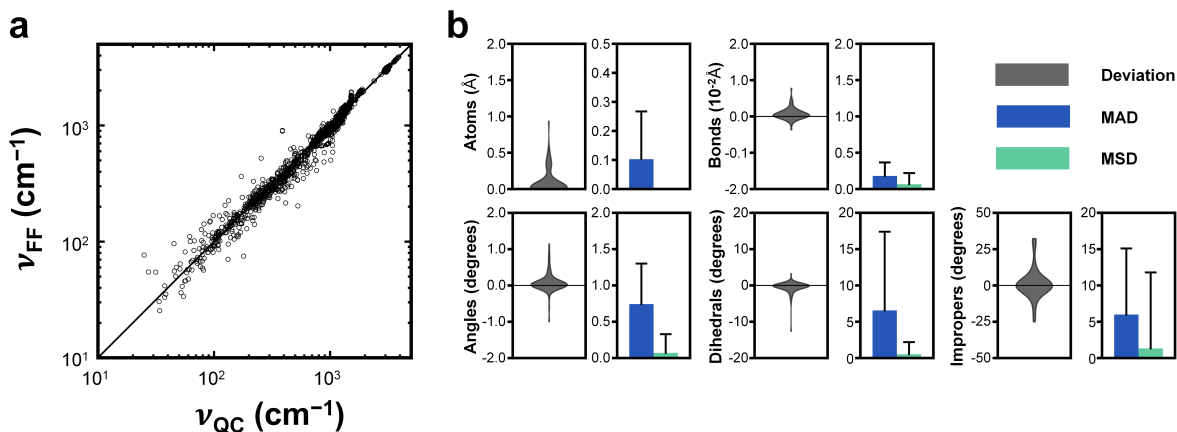


Figure 2: (a) Comparison of the TAFFI-gen and ω B97X-D3/def2-PVTZ (DFT) normal mode frequencies for the benchmark compounds. (b) The distributions of signed deviations ($x_{\text{TAFI}} - x_{\text{DFT}}$) for selected structural features over all benchmarked compounds are shown in each violin plot. The distribution of atom deviations corresponds to the MAD in the atomic positions after alignment of the TAFFI-gen and DFT optimized structures. The other distributions correspond to the signed differences in the bond lengths, bending angles, dihedral angles, and improper angles in the optimized TAFFI-gen geometries and in the optimized DFT geometries. The mean and standard deviation of the mean absolute differences (MAD, blue) and mean signed differences (MSD, green) for each quantity calculated across all benchmark compounds are shown in the bar plots. Improvers are only included for 3-coordinate atoms.

Table 1: Summary of TAFFI-gen performance in reproducing the DFT normal mode frequencies and structural features of the 87 molecules in the benchmark set.

Structure	MAD	MSD	N	Molecules
Normal Modes (cm^{-1})	52.1	-14.7	2908	87
Atoms (\AA)	0.103	^a	1151	87
Bonds (\AA)	0.00181	0.000665	1064	87
Angles (degrees)	0.743	0.0714	1842	87
Dihedrals (degrees)	6.56	-0.520	1919	80
Improvers (degrees)	5.99	1.32	58	40

^aTrivially zero due to structural alignment.

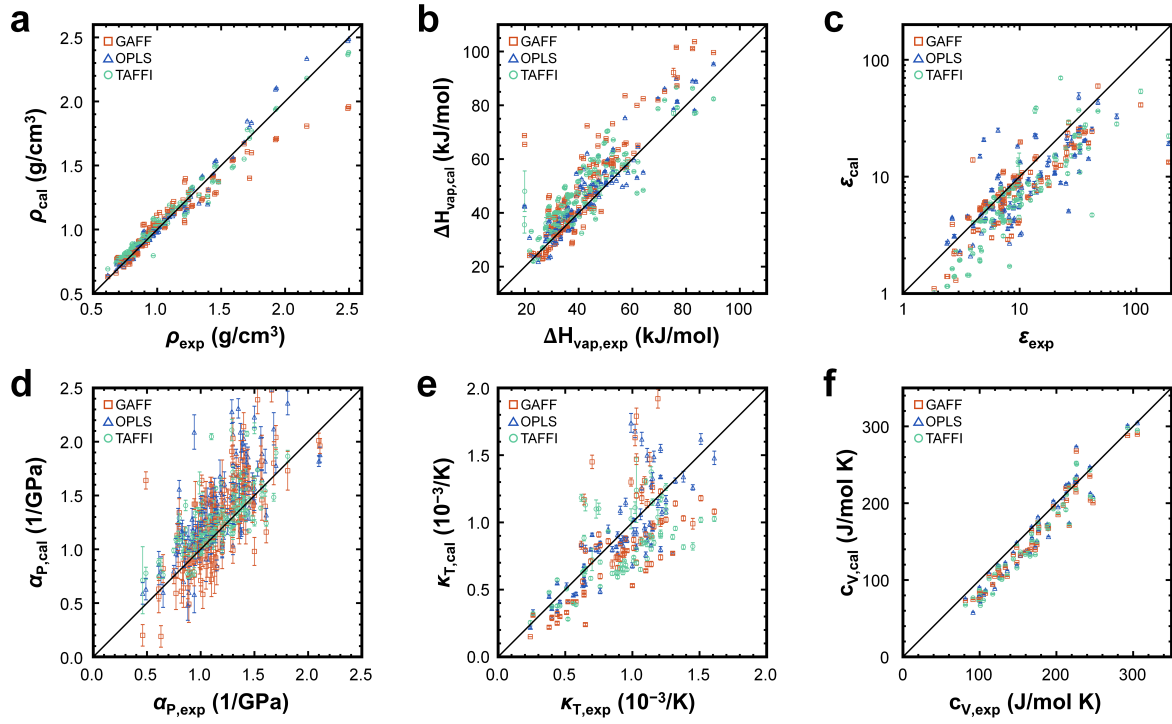


Figure 3: Comparisons of the experimental values for (a) densities, (b) enthalpies of vaporization, (c) static dielectric constants, (d) volumetric thermal expansion coefficients (e) isothermal compressibilities, and (f) quantum-corrected heat capacities at constant volume with those predicted by GAFF (red), OPLS-AA (blue), and TAFFI-gen (green).

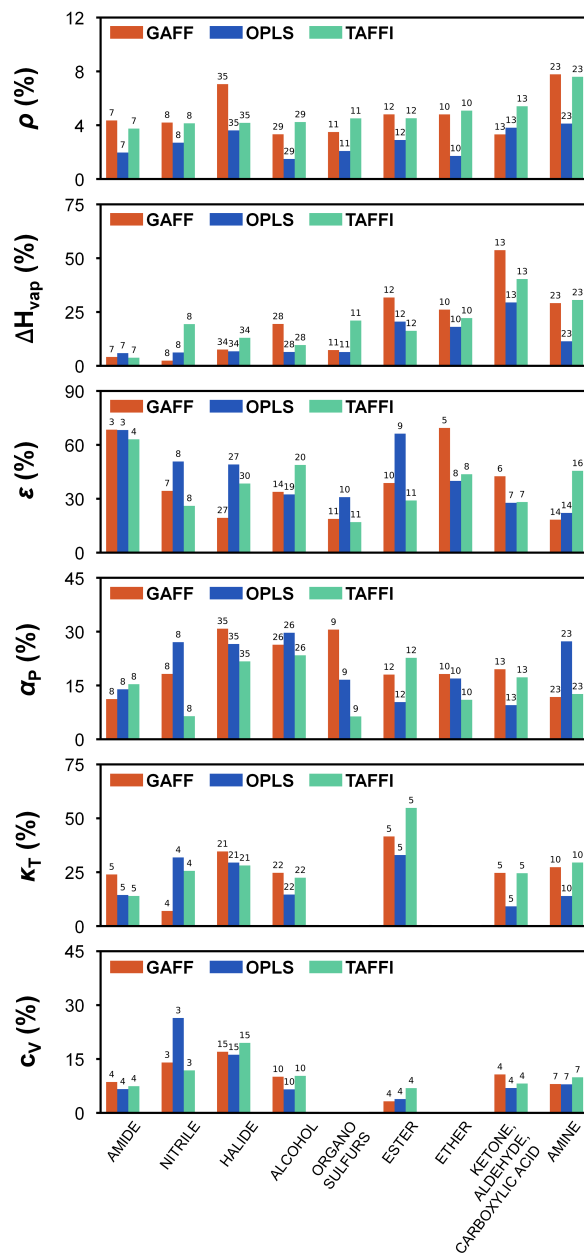


Figure 4: Mean absolute percent difference (MAPD) of the liquid properties for each functional group. The benchmark molecules are classified by the functional groups exhibited by each molecule. Each bar represents the average of the MAPD for all molecules belonging to each group. The numbers of molecules in each case are indicated above the bars and properties with less than three values have been omitted. GAFF (red) and OPLS-AA (blue) data are from reference⁶⁵ whereas TAFFI-gen (green) data are from MD simulations performed in the current study.

Table 2: Comparison of the errors in the liquid properties for the GAFF, OPLS-AA, and TAFFI-gen force-fields. The mean absolute difference (MAD), mean signed difference (MSD), mean absolute percent difference (MAPD), the mean signed percent difference (MSPD), the root mean square deviation (RMSD) from experimental values, and the correlation coefficient R^2 are reported.

Force-field	MAD ^a	MSD ^a	MAPD ^b	MSPD ^b	RMSD ^a	R^2 ^b	N
ρ (g/cm ³)							
GAFF	0.0590	-0.0060	5.0970	0.8421	1.00	94.17	145
OPLS-AA	0.0311	0.0114	2.9424	1.2046	0.48	98.24	145
TAFFI-gen	0.0484	0.0231	5.0971	3.2570	0.58	97.94	145
ΔH_{vap} (kJ/mol)							
GAFF	7.7691	6.4625	19.7032	16.0226	11.10	78.69	143
OPLS-AA	4.3424	2.9003	11.2727	7.7738	6.18	87.17	143
TAFFI-gen	7.3489	5.9987	19.5204	16.9972	8.89	78.62	143
ε							
GAFF	6.1100	-4.9042	30.1701	-13.9654	19.90	29.84	97
OPLS-AA	6.9686	-4.7846	40.7308	-9.5976	18.67	25.60	103
TAFFI-gen	7.2708	-5.2487	37.8468	-25.1088	19.03	30.39	113
α_P (10 ⁻³ /K)							
GAFF	0.2411	0.1124	21.9688	9.5985	0.34	50.00	140
OPLS-AA	0.2528	0.1906	22.2424	16.9217	0.33	54.80	140
TAFFI-gen	0.1821	0.1308	16.5202	12.5512	0.27	58.43	140
κ_T (1/GPa)							
GAFF	0.2475	-0.0577	27.6643	-6.8676	0.31	43.49	73
OPLS-AA	0.1875	0.0273	20.3002	2.8656	0.29	52.02	73
TAFFI-gen	0.2584	-0.0811	27.5593	-5.7311	0.38	22.18	73
c_V (J/mol K)							
GAFF	17.7962	-15.4722	11.5785	-10.5375	21.01	93.89	50
OPLS-AA	16.5314	-12.0042	11.0421	-8.9901	20.48	93.51	50
TAFFI-gen	18.4626	-15.7177	12.3048	-11.1073	21.68	94.21	50

^aIn indicated units

^bIn units of %

References

- [1] Jorgensen, W. L. & Tirado-Rives, J. Potential energy functions for atomic-level simulations of water and organic and biomolecular systems. *Proc. Natl. Acad. Sci. USA* **102**, 6665–6670 (2005).
- [2] Huang, J. & MacKerell Jr, A. D. Force field development and simulations of intrinsically disordered proteins. *Curr. Opin. Struct. Biol.* **48**, 40–48 (2018).
- [3] Nerenberg, P. S. & Head-Gordon, T. New developments in force fields for biomolecular simulations. *Curr. Opin. Struct. Biol.* **49**, 129–138 (2018).
- [4] Lemkul, J. A., Huang, J., Roux, B. & MacKerell Jr, A. D. An empirical polarizable force field based on the classical drude oscillator model: development history and recent applications. *Chem. Rev.* **116**, 4983–5013 (2016).
- [5] Liang, T. *et al.* Reactive potentials for advanced atomistic simulations. *Annu. Rev. Mater. Res.* **43**, 109–129 (2013).
- [6] Xu, P., Guidez, E. B., Bertoni, C. & Gordon, M. S. Perspective: Ab initio force field methods derived from quantum mechanics. *J. Chem. Phys.* **148**, 090901 (2018).
- [7] Riniker, S. Fixed-charge atomistic force fields for molecular dynamics simulations in the condensed phase: An overview. *J. Chem. Inf. Model.* **58**, 565–578 (2018).
- [8] Cornell, W. *et al.* A second generation force field for the simulation of proteins, nucleic acids, and organic molecules. *J. Am. Chem. Soc.* **117**, 5179–5197 (1995).
- [9] Debiec, K. *et al.* Further along the road less traveled: Amber ff15ipq, an original protein force field built on a self-consistent physical model. *J. Chem. Theory Comput.* **12**, 3926–3947 (2016).
- [10] MacKerell Jr, A. *et al.* All-atom empirical potential for molecular modeling and dynamics studies of proteins. *J. Phys. Chem. B* **102**, 3586–3616 (1998).

- 663 [11] Huang, J. *et al.* Charmm36m: an improved force field for folded and intrinsically
664 disordered proteins. *Nat. Methods.* **14**, 71–73 (2016).
- 665 [12] Daura, X., Mark, A. & van Gunsteren, W. Parametrization of aliphatic chn united
666 atoms of gromos96 force field. *J. Comput. Chem.* **19**, 535–547 (1998).
- 667 [13] Horta, B. *et al.* A gromos-compatible force field for small organic molecules in the
668 condensed phase: The 2016h66 parameter set. *J. Chem. Theory Comput.* **12**, 3825–
669 3850 (2016).
- 670 [14] Jorgensen, W. L. & Tirado-Rives, J. The opls [optimized potentials for liquid sim-
671 ulations] potential functions for proteins, energy minimizations for crystals of cyclic
672 peptides and crambin. *J. Am. Chem. Soc.* **110**, 1657–1666 (1988).
- 673 [15] Harder, E. *et al.* Opls3: A force field providing broad coverage of drug-like small
674 molecules and proteins. *J. Chem. Theory Comput.* **12**, 281–296 (2016).
- 675 [16] Behler, J. & Parrinello, M. Generalized neural-network representation of high-
676 dimensional potential-energy surfaces. *Phys. Rev. Lett.* **98**, 583–4 (2007).
- 677 [17] Bartók, A. P., Payne, M. C., Kondor, R. & Csányi, G. Gaussian approximation poten-
678 tials: The accuracy of quantum mechanics, without the electrons. *Phys. Rev. Lett.* **104**,
679 136403 (2010).
- 680 [18] Artrith, N. & Urban, A. An implementation of artificial neural-network potentials
681 for atomistic materials simulations: Performance for tio2. *Comput. Mater. Sci.* **114**,
682 135–150 (2016).
- 683 [19] Khorshidi, A. & Peterson, A. A. Amp: A modular approach to machine learning in
684 atomistic simulations. *Comput. Phys. Commun.* **207**, 310–324 (2016).
- 685 [20] Smith, J. S., Isayev, O. & Roitberg, A. E. Ani-1: an extensible neural network potential
686 with dft accuracy at force field computational cost. *Chem. Sci.* **8**, 3192–3203 (2017).

- 687 [21] Zhang, L., Han, J., Wang, H., Car, R. & Weinan, E. Deep potential molecular dynamics:
688 a scalable model with the accuracy of quantum mechanics. *Phys. Rev. Lett.* **120**, 143001
689 (2018).
- 690 [22] Yao, K., Herr, J. E., Toth, D. W., Mckintyre, R. & Parkhill, J. The tensormol-0.1
691 model chemistry: a neural network augmented with long-range physics. *Chem. Sci.* **9**,
692 2261–2269 (2018).
- 693 [23] Schütt, K. T., Arbabzadah, F., Chmiela, S., Müller, K. R. & Tkatchenko, A. Quantum-
694 chemical insights from deep tensor neural networks. *Nat. Commun.* **8**, 1–8 (2017).
- 695 [24] Schütt, K. T., Sauceda, H. E., Kindermans, P.-J., Tkatchenko, A. & Müller, K.-R.
696 Schnet—a deep learning architecture for molecules and materials. *J. Chem. Phys.* **148**,
697 241722 (2018).
- 698 [25] Chmiela, S., Sauceda, H. E., Müller, K.-R. & Tkatchenko, A. Towards exact molecular
699 dynamics simulations with machine-learned force fields. *Nat. Commun.* **9**, 1–10 (2018).
- 700 [26] Anisimov, V. *et al.* Determination of electrostatic parameters for a polarizable force field
701 based on the classical drude oscillator. *J. Chem. Theory Comput.* **1**, 153–168 (2005).
- 702 [27] Lemkul, J., Huang, J., Roux, B. & Mackerell, A. An empirical polarizable force field
703 based on the classical drude oscillator model: Development history and recent applica-
704 tions. *Chem. Rev.* **116**, 4983–5013 (2016).
- 705 [28] McDaniel, J. & Schmidt, J. Physically-motivated force fields from symmetry-adapted
706 perturbation theory. *J. Phys. Chem. A* **117**, 2053–2066 (2015).
- 707 [29] McDaniel, J., Choi, E., Son, C., Schmidt, J. & Yethiraj, A. Conformational and dynamic
708 properties of poly (ethylene oxide) in an ionic liquid: Development and implementation
709 of a first-principles force field. *J. Phys. Chem. B* **120**, 231–243 (2016).

- 710 [30] Duan, Y. *et al.* A point-charge force field for molecular mechanics simulations of proteins
711 based on condensed-phase quantum mechanical calculations. *J. Comput. Chem.* **24**,
712 1999–2012 (2003).
- 713 [31] Hornak, V. *et al.* Comparison of multiple amber force fields and development of im-
714 proved protein backbone parameters. *Proteins* **65**, 712–725 (2006).
- 715 [32] Lindorff-Larsen, K. *et al.* Improved side-chain torsion potentials for the amber ff99sb
716 protein force field. *Proteins* **78**, 1950–1958 (2010).
- 717 [33] Maier, J. A. *et al.* ff14sb: improving the accuracy of protein side chain and backbone
718 parameters from ff99sb. *J. Chem. Theory Comput.* **11**, 3696–3713 (2015).
- 719 [34] Wang, J. & Hou, T. Application of molecular dynamics simulations in molecular prop-
720 erty prediction. 1. density and heat of vaporization. *J. Chem. Theory Comput.* **7**,
721 2151–2165 (2011).
- 722 [35] Kaminski, G. A., Friesner, R. A., Tirado-Rives, J. & Jorgensen, W. L. Evaluation and
723 reparametrization of the opls-aa force field for proteins via comparison with accurate
724 quantum chemical calculations on peptides. *J. Phys. Chem. B* **105**, 6474–6487 (2001).
- 725 [36] Mackerell Jr, A. D., Feig, M. & Brooks III, C. L. Extending the treatment of back-
726 bone energetics in protein force fields: Limitations of gas-phase quantum mechanics in
727 reproducing protein conformational distributions in molecular dynamics simulations. *J.*
728 *Comput. Chem.* **25**, 1400–1415 (2004).
- 729 [37] Best, R. B. *et al.* Optimization of the additive charmm all-atom protein force field
730 targeting improved sampling of the backbone ϕ , ψ and side-chain χ_1 and χ_2 dihedral
731 angles. *J. Chem. Theory Comput.* **8**, 3257–3273 (2012).
- 732 [38] Cacelli, I. & Prampolini, G. Parametrization and validation of intramolecular force
733 fields derived from dft calculations. *J. Chem. Theory Comput.* **3**, 1803–1817 (2007).

- 734 [39] Horton, J. T., Allen, A. E., Dodda, L. S. & Cole, D. J. Qubekit: automating the
735 derivation of force field parameters from quantum mechanics. *J. Chem. Inf. Model.* **59**,
736 1366–1381 (2019).
- 737 [40] Grimme, S. A general quantum mechanically derived force field (qmdff) for molecules
738 and condensed phase simulations. *J. Chem. Theory Comput.* **10**, 4497–4514 (2014).
- 739 [41] Vanommeslaeghe, K., Raman, E. & MacKerell, A. Automation of the charmm general
740 force field (cgenff) ii: Assignment of bonded parameters and partial atomic charges. *J.*
741 *Chem. Inf. Model* **52**, 3155–3168 (2012).
- 742 [42] Shivakumar, D., Harder, E., Damm, W., Friesner, R. & Sherman, W. Improving the
743 prediction of absolute solvation free energies using the next generation opl force field.
744 *J. Chem. Theory Comput.* **8**, 2553–2558 (2012).
- 745 [43] Boyd, N. & Wilson, M. Optimization of the gaff force field to describe liquid crystal
746 molecules: the path to a dramatic improvement in transition temperature predictions.
747 *Phys. Chem. Chem. Phys.* **17**, 24851–24865 (2015).
- 748 [44] Doherty, B., Zhong, X., Gathiaka, S., Li, B. & Acevedo, O. Revisiting opl force field
749 parameters for ionic liquid simulations. *J. Chem. Theory Comput.* **13**, 6131–6145 (2017).
- 750 [45] Jin, Z. *et al.* Hierarchical atom type definitions and extensible all-atom force fields. *J.*
751 *Comput. Chem.* **37**, 653–664 (2016).
- 752 [46] Mobley, D. L. *et al.* Escaping atom types in force fields using direct chemical perception.
753 *J. Chem. Theory Comput.* **14**, 6076–6092 (2018).
- 754 [47] Wang, J., Wolf, R., Caldwell, J., Kollman, P. & Case, D. Development and testing of a
755 general amber force field. *J. Comput. Chem.* **25**, 1157–1174 (2004).
- 756 [48] Case, D. *et al.* Amber 2016. *University of California, San Francisco* (2016).

- 757 [49] Vanommeslaeghe, K. *et al.* Charmm general force field: A force field for drug-like
758 molecules compatible with the charmm all-atom additive biological force fields. *J. Comput. Chem.* **31**, 671–690 (2010).
759
- 760 [50] Mayne, C., Saam, J., Schulten, K., Tajkhorshid, E. & Gumbart, J. Rapid parameteri-
761 zation of small molecules using the force field toolkit. *J. Comput. Chem.* **34**, 2757–2770
762 (2013).
- 763 [51] Jorgensen, W., Maxwell, D. & Tirado-Rives, J. Development and testing of the opls
764 all-atom force field on conformational energetics and properties of organic liquids. *J. Am. Chem. Soc.* **118**, 11225–11236 (1996).
765
- 766 [52] Roos, K. *et al.* Opls3e: Extending force field coverage for drug-like small molecules. *J. Chem. Theory Comput.* **15**, 1863–1874 (2019).
767
- 768 [53] Savoie, B. M., Webb, M. A. & Miller III, T. F. Enhancing cation diffusion and sup-
769 pressing anion diffusion via lewis-acidic polymer electrolytes. *J. Phys. Chem. Lett.* **8**,
770 641–646 (2017).
- 771 [54] Khot, A., Shiring, S. B. & Savoie, B. M. Evidence of information limitations in coarse-
772 grained models. *J. Chem. Phys.* **151**, 244105 (2019).
- 773 [55] Zhao, Q. & Savoie, B. M. Self-consistent component increment theory for predicting
774 enthalpy of formation. *J. Chem. Inf. Model.* **60**, 2199–2207 (2020).
- 775 [56] Sanderson, R. T. Electronegativity and bond energy. *J. Am. Chem. Soc.* **105**, 2259–2261
776 (1983).
- 777 [57] Sanderson, R. *Chemical bonds and bonds energy*, vol. 21 (Elsevier, 2012).
- 778 [58] Rappe, A., Casewit, C., Colwell, K., Goddard, W. & Skiff, W. Uff, a full periodic table
779 force field for molecular mechanics and molecular dynamics simulations. *J. Am. Chem. Soc.* **114**, 10024–10035 (1992).
780

- 781 [59] Bayly, C. I., Cieplak, P., Cornell, W. & Kollman, P. A. A well-behaved electrostatic
782 potential based method using charge restraints for deriving atomic charges: the RESP
783 model. *J. Phys. Chem.* **97**, 10269–10280 (1993).
- 784 [60] Plimpton, S. Fast parallel algorithms for short-range molecular dynamics. *J. Comput.*
785 *Phys.* **117**, 1–19 (1995).
- 786 [61] Neese, F. The orca program system. *WIREs Comput Mol Sci* **2**, 73–78 (2012).
- 787 [62] Lin, Y.-S., Li, G.-D., Mao, S.-P. & Chai, J.-D. Long-range corrected hybrid density
788 functionals with improved dispersion corrections. *J. Chem. Theory Comput.* **9**, 263–272
789 (2013).
- 790 [63] Schäfer, A., Horn, H. & Ahlrichs, R. Fully optimized contracted gaussian basis sets for
791 atoms li to kr. *J. Chem. Phys.* **97**, 2571–2577 (1992).
- 792 [64] Weigend, F. & Ahlrichs, R. Balanced basis sets of split valence, triple zeta valence and
793 quadruple zeta valence quality for h to rn: Design and assessment of accuracy. *Phys.*
794 *Chem. Chem. Phys.* **7**, 3297–3305 (2005).
- 795 [65] Caleman, C. *et al.* Force field benchmark of organic liquids: Density, enthalpy of
796 vaporization, heat capacities, surface tension, isothermal compressibility, volumetric
797 expansion coefficient, and dielectric constant. *J. Chem. Theory Comput.* **8**, 61–74 (2012).
- 798 [66] Berens, P. H., Mackay, D. H. J., White, G. M. & Wilson, K. R. Thermodynamics and
799 quantum corrections from molecular dynamics for liquid water. *J. Chem. Phys.* **79**,
800 2375–2389 (1983).
- 801 [67] Pascal, T. A., Lin, S.-T. & Goddard III, W. A. Thermodynamics of liquids: standard
802 molar entropies and heat capacities of common solvents from 2pt molecular dynamics.
803 *Phys. Chem. Chem. Phys.* **13**, 169–181 (2011).

- 804 [68] Chen, M. *et al.* Ab initio theory and modeling of water. *Proc. Natl. Acad. Sci. USA*
805 **114**, 10846–10851 (2017).
- 806 [69] Yao, Y. & Kanai, Y. Free energy profile of nacl in water: first-principles molecular
807 dynamics with scan and ω b97x-v exchange–correlation functionals. *J. Chem. Theory*
808 *Comput.* **14**, 884–893 (2018).
- 809 [70] Seeger, Z. L. & Izgorodina, E. I. A systematic study of dft performance for geometry
810 optimizations of ionic liquid clusters. *J. Chem. Theory Comput.* **16**, 6735–6753 (2020).
- 811 [71] Sure, R. & Grimme, S. Comprehensive benchmark of association (free) energies of
812 realistic host–guest complexes. *J. Chem. Theory Comput.* **11**, 3785–3801 (2015).
- 813 [72] Lao, K. U., Schäffer, R., Jansen, G. & Herbert, J. M. Accurate description of inter-
814 molecular interactions involving ions using symmetry-adapted perturbation theory. *J.*
815 *Chem. Theory Comput.* **11**, 2473–2486 (2015).
- 816 [73] Mobley, D. L. *et al.* Open force field consortium: Escaping atom types using direct
817 chemical perception with smirnoff v0. 1. *BioRxiv* 286542 (2018).
- 818 [74] Fennell, C. J., Wymer, K. L. & Mobley, D. L. A fixed-charge model for alcohol po-
819 larization in the condensed phase, and its role in small molecule hydration. *J. Phys.*
820 *Chem. B* **118**, 6438–6446 (2014).
- 821 [75] Mobley, D. L., Bayly, C. I., Cooper, M. D., Shirts, M. R. & Dill, K. A. Small molecule
822 hydration free energies in explicit solvent: an extensive test of fixed-charge atomistic
823 simulations. *J. Chem. Theory Comput.* **5**, 350–358 (2009).
- 824 [76] Sun, H. *et al.* Compass ii: extended coverage for polymer and drug-like molecule
825 databases. *J. Mol. Model.* **22**, 47 (2016).
- 826 [77] Kramer, C., Spinn, A. & Liedl, K. R. Charge anisotropy: where atomic multipoles
827 matter most. *J. Chem. Theory Comput.* **10**, 4488–4496 (2014).

- 828 [78] Harder, E., Anisimov, V. M., Whitfield, T., MacKerell, A. D. & Roux, B. Understanding
829 the dielectric properties of liquid amides from a polarizable force field. *J. Phys. Chem.*
830 *B* **112**, 3509–3521 (2008).
- 831 [79] Murray, J. S., Lane, P., Clark, T. & Politzer, P. σ -hole bonding: molecules containing
832 group vi atoms. *J. Mol. Model.* **13**, 1033–1038 (2007).
- 833 [80] Clark, T., Hennemann, M., Murray, J. S. & Politzer, P. Halogen bonding: the σ -hole.
834 *J. Mol. Model.* **13**, 291–296 (2007).
- 835 [81] Ibrahim, M. A. Molecular mechanical study of halogen bonding in drug discovery. *J.*
836 *Comput. Chem.* **32**, 2564–2574 (2011).
- 837 [82] Rendine, S., Pieraccini, S., Forni, A. & Sironi, M. Halogen bonding in ligand–receptor
838 systems in the framework of classical force fields. *Phys. Chem. Chem. Phys.* **13**, 19508–
839 19516 (2011).
- 840 [83] Kolář, M. & Hobza, P. On extension of the current biomolecular empirical force field
841 for the description of halogen bonds. *J. Chem. Theory Comput.* **8**, 1325–1333 (2012).
- 842 [84] Jorgensen, W. L. & Schyman, P. Treatment of halogen bonding in the opls-aa force
843 field: application to potent anti-hiv agents. *J. Chem. Theory Comput.* **8**, 3895–3901
844 (2012).
- 845 [85] Gutiérrez, I. S. *et al.* Parametrization of halogen bonds in the charmm general force field:
846 Improved treatment of ligand–protein interactions. *Bioorg. Med. Chem.* **24**, 4812–4825
847 (2016).
- 848 [86] Bereau, T., Kramer, C. & Meuwly, M. Leveraging symmetries of static atomic multipole
849 electrostatics in molecular dynamics simulations. *J. Chem. Theory Comput.* **9**, 5450–
850 5459 (2013).

- 851 [87] Mu, X. *et al.* Modeling organochlorine compounds and the σ -hole effect using a polar-
852 izable multipole force field. *J. Phys. Chem. B* **118**, 6456–6465 (2014).
- 853 [88] Du, L., Gao, J., Bi, F., Wang, L. & Liu, C. A polarizable ellipsoidal force field for
854 halogen bonds. *J. Comput. Chem.* **34**, 2032–2040 (2013).
- 855 [89] Lin, F.-Y. & MacKerell Jr, A. D. Polarizable empirical force field for halogen-containing
856 compounds based on the classical drude oscillator. *J. Chem. Theory Comput.* **14**, 1083–
857 1098 (2018).
- 858 [90] Carter, M., Rappé, A. K. & Ho, P. S. Scalable anisotropic shape and electrostatic
859 models for biological bromine halogen bonds. *J. Chem. Theory Comput.* **8**, 2461–2473
860 (2012).
- 861 [91] Jedlovszky, P. & Turi, L. A new five-site pair potential for formic acid in liquid simu-
862 lations. *J. Phys. Chem. A* **101**, 2662–2665 (1997).
- 863 [92] Qian, W. & Krimm, S. Electrostatic model for the interaction force constants of the
864 formic acid dimer. *J. Phys. Chem. A* **102**, 659–667 (1998).
- 865 [93] Ramón, J. M. H. & Rios, M. A. A new intermolecular polarizable potential for cis-formic
866 acid. introduction of many-body interactions in condensed phases. *Chem. Phys.* **250**,
867 155–169 (1999).
- 868 [94] Roszak, S., Gee, R. H., Balasubramanian, K. & Fried, L. E. New theoretical insight
869 into the interactions and properties of formic acid: Development of a quantum-based
870 pair potential for formic acid. *J. Chem. Phys.* **123**, 144702 (2005).
- 871 [95] Schnabel, T., Cortada, M., Vrabec, J., Lago, S. & Hasse, H. Molecular model for formic
872 acid adjusted to vapor–liquid equilibria. *Chem. Phys. Lett.* **435**, 268–272 (2007).