

1 **Application of machine-learning algorithms to predict the transport properties of**
2 **Mie fluids**

3 Justinas Šlepavičius,¹ Alessandro Patti,^{1,2} James L. McDonagh,³ and Carlos
4 Avendaño^{1, a)}

5 ¹⁾*Department of Chemical Engineering, School of Engineering,*
6 *The University of Manchester, Oxford Road, Manchester, M13 9PL,*
7 *United Kingdom*

8 ²⁾*Department of Applied Physics, University of Granada, Fuente Nueva s/n,*
9 *18071 Granada, Spain*

10 ³⁾*IBM Research Europe, The Hartree Centre STFC Laboratory Sci-Tech Daresbury,*
11 *Warrington, United Kingdom ^{b)}*

12 (Dated: 9 June 2023)

The ability to predict transport properties of fluids, such as the self-diffusion coefficient and viscosity, has been an ongoing effort in the field of molecular modelling. While there are theoretical approaches to predict the transport properties of simple systems, they are typically applied in the dilute gas regime and are not directly applicable to more complex systems. Other attempts to predict transport properties are done by fitting available experimental or molecular simulation data to empirical or semi-empirical correlations. Recently, there have been attempts to improve the accuracy of these fittings through the use of Machine Learning (ML) methods. In this work, the application of ML algorithms to represent the transport properties of systems comprising spherical particles interacting *via* the Mie potential is investigated. To this end, the self-diffusion coefficient and shear viscosity of 54 potentials are obtained at different regions of the fluid-phase diagram. This data set is used together with three ML algorithms, namely *k*-Nearest Neighbours, Artificial Neural Network and Symbolic Regression, to find correlations between the parameters of each potential and the transport properties at different densities and temperatures. It is shown that ANN and KNN perform to a similar extent, followed by SR, which exhibits larger deviations. Finally, the application of the three ML models to predict the self-diffusion coefficient of small molecular systems, such as krypton, methane and carbon dioxide is demonstrated using molecular parameters derived from the so-called SAFT-VR Mie equation of state [J. Chem. Phys. **139**, 154504 (2013)] and available experimental vapour-liquid coexistence data.

^{a)}Electronic mail: carlos.avendano@manchester.ac.uk

^{b)}Current address: Ladder Therapeutics doing business as Serna Bio, Lab F37, Stevenage Bioscience Catalyst, Gunnels Wood Road, Stevenage, Hertfordshire, SG1 2FX, United Kingdom

13 I. INTRODUCTION

14 Accurate representation of thermodynamic and transport properties of molecular systems is
15 key for the design and optimization of chemical and biochemical processes and plays an important
16 role in several areas of science and technology^{1,2}. For many decades, most attention has been
17 drawn to the development of robust thermodynamic models of fluids based on activity coefficient
18 models³, cubic equations of state^{4,5}, and molecular-based models⁶⁻⁹. The development of models
19 for the description of transport properties has also been reported in recent years, particularly for
20 high-density liquids. Many of these approaches are based on kinetic theory, mode-coupling theory,
21 free-volume theory, and friction theory¹⁰⁻¹³.

22 Molecular dynamics is a powerful technique to determine the transport properties of molecular
23 systems using either Green-Kubo or Einstein relations^{14,15}, provided a suitable force field is avail-
24 able. The simulations can be used to test theoretical predictions of transport properties. However,
25 molecular simulations can be computationally expensive and not suitable for fast property predic-
26 tion required in many industrial applications. For this reason, simple empirical and semi-empirical
27 correlations are still commonly used for many engineering calculations, but these equations have
28 a limited range of applicability making them unreliable for the design of novel technologies.

29 There have been previous attempts to obtain theoretical or semi-empirical expressions for the
30 self-diffusion coefficients of simple systems, such as hard spheres, Lennard-Jones, WCA, and
31 molecular and atomic species¹⁶⁻²⁷. Still, large amounts of data are required for fitting the expres-
32 sions that typically have complicated functional forms with multiple fitting parameters. A similar
33 situation has been observed in the developments of semi-empirical equations for the viscosity²⁸⁻³¹,
34 however, new advances in the use of equations of state coupled with methods such as the free-
35 volume, kinetic, and friction theories to predict shear viscosity of complex molecular systems
36 have been reported, particularly for the prediction of the shear viscosity³²⁻³⁴.

37 To facilitate the way the predictions about the behaviour of transport properties are made, at-
38 tention has been turning to machine learning (ML) methods. These methods are able to take the
39 large amounts of data available and explore possible correlations between the system parameters
40 and properties of interest. This change towards ML methods has already been taking place in
41 similar fields, such as in the prediction of physicochemical properties³⁵⁻⁵⁴. For a recent account
42 on this topic, the reader is directed to the recent reviews by Schmidt *et al.*⁵⁵, Moud *et al.*⁵⁶ and
43 Ihme *et al.*⁵⁷ The use of ML to predict transport properties, such as the self-diffusion coefficient,

44 has been studied to a much lesser extent⁵⁸. ML algorithms to study transport properties have been
45 used to predict diffusion in binary hydrocarbons⁵⁹, water mixtures⁶⁰, polar and non-polar binary
46 gases⁶¹, organic molecules^{62,63}, and CO₂⁶⁴. Recently, the application of ML algorithms to study
47 the self-diffusion coefficient of Lennard-Jones systems has been reported^{65–68}, where the accuracy
48 of the predictions was increased by more than one order of magnitude compared to semi-empirical
49 equations.

50 In this work, ML methods are applied to predict the transport properties, particularly self-
51 diffusion coefficient and viscosity, of spherical particles interacting *via* the Mie potential⁶⁹. The
52 Mie potential is a generalized form of the Lennard-Jones potential but with the advantage that
53 the repulsive and attractive contributions can be tuned to represent accurately the properties of
54 real molecular systems. This potential has been used to study successfully the behaviour of small
55 molecular systems (H₂O, CO₂, SF₆, and CF₄)^{70–73}, and is also the underlying intermolecular
56 potential used in several force fields and molecular-based theories for the description of properties
57 of complex molecular systems^{9,74,75}. The different combinations of the repulsive and attractive
58 contributions impart a large variability of the potential, which also means that there is not a large
59 amount of published data, especially of transport properties, that can be used in this approach
60 using ML methods. Only a small amount of published research in the literature report data on
61 many different types of Mie potentials since most reports concentrate on modelling real particles
62 using a specific combination of potential parameters, hence only a fraction of the potential space
63 is explored and published in the literature. Therefore, in this work, the data required to fit the ML
64 models are also determined.

65 This paper is organised as follows. In Sec. II, the simulation methods employed to obtain the
66 self-diffusion coefficient and shear viscosity data required for the ML algorithms are described.
67 In the same section, the different ML algorithms applied in this work are discussed. In Sec. III,
68 the performance of the ML algorithms in each of the cases studied is reported, and the appli-
69 cation of the methods to describe the self-diffusion coefficient of three quasi-spherical systems,
70 namely krypton Kr, methane CH₄, and carbon dioxide CO₂ is discussed. Finally, in Sec. IV, the
71 conclusions of the work are presented.

72 II. METHODS

73 A. Molecular-dynamics simulations

74 The molecular model consists of spherical particles interacting via the Mie potential given by⁶⁹

$$\phi(r) = \mathcal{C}\varepsilon \left[\left(\frac{\sigma}{r} \right)^n - \left(\frac{\sigma}{r} \right)^m \right], \quad (1)$$

75 where r is the interparticle distance, ε is the potential well-depth, σ is the diameter of the particles,
76 and the exponents n and m describe the range of the repulsive and attractive contributions of the
77 potential, respectively. The coefficient \mathcal{C} in Equation 1 is given by

$$\mathcal{C} \equiv \left(\frac{n}{n-m} \right) \left(\frac{n}{m} \right)^{\frac{m}{n-m}}, \quad (2)$$

78 and is defined such that the minimum of the potential is set at $-\varepsilon$. The Mie(n, m) potential is a
79 generalised form of the well-known Lennard-Jones (LJ) potential, which is obtained when $n = 12$
80 and $m = 6$. It has been demonstrated that the use of different repulsive and attractive exponents
81 results in a much better description of the intermolecular interactions of complex systems^{9,70,71,75}.
82 While there is freedom in selecting arbitrary values of (n, m) to represent a particular molecular
83 system, it has been shown that there are different combinations of exponents that can lead to the
84 same critical point, albeit different triple points⁷⁶. In particular, Ramrattan *et al.* showed that
85 combinations of exponents with the same cohesive parameter α , which is defined as

$$\alpha \equiv \mathcal{C} \left(\frac{1}{m-3} - \frac{1}{n-3} \right), \quad (3)$$

86 exhibit identical critical points⁷⁶. In other words, fluids with the same cohesive parameter are
87 conformal. Due to this observation, it is hypothesised that pair of exponents (n, m) with the same
88 cohesive parameter α should also exhibit the same transport properties in the fluid region akin to
89 the principle of corresponding states for transport properties^{77,78}. To determine the critical point
90 and fluid phase diagram associated with a particular combination of exponents (n, m) , the SAFT-
91 VR Mie equation of state (EoS) is employed⁹. This EoS has been shown to represent the critical
92 region accurately due to the high-order terms considered in the Barker and Henderson perturbation
93 theory⁷⁹.

94 Throughout this work, reduced units are employed to describe thermodynamic and trans-
95 port properties using the Mie potential parameters: number density $\rho^* = N\sigma^3/V$, tempera-
96 ture $T^* = k_B T/\varepsilon$, pressure $p^* = p\sigma^3/\varepsilon$, time $t^* = t[\varepsilon/(m\sigma^2)]^{1/2}$, self-diffusion coefficient

97 $D^* = D[m/(\sigma^2\varepsilon)]^{1/2}$, and viscosity $\eta^* = \eta\sigma^2/(\varepsilon m)^{1/2}$, where N is the total number of parti-
 98 cles, V is the volume of the system, T is the absolute temperature, k_B is the Boltzmann constant, p
 99 is the absolute pressure, t is the time, m is the mass of a spherical particle, D is the self-diffusion
 100 coefficient, and η is the shear viscosity. Similarly, all distances are given in units of σ . To deter-
 101 mine the transport properties of the Mie potential, molecular-dynamics simulations are performed
 102 in systems comprising $N = 10^4$ particles in the canonical NVT ensemble using the Nosé-Hoover
 103 thermostat⁸⁰. For the simulations, the Mie potential is truncated and shifted to zero using a cut-off
 104 of $r_c^* = 6$. All the simulations are performed using the LAMMPS package⁸¹. The equations of
 105 motion are integrated using the velocity-Verlet algorithm with a time step of $\Delta t = 0.001\tau$. For
 106 each state point, five independent simulations of 10^7 time steps are performed to collect averages.
 107 The self-diffusion coefficient has been calculated using the Einstein equation given by

$$D^* = \frac{1}{6t^*} \langle \Delta r^{*2}(t^*) \rangle, \quad (4)$$

108 where $\langle \Delta r^{*2}(t^*) \rangle$ is the ensemble average of the mean-square displacement (MSD) given by

$$\langle \Delta r^{*2}(t^*) \rangle = \frac{1}{N} \left\langle \sum_{j=1}^N [\mathbf{r}_j^*(t^*) - \mathbf{r}_j^*(0)]^2 \right\rangle, \quad (5)$$

109 and $\mathbf{r}_j^*(t^*)$ is the position of particle j at time t^* . The trajectories used to calculate the MSD are
 110 considered independent when the system enters the diffusive regime from the ballistic, i.e., when
 111 $\langle \Delta r^{*2} \rangle \propto t^{*2}$ changes to $\langle \Delta r^{*2} \rangle \propto t^*$. In this systems, the diffusive regime starts at $t^* \approx 1$, meaning
 112 at least 1000 timesteps are required to reach the change in regime. To ensure our trajectories are
 113 uncorrelated, new simulations are started after 10^5 timesteps have elapsed, which is 100 times
 114 longer than the average time to decorrelation.

115 The shear viscosity is calculated through the Green-Kubo relation of the time-correlation of the
 116 off-diagonal elements of the pressure tensor given by^{15,82,83}

$$\eta^* = \frac{V^*}{T^*} \int_0^\infty \langle p_{\alpha\beta}^*(t^*) p_{\alpha\beta}^*(t_0^*) \rangle dt^* \quad (6)$$

117 where $p_{\alpha\beta}^*(t^*)$, with $\alpha \neq \beta$, is the off-diagonal component of the pressure tensor at time t^* .

118 The training data for the ML models have been collected from 54 different Mie(n, m) potentials.
 119 For each potential, 9 state points are sampled in the supercritical region, 5 state points in the liquid
 120 phase, and also 5 states in the vapour phase. The state points for the supercritical phase are taken
 121 at temperatures of $T^*/T_c^* = \{1.05, 1.25, 1.5\}$ and densities of $\rho^*/\rho_c^* = \{0.2, 1, 2\}$ for each potential,
 122 where T_c and ρ_c are the critical temperature and critical density obtained using the SAFT-VR Mie

123 EoS⁹, respectively. For liquid and vapour phases, 5 state points per potential are chosen at random
124 to ensure that there is no bias in the data collection. For the liquid phase, the state points are
125 chosen by taking random points in the liquid side of the VLE, in the range of $T^*/T_c^* = 0.8 - 0.95$
126 and $\rho^*/\rho_c^* = 1.7 - 2.5$, while ensuring that the resulting points are at least 5% away from the
127 saturated liquid density. A similar approach is employed to study the vapour phase using the
128 same temperature range ($T^*/T_c^* = 0.8 - 0.95$) for densities $\rho^*/\rho_c^* < 0.7$, and ensuring that the
129 selected density is $\rho^* > 0.005$ and at least 5% away from the saturated vapour density. Note
130 that the shear viscosity is only computed and analyzed for the liquid-state region. The following
131 combinations of Mie exponents are studied: $\{n = 12 - 14, m = 6 - 8\}$, $\{n = 15 - 17, m = 6 - 11\}$
132 and $\{n = 18 - 20, m = 6 - 14\}$. These combinations ensure that at least $n - m \geq 4$ and also cover
133 enough parameter space for the study of real molecular systems. All the transport properties of
134 the Mie potentials determined in this work as well as the ML model files are reported in the
135 Supplementary Information (SI) and are available on [GitHub](#).

136 B. Machine Learning

137 In this work, three algorithms to predict the transport properties of Mie fluids are assessed: k -
138 nearest neighbours (KNN), artificial neural network (ANN), and symbolic regression (SR). These
139 algorithms are chosen due to their different levels of complexity and interpretability. KNN and
140 ANN are both “black box” methods, which makes them very difficult to interpret, without studying
141 the resulting algorithm in depth, whereas SR provides a straightforward correlation equation. For
142 all the algorithms, 80% of the data are used to train the model and 20% to test the model. For
143 the case of ANN and KNN, the data is normalized to be in the interval $[-1, 1]$. Before training
144 any model, the self-diffusion coefficient and the shear viscosity are normalized using Chapman-
145 Enskog expressions for the self-diffusion coefficient D_0^* and the shear viscosity η_0^* of a dilute gas
146 of hard spheres, respectively, given by^{10,20}

$$D_0^* = \frac{3}{8} \left(\frac{T^*}{\pi} \right)^{1/2} \frac{1}{\rho^*} \quad (7)$$

147 and

$$\eta_0^* = \frac{5}{16} \left(\frac{T^*}{\pi} \right)^{1/2} \quad (8)$$

148 In other words, all ML methods explored in this work are trained in D/D_0 for the self-diffusion,
149 and in η/η_0 for the viscosity. This semi-empirical approach has also been used by other authors to

150 study systems such as hard spheres¹⁹ and Lennard-Jones fluids⁸⁴. The quantification of errors and
151 accuracy of the models, however, are reported with respect to D^* and η^* . The heat maps of the
152 errors observed using the three different ML in the phase space are presented in the Supplementary
153 Information, as well as the prediction of our method to predict available transport properties of LJ
154 particles.

155 The performance of the ML algorithms with respect to the testing data has been quantified
156 using the coefficient of determination R^2 as well as the absolute average relative deviation (AARD)
157 defined as

$$\text{AARD} = \frac{1}{n} \sum_{i=1}^n \frac{|(y_i - \hat{y}_i)|}{y_i} \times 100\%, \quad (9)$$

158 where n is the number of samples, and y_i and \hat{y}_i indicate the true and predicted values of the sample
159 i , respectively. Both AARD and R^2 are used in assessing the performance of ML algorithms and,
160 in the case of ANN and KNN, are also used in determining the accuracy of a learning method in
161 itself. The calculation of method-training accuracy is done using 10-fold cross-validation (CV10),
162 which entails dividing the test data into 10 different randomly selected and equal-sized sections,
163 training the model with 9 of the sections and validating with the 10th section. This process is
164 repeated until the algorithm uses all sections as validation data, thus allowing for the calculation
165 of the performance using AARD. Both ANN and KNN methods have been implemented in Python
166 3 using the *scikit-learn* library version 1.2.2⁸⁵.

167 The complexity and performance of the ANN are largely influenced by the number of hidden
168 layers and the number of neurons each layer has. Additional layers provide the ability for the
169 ANN to capture more complex input/output dependencies but require more time and data to train.
170 Additionally, more complex ANN architectures may lead to data over-fitting if the underlying
171 correlations are less complex, so the choice of architecture needs to be optimized. In this work,
172 the performances of networks with different numbers of hidden layers and numbers of neurons
173 in the hidden layers are quantified using AARD as the metric for comparison. It is found that
174 a neural network with a single hidden layer consisting of 28 nodes is sufficient and increasing
175 the complexity does not significantly improve the performance. The activation function used is
176 the ReLU function, the training is done for 1000 epochs and the *lbfgs* solver is used for back-
177 propagation.

178 The KNN algorithm can be used for both classification and regression, and uses a number of
179 k nearest neighbours to perform interpolation to predict a new state. In classification, the method

180 uses the distances to these k neighbours to classify the test point, while in regression problems, the
181 values of the neighbours are weighted by the distances to the neighbours to predict the value of
182 the test point. The hyperparameters of the KNN algorithm play a primary role in the performance
183 and efficiency of the algorithm. To select the most appropriate hyperparameters, CV10 is used
184 to find the combination of hyperparameters that lead to the lowest AARD. The hyperparameters
185 validated are the value of k , the weights assigned to the data points, the algorithm to find the
186 closest points, and the power parameter for the Minkowski metric. During the validation, none
187 of the hyperparameters tested shows a significant drop in performance in both quality and speed,
188 with the largest difference in performance only observed when changing the value of k . The
189 hyperparameters that provide the best performance are $k=4$, the neighbouring points weighted
190 by distance, and the power parameter of 4 in the Minkowski metric. The algorithm used for the
191 closest search point does not make any appreciable difference as long as the brute-force search is
192 not used.

193 SR is a very different ML algorithm compared to KNN and ANN. Rather than learning the
194 connections between outputs and inputs in the traditional ML sense, it attempts to find a mathe-
195 matical expression that correlates the behaviour of the output given the inputs by building a binary
196 tree, in which leaves are constants or inputs and branches are mathematical operations. To find
197 the most appropriate equation, the method starts with naive random guesses and uses evolution
198 and mutations to obtain more accurate equations. In this work, the SR implementation from the
199 *gplearn* library version 0.4.1 is used⁸⁶. The hyperparameter space for SR is very large, as the mu-
200 tation and evolution steps can be changed to allow for more or less variability between different
201 members of the population and between parents and offspring equations. Other hyperparameters
202 control the population size, the number of generations of evolution, the types of functions that can
203 be used to connect the nodes, the parsimony coefficient, which is a measure of the 'complexity'
204 of the final equation, and the metric by which the performance of each member of the population
205 is evaluated. Each hyperparameter needed to be checked individually not only for accuracy but
206 also for variability in obtained equations in different random states and the length of the evolution.
207 After some initial testing, it is found that a population size of 5000, 50 generations, the AARD
208 metric and a parsimony coefficient of 0.3 are the hyperparameters that affect the algorithm the
209 most. The operation set used in this work comprises additions, subtractions, multiplications, di-
210 visions, exponential, square roots of the absolute values and the natural logarithm of the absolute
211 values.

212 III. RESULTS

213 A. Self-diffusion coefficient

214 1. Selection of relevant features for the ML methods and conformality of the Mie potential

215 The correct selection of features is critical in ML. For fluids comprising spherical particles
216 interacting *via* the Mie potential, the obvious features that can be used are temperature T^* and
217 density ρ^* to define the thermodynamic state, and both repulsive n and attractive m exponents to
218 define the intermolecular potential. Neither σ nor ε is needed since all the properties are expressed
219 in reduced units. However, as discussed in Section II A, it has been demonstrated that Mie poten-
220 tials with different pairs of exponents (n, m) that lead to the same value of the cohesive parameter
221 α (Equation 2) are conformal and exhibit the same critical points and, potentially, should lead to
222 the same transport properties in the fluid phase. To corroborate this hypothesis, the self-diffusion
223 coefficients obtained from the MD simulations of three triplets of Mie potentials are studied. Here,
224 each triplet has the same value of the cohesive parameter α corresponding to $\alpha = 0.899, 0.585,$
225 and 0.254 but within each set, the repulsive and attractive exponents differ considerably from one
226 another. These sets of potentials have been studied at a subcritical isotherm ($T^*/T_c^* = 0.95$) and at
227 a supercritical isotherm ($T^*/T_c^* = 1.2$) for a wide range of densities, and the results are presented
228 in Figure 1. Despite the small deviations observed at very high densities near the freezing points,
229 the agreement of the self-diffusion coefficient within a set of potentials with an identical value of
230 the cohesive parameter is remarkable, particularly for densities corresponding to $\rho^* < 0.7$, which
231 is to the region in which most of the simulation data has been collected.

232 The predictive power of the cohesive parameter α to describe the self-diffusion coefficient of
233 different potentials reaffirms the need to use it as a feature in the ML algorithm training rather than
234 using the individual exponents n and m . The relation between the cohesive parameter α and the
235 exponents is non-linear, as shown in Equation 3, hence using the individual exponents introduces a
236 layer of complexity between inputs and outputs, that may result in a lower ability of ML algorithms
237 to predict transport properties.

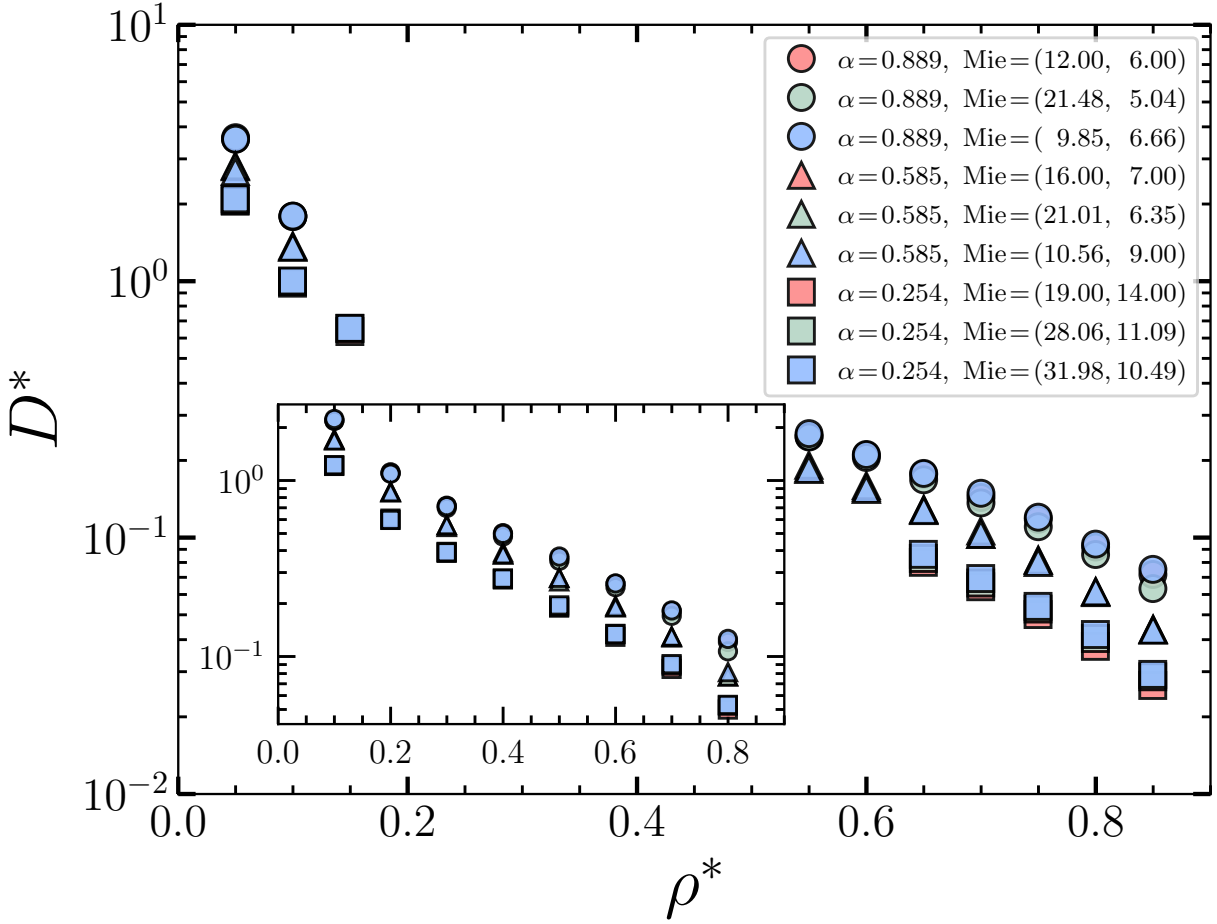


FIG. 1. MD simulation results of the self-diffusion coefficient D^* as a function of the density ρ^* obtained for three triplets of potentials with different cohesive parameter α . The results for the subcritical isotherm $T^*/T_c^* = 0.95$ are presented in the main figure, while the results for the supercritical isotherm $T^*/T_c^* = 1.2$ are shown in the inset.

238 2. *k*-Nearest Neighbours and Artificial Neural Networks

239 The training of the ANN and KNN algorithms using the MD results for the self-diffusion co-
 240 efficient has been performed first according to three regions in the phase diagram, namely the
 241 subcritical vapour phase, the subcritical liquid phase and the supercritical fluid phase, and then
 242 analysed the entire set as a whole. As explained in the previous section, the features used in both
 243 algorithms are the temperature T^* , the density ρ^* , and the cohesive parameter α and the output is
 244 the Chapman-Eskog regularised self-diffusion coefficient (D^*/D_0^*).

245 First, the application of the ANN and KNN algorithms on the vapour phase data set is discussed.

TABLE I. Summary of the AARD and R^2 descriptors of different ML methods applied in this work.

Property	State	KNN		ANN		SR	
		AARD	R^2	AARD	R^2	AARD	R^2
D^*	Vapour phase	0.42%	0.9998	0.87%	0.9998	3.5%	0.995
D^*	Liquid phase	3.3%	0.988	1.8%	0.998	3.5%	0.992
D^*	Supercritical phase	3.4%	0.989	8.4%	0.977	7.7%	0.970
D^*	All phases	2.8%	0.9997	6.4%	0.998	30%	0.995
η^*	Liquid phase	2.8%	0.977	3.3%	0.988	6.3%	0.964

246 The results are presented in Figure 2(a) in the form of parity plots. It is evident from this Figure
247 that in the vapour phase the ANN algorithm performs very similarly to KNN as seen from the
248 summary of the AARD and R^2 presented in Table I. The AARD and R^2 for the ANN are 0.87% and
249 0.9998 for the vapour phase, respectively, compared to 0.42% and 0.9998 obtained using KNN.
250 The algorithms differ in the accuracy of the prediction, to which AARD is more sensitive. The fact
251 that KNN has a lower AARD may be explained by the fact that KNN uses interpolation to learn
252 points to predict new points, rather than predicting the values based on inputs. The differences are
253 minor, however, showing a good ability of both methods to predict the self-diffusion coefficient.

254 For the liquid phase, the algorithms do not predict D^* as well as for the vapour phase, with all
255 metrics decreasing for both methods. The results are presented in Figure 2(b). The decrease in
256 accuracy of the methods is contributed by an increased variation of the range of D^*/D_0^* in the liquid
257 phase. Since the magnitude of D^* in the liquid phase is about two orders of magnitude smaller
258 than in the vapour phase, the values of AARD are larger due to the amplification of the errors. It
259 is observed that R^2 also decreases for the liquid phase, which is likely due to the normalization
260 of D^*/D_0^* between [-1,1] where any deviations in the prediction of the methods will be amplified
261 due to the increased range of the data used ($D^*/D_0^* = [0.08 - 0.56]$ for liquid phase compared
262 to $D^*/D_0^* = [0.59 - 0.77]$ for the vapour phase). Additionally, it is observed that the decrease in
263 accuracy is larger for KNN, particularly in the region of low values of D^* and around $D^* \approx 0.1$
264 where data points are further away from the main diagonal in the parity plots.

265 The MD simulation data for the supercritical fluid phase has been collected at specific values
266 of T^*/T_c^* and ρ^*/ρ_c^* . While the value of T_c^* varies greatly with the cohesive parameter α , the
267 critical density remains mostly constant⁸⁷. Therefore, the ranges of the self-diffusion coefficient

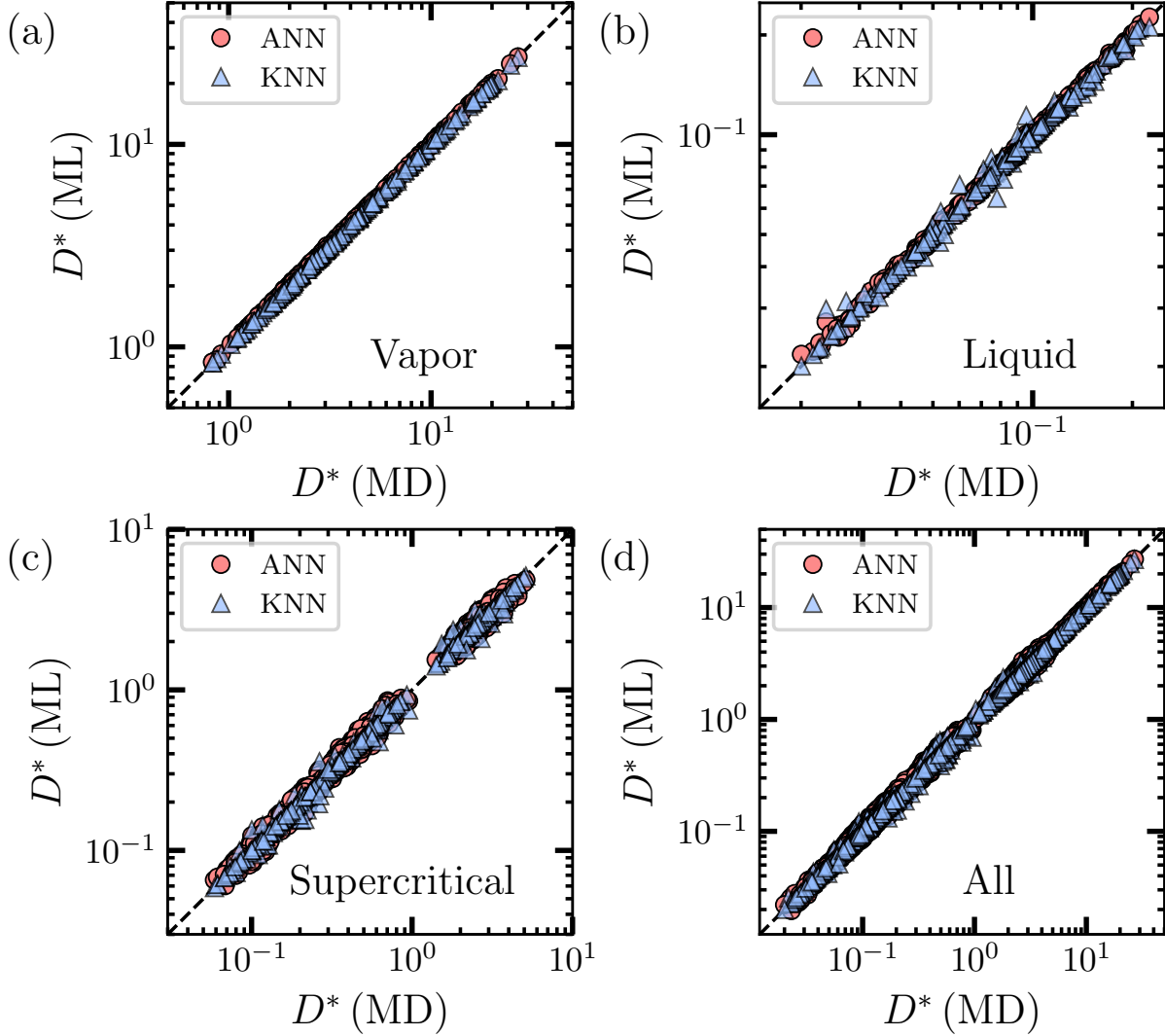


FIG. 2. Parity plots describing the performance of the ANN (red circles) and KNN (blue triangles) on the testing set of the self-diffusion coefficient for three different regions in the phase diagram, as well as for the entire data set as a whole. The results correspond to (a) subcritical vapour, (b) subcritical liquid, (c) supercritical fluid, and (d) all states. D^* (MD) denotes the self-diffusion coefficients obtained from MD simulations, while D^* (ML) denotes the predictions using machine learning.

268 are much more well-defined. The performance of ANN decreases for the supercritical phase with
 269 respect to the liquid phase, while KNN performs similarly, as shown Table I. It is worth noting
 270 that the performance of the methods fitting D^*/D_0^* is slightly lower than when trained on D^* alone
 271 only for the supercritical phase (data not shown).

272 Finally, the application of both ANN and KNN algorithms to describe the entire data set in

273 the three regions of interest is discussed. Taking the entire data set vastly increases the amount
274 of training data available to the algorithms, while increasing the range of output values as well.
275 The performance of the methods is observed in the parity plot in Figure 2(d), where it can be
276 observed that both KNN and ANN are able to similarly predict the value of D^* with KNN slightly
277 outperforming the ANN algorithm.

278 3. *Symbolic Regression*

279 Symbolic regression produces an equation, i.e., a correlation, by applying different mathemat-
280 ical operations to the training data. This methodology produces a good fitting of a mathematical
281 expression, but the drawback is that the algorithm can generate correlations without any physical
282 meaning. Therefore, the use of Chapman-Enskog to normalize the data in the form D^*/D_0^* allows
283 the obtained correlation to be semi-empirical and the correlations obtained from SR simply quan-
284 tify the deviation of D^* with respect to the value of this property for a reference dilute gas of hard
285 spheres.

286 The correlations obtained from the SR of the vapour, liquid, and supercritical states, as well as
287 for the entire data set are presented in Table II. The performance of these correlations is shown
288 in the parity plots presented in Figure 3, and the summary of the values of AARD and R^2 are
289 presented in Table I. The equation for the self-diffusion coefficient of the vapour phase, Equation
290 10, is the simplest of all models and indicates that D^* is approximately 65% of the value of D_0^* ,
291 which makes sense since D_0^* is derived for fluids of infinite dilution. Moreover, the correlation
292 does not include any dependence on the cohesive parameter α , meaning that the expression is
293 independent of the intermolecular potential. Despite the simplicity, this expression is sufficient to
294 predict the self-diffusion coefficient to great accuracy as observed in Figure 3(a), in which both
295 training and testing data sets lie on top of the diagonal in the parity plot. The AARD for this model
296 is 3.5% and $R^2 = 0.995$. The simplicity of the equation is due to the small range in D^*/D_0^* that is
297 observed in the vapour phase.

298 Contrary to the equation obtained for the vapour phase, the one for the liquid phase, Equation
299 11, is complex and introduces new additional terms. One important distinction, however, is the
300 appearance of the cohesive parameter α in the expression suggesting the relevance of the shape
301 of the intermolecular potential in this high-density region and at low temperatures where particle-
302 particle correlations are important. This expression performs equally well when compared to

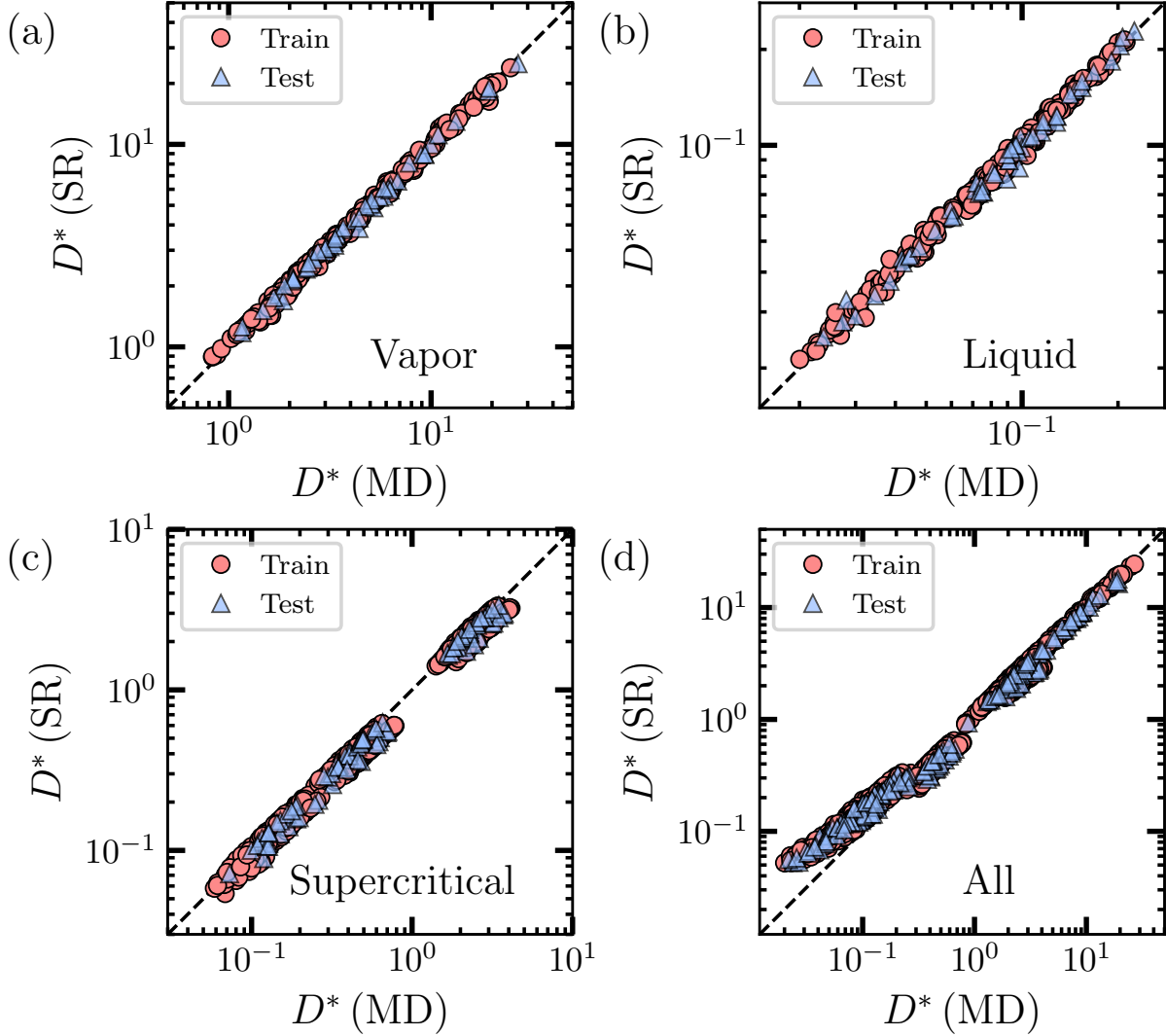


FIG. 3. Parity plots describing the performance of the equations obtained from symbolic regression (SR) for the description of the self-diffusion coefficient for three different regions in the phase diagram, as well as for the entire data set as a whole. The results correspond to (a) subcritical vapour, (b) subcritical liquid, (c) supercritical fluid, and (d) all states. Training data are shown in red circles while the testing data are shown in blue triangles. D^* (MD) denotes the self-diffusion coefficients obtained from MD simulations, while D^* (SR) denotes the predictions using symbolic regression.

303 the KNN algorithm but underperforms in comparison to the ANN algorithm when considering
 304 the performance metrics. The strength of the SR algorithm is its interpretability, as the obtained
 305 equation is very simple to implement, whereas the inner working of the ANN and KNN algorithms
 306 are ‘black boxes’.

TABLE II. Semi-empirical equations obtained from the symbolic regression of the self-diffusion coefficient D^* and the viscosity η^* at different states.

Property	State	Equation
D^*	Vapour phase	$\frac{D^*}{D_0^*} = 0.649$ (10)
D^*	Liquid phase	$\frac{D^*}{D_0^*} = [T^*(1 + \exp(-\alpha) - 0.694)] \frac{\exp(-\rho^*)}{\sqrt{\alpha}}$ (11)
D^*	Supercritical phase	$\frac{D^*}{D_0^*} = \sqrt{(0.664 - \rho^{*2})} \sqrt{\frac{2}{3} T^*}$ (12)
D^*	All states	$\frac{D^*}{D_0^*} = \exp(-\rho^*/T^*) \exp(- 0.273 \log(\log(1/\rho^*)))$ (13)
η^*	Liquid phase	$\frac{\eta^*}{\eta_0^*} = \frac{\exp(\rho^*) [\exp(\rho^*) + \log(\alpha)]}{T^* - 0.405}$ (14)
η^*	Liquid phase	$\eta^* = 0.892 \alpha \rho^* \left(1 + \frac{\rho^*}{T^*(T^* - 0.381)} \right)$ (15)

307 The performance of the equation obtained for the supercritical phase is shown in Figure 3(c)
308 and the expression is presented in Equation 12. This expression is relatively simple compared to
309 the equation for the liquid phase. Due to the lack of very high-density state points in the data
310 set, this equation is only valid for $\rho^* < \sqrt{0.664}$ to ensure that the square root function produces a
311 real value. The lack of high-density state points is also a contributing factor to the absence of any
312 dependence on the cohesive parameter α . In comparison to the algorithms discussed previously
313 for this data set, SR is outperformed by the KNN algorithm, but exhibiting nearly equivalent
314 performance to ANN.

315 Finally, the performance of the SR equation obtained for the entire data set can be observed
316 in Figure 3(d) and is given by Equation 13. This expression exhibits a very good value of R^2 ,
317 but a poor AARD. While it is clear that the equation is unphysical, it has very good predictive
318 power, having a value of R^2 equivalent to the black box ML methods. Much of the error in this
319 equation comes from the overprediction of low values of D^* as can be seen in Figure 3(d), which
320 is the location of the self-diffusion coefficients obtained from the liquid phase simulations. As the
321 SR equation for all phases does not take into account the cohesive parameter α of the potential,
322 the liquid self-diffusion coefficients are not predicted accurately, where the value of the cohesion
323 parameter has a larger influence on the diffusion coefficient due to the large proximity of particles
324 and low temperature of the system. Additionally, the liquid phase data only makes 25% of the
325 full data set, which makes the effects of the cohesion parameter more difficult to be captured for

326 the SR algorithm. The lack of accuracy at low self-diffusion coefficients is balanced by the high
327 precision of high-self diffusion coefficients, which are not as susceptible to changes in the cohesive
328 parameter of the potential.

329 It is also important to note that depending on the value of the initial seed set in the random
330 number generator used in the SR algorithm, there is a larger variability of the expression obtained
331 (except for the vapour phase). The expressions presented in Table II are only a subset of the
332 expressions found by the algorithm.

333 B. Shear viscosity

334 As for the case of the self-diffusion coefficient, the ANN and KNN algorithms, as well as SR,
335 have been applied to represent the shear viscosity η^* of particles interacting *via* the Mie potential.
336 However, only the liquid region is analysed in this work as it is the main region of interest for
337 fluid flow applications. The performance of both ANN and KNN algorithms is shown in Figure
338 4(a). Similar to the case of the self-diffusion coefficients, the ANN exhibits a better performance
339 to predict the shear viscosity of the liquid region with AARD of 3.3% and $R^2 = 0.988$. In contrast,
340 the AARD and R^2 for the KNN algorithm are 2.8% and 0.977, respectively. This difference in
341 predicting power between the measurements shows that KNN is more accurate for low values of
342 η^* , while ANN is more accurate for high values. Interestingly, both methods improve when they
343 are trained using non-normalised data, that is, when the models are train with respect to D^* instead
344 of trained with respect to D^*/D_0^* . This enhancement in their performance is shown in Figure 4(b),
345 especially for ANN, with AARD= 2.21% and $R^2 = 0.996$. KNN receives a smaller improvement
346 with AARD= 2.6% and $R^2 = 0.977$.

347 For the case of the SR applied to the shear viscosity, the performance of the methodology also
348 follows a similar behaviour as in the case of the self-diffusion coefficient. The semi-empirical
349 equation obtained for SR has been obtained by fitting η^*/η_0^* , Equation 14, and the results are
350 presented in Figure 5(a) and in Table I. It is clear from these results that the SR performs worse
351 than the KNN and ANN algorithms, in terms of both R^2 and AARD. Furthermore, the equation
352 is complex and has a dependency on the cohesion parameter α , mirroring the equation obtained
353 for the predictions of the self-diffusion coefficient for the liquid phase. This shows that at high
354 densities there is a more pronounced effect of the exact inter-particle potential and, hence, the
355 requirement for the inclusion of cohesion parameter becomes more apparent. For higher values of

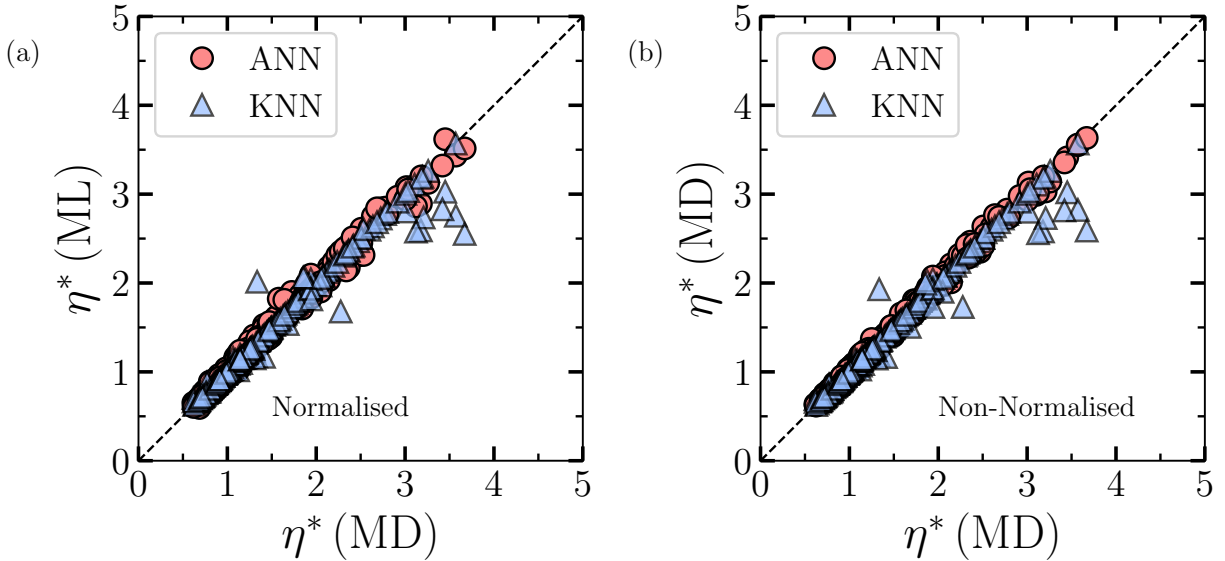


FIG. 4. Parity plot describing the performance of the ANN (red circles) and KNN (blue triangles) on the testing set of the shear viscosity η^* of the liquid phase for (a) normalised and (b) non-normalised (right). η^* (MD) denotes the shear viscosity obtained from MD simulations, while η^* (ML) denotes the predictions using machine learning.

356 η , it is observed that there is an increase in deviations of the data from the parity line. While this
 357 semi-empirical equation performs well, it can be observed in Figure 5(a) that the slope of the data
 358 at low values of η^* also deviate from the main diagonal, implying that there are missing factors in
 359 the equation. For comparison, the SR algorithm has also been trained with respect to η^* instead
 360 of η^*/η_0^* to obtain a fully empirical equation, given by Equation 15, for the shear viscosity and
 361 the results are presented in Figure 5(b). This equation exhibits an AARD=3.8% and a $R^2 = 0.987$,
 362 which shows that the empirical model has slightly better performance in both metrics compared to
 363 the semi-empirical equation, with the value of R^2 outperforming the value for KNN and equaling
 364 the values of ANN. The shape of the points on the fully empirical equation also follows the parity
 365 line more closely, especially at low viscosity values.

366 C. Application to real fluids

367 In this section, the application of the ML models to predict the self-diffusion coefficients of
 368 real fluids is demonstrated. The selected systems correspond to krypton Kr, methane CH₄ and
 369 carbon dioxide CO₂. These systems have been chosen as they can be represented approximately

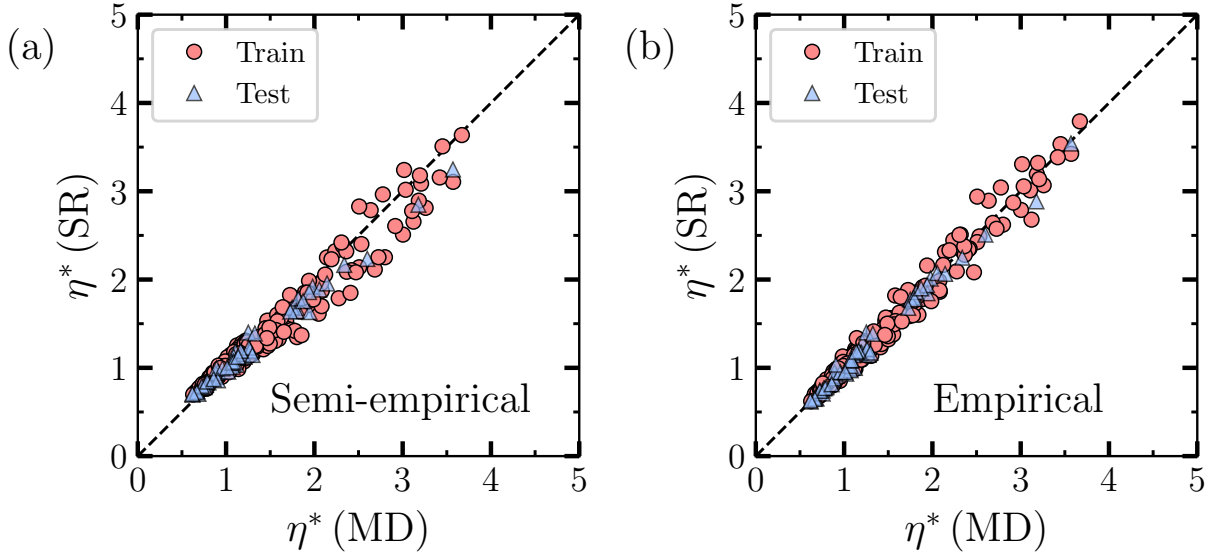


FIG. 5. Parity plots describing the performance of the equations obtained from symbolic regression for the description of the shear viscosity η^* for the liquid phase. The results correspond to (a) the semi-empirical equation (Equation 14) and (b) the empirical equation (Equation 15). Training data are shown in red circles while testing data are shown in blue triangles. η^* (MD) denotes the shear viscosity obtained from MD simulations, while η^* (SR) denotes the predictions using symbolic regression.

370 by a single coarse-grained sphere and exhibit a wide variety of intermolecular interactions that the
 371 Mie potential can capture. The representation of CO₂ as a single coarse-grained sphere has been
 372 previously demonstrated⁷¹ and is used in this work. To find the best Mie potential parameters that
 373 describe the other fluids (Kr, CH₄), the highly accurate SAFT-VR Mie EoS is used⁹. This equation
 374 is based on a high-order Barker and Henderson perturbation theory and is able to represent with
 375 high accuracy the vapour-liquid equilibria, including an excellent representation of the critical
 376 region.

377 The methodology to predict the self-diffusion coefficients of these systems is as follows. First,
 378 experimental vapour-liquid equilibria and vapour pressures have been obtained from the NIST
 379 Chemistry WebBook⁸⁸ and fitted to the SAFT-VR Mie EoS to find the intermolecular parameters
 380 of the Mie potential, i.e. σ , ε , n and m , that best represent the vapour and liquid saturation
 381 densities, as well as the vapour pressure. Using the calculated cohesive parameter α from the
 382 optimized exponents, the ML algorithms optimized for all the fluid regions in the phase diagram
 383 can be deployed to predict the self-diffusion coefficient at any temperature and density. Finally,

TABLE III. Results for the optimized Mie intermolecular parameters ϵ , σ , n and m obtained using the SAFT-VR Mie EoS⁹. The cohesive parameter α is obtained from Equation 2.

System	$\epsilon/(k_B/K)$	$\sigma/\text{\AA}$	n	m	α
Kr	176.34	3.663	14.28	5.98	0.795
CH ₄	160.53	3.754	14.15	5.98	0.800
CO ₂	353.55	3.741	23.0	6.66	0.358

384 the optimized values of σ and ϵ are used to convert from reduced units to real units and the results
 385 are compared with experimental data collected recently by Allers *et al.*⁶³.

386 The optimized Mie intermolecular parameters obtained from the fitting of the SAFT-RV EoS,
 387 and the values used for CO₂ from Avendaño *et al.*⁷¹, are presented in Table III, and the corre-
 388 sponding prediction of the saturated vapour and liquid densities are shown in Figures 6(a,c,e) for
 389 Kr, CH₄, and CO₂, respectively. As can be observed in these Figures, the SAFT-VR Mie provides
 390 an excellent representation of the saturated densities when compared to the experimental results
 391 reported in NIST⁸⁸. Using the parameters reported in Table III, the self-diffusion coefficients for
 392 the three substances are predicted using, for example, the optimized ANN model and can be com-
 393 pared with available experimental data⁶³. The results of this comparison are presented in the parity
 394 plots shown in Figures 6(b,d,f) for Kr, CH₄, and CO₂, respectively. The calculated metrics AARD
 395 and R^2 for the three substances are also reported in Table IV. The predictions using KNN and SR
 396 are also presented in the SI.

397 For KNN and ANN the predictions are equivalent for all three fluids. The predictions are best
 398 for CO₂ in both metrics. When studying the parity plots presented in Figure 6, the parity plot
 399 of CO₂ looks the most accurate over the range of values studied, while for CH₄ and Kr there
 400 are strong deviations at low values of D^* . These deviations may be attributed directly to the
 401 lack of training data in very dense liquids, where the intermolecular forces play a larger role.
 402 Additionally, the coarse-graining applied to the fluid particles treats them as perfect spheres. For
 403 example., CO₂ is neither spherical nor can be represented by simple dispersion forces due to the
 404 strong quadrupolar moments of the molecule. The same effect has been reported by Aimoli *et*
 405 *al.* in their study of transport properties of CO₂ using the same coarse-grained model^{71,73}. The
 406 deviations at low values of D^* are more pronounced in the SR predictions, where the cohesion
 407 parameter is ignored entirely.

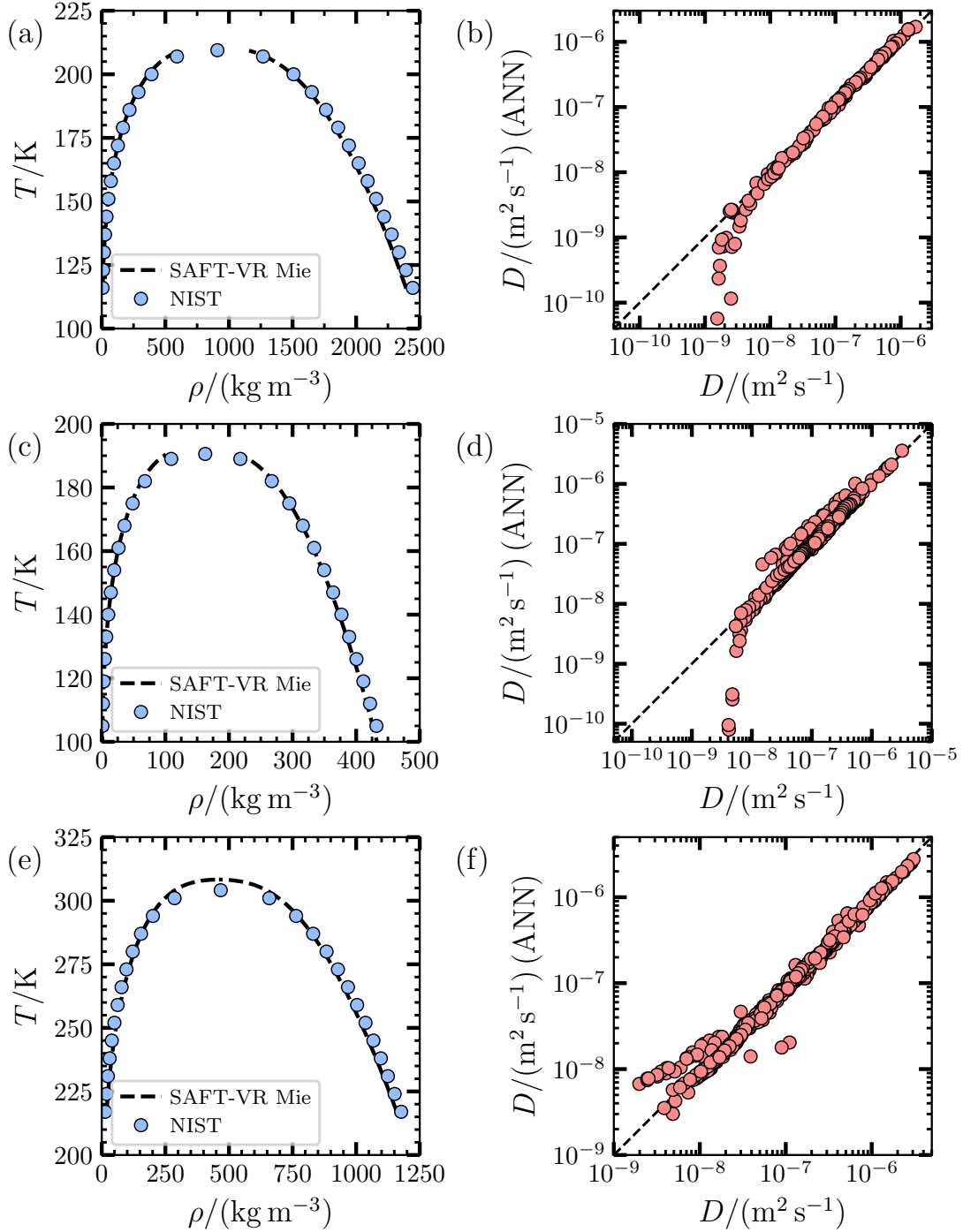


FIG. 6. Results for the vapour-liquid equilibria and self-diffusion coefficients of (a,b) krypton Kr, (c,d) methane CH_4 , and (e,f) carbon dioxide CO_2 . In figures (a, c, d), the blue circles correspond to the experimental results obtained from NIST⁸⁸ while the dashed curves are the predictions using the SAFT-VR Mie Eos⁹ using the parameters reported in Table III. The parity plots shown in Figures (b,d,f) show the comparison of the prediction of the ANN developed in this work with the experimental results collected by Allers *et al.*⁶³.

TABLE IV. The results for the metrics AARD and R^2 obtained from the prediction of the self-diffusion coefficient using the ML methods optimized for all states.

System	KNN		ANN		SR	
	AARD	R^2	AARD	R^2	AARD	R^2
Kr	26.6%	0.970	23.4%	0.968	45.8%	0.906
CH ₄	24.7%	0.970	24.9%	0.971	40.6%	0.902
CO ₂	22.9%	0.972	22.7%	0.977	24.1%	0.944

408 IV. CONCLUSION

409 In summary, the ability of three different ML algorithms to predict transport properties of
 410 spheres interacting *via* the Mie potential has been investigated. The main transport property stud-
 411 ied in this work is the self-diffusion coefficient, which has been obtained from the calculations
 412 of the mean-squared displacement from MD simulations. It has been found that KNN and ANN
 413 perform equivalently over all the fluid phases considered, with KNN having slightly better perfor-
 414 mance in general, especially in terms of AARD. This is attributed to the interpolation used in the
 415 KNN method, which allows for a more accurate prediction of small D^* values. The region stud-
 416 ied in the supercritical regime was trained on a biased data set, resulting in reduced performance
 417 from all methods employed. The SR model performs slightly worse, in general, than the methods
 418 discussed previously. The worst performance from SR is the high AARD in the all-phase data set,
 419 where the lack of the cohesion parameter in the SR expression leads to reduced performance in
 420 the low D^* region. Additionally, for all the ML methods used, viscosity was better predicted when
 421 training without the use of η_0 .

422 The equations obtained from applying the SR algorithm have shown an unexpected result,
 423 where the cohesive parameter, which defines the potential’s fluid phase diagram, is absent for
 424 all the fluid phase regions except for the liquid phase. The prediction of the shear viscosity has
 425 reaffirmed the comparisons for the performance of the algorithms, as the same trends have been
 426 observed. To increase the generalization and applicability of the ML models, the all-phase models
 427 have been applied to predict the self-diffusion coefficients of krypton, methane and carbon dioxide.
 428 To that end, the Mie potential parameters for each fluid have been calculated using SAFT-VR Mie
 429 equation of state using available vapour-liquid experimental data. These parameters have been

430 used along with the previously presented models to predict the self-diffusion coefficient of the
431 three selected systems. While a good overall agreement has been observed, the results have shown
432 some deviations that are attributed to the lack of training data at very high liquid density regime
433 as well as the simplicity of single-sphere coarse-grained representation.

434 V. SUPPLEMENTARY INFORMATION

435 The supplementary information contains a discussion on the uncertainties in the calculation of
436 the transport properties, heat maps representing the accuracy of the models in the phase space,
437 prediction of self-diffusion coefficients of Lennard-Jones particles from published literature, and
438 the predicted results of the self-diffusion coefficient for Kr, CH₄, and CO₂ using the three ML
439 models.

440 VI. ACKNOWLEDGEMENTS

441 This work was supported by the UK Engineering and Physical Sciences Research Council (EP-
442 SRC) via an Industrial Cooperative Award in Science & Technology (ICASE) co-funded by IBM,
443 project ID 2327699 - EP/T517689/1. The authors would like to acknowledge the assistance given
444 by Research IT and the use of the Computational Shared Facility at The University of Manchester.
445 This work was also supported by the Hartree National Centre for Digital Innovation, a collabo-
446 ration between STFC and IBM. A.P. is supported by a “Maria Zambrano Senior” distinguished
447 researcher fellowship, financed by the European Union within the NextGenerationEU program
448 and the Spanish Ministry of Universities.

449 REFERENCES

- 450 ¹E. Hendriks, G. M. Kontogeorgis, R. Dohrn, J.-C. de Hemptinne, I. G. Economou, L. F. Zvilnik,
451 and V. Vesovic, [Ind. Eng. Chem. Res.](#) **49**, 11131 (2010).
- 452 ²K. E. Gubbins, Y.-C. Liu, J. D. Moore, and J. C. Palmer, [Phys. Chem. Chem. Phys.](#) **13**, 58
453 (2011).
- 454 ³B. E. Poling, J. M. Prausnitz, and J. O’Connell, *The properties of gases and liquids*, 5th ed.
455 (McGraw-Hill Professional, New York, NY, 2001).
- 456 ⁴J. O. Valderrama, [Ind. Eng. Chem. Res.](#) **42**, 1603 (2003).

- 457 ⁵M. Michelsen and J. Mollerup, *Thermodynamic Models: Fundamentals & Computational As-*
458 *pects* (Tie-Line Publications, 2004).
- 459 ⁶W. G. Chapman, K. E. Gubbins, G. Jackson, and M. Radosz, *Ind. Eng. Chem. Res.* **29**, 1709
460 (1990).
- 461 ⁷A. Gil-Villegas, A. Galindo, P. J. Whitehead, S. J. Mills, G. Jackson, and A. N. Burgess, *J.*
462 *Chem. Phys.* **106**, 4168 (1997).
- 463 ⁸J. M. Prausnitz, R. N. Lichtenthaler, and E. G. de Azevedo, *Molecular thermodynamics of*
464 *fluid-phase equilibria* (Pearson Education, Philadelphia, PA, 1998).
- 465 ⁹T. Lafitte, A. Apostolakou, C. Avedaño, A. Galindo, C. S. Adjiman, E. A. Müller, and G. Jack-
466 son, *J. Chem. Phys.* **139**, 154504 (2013).
- 467 ¹⁰D. A. McQuarrie, *Statistical Mechanics*, Harper's chemistry series (Longman, London, England,
468 1976).
- 469 ¹¹S. E. Quiñones-Cisneros, C. K. Zéberg-Mikkelsen, and E. H. Stenby, *Fluid Phase Equilib.* **169**,
470 249 (2000).
- 471 ¹²D. R. Reichman and P. Charbonneau, *J. Stat. Mech.* **2005**, P05013 (2005).
- 472 ¹³A. Mulero, ed., *Theory and simulation of hard-sphere fluids and related systems*, 2008th ed.,
473 *Lecture notes in physics* (Springer, Berlin, Germany, 2008).
- 474 ¹⁴D. Frenkel and B. Smit, *Understanding Molecular Simulation*, 2nd ed. (Elsevier, 2002).
- 475 ¹⁵M. P. Allen and D. J. Tildesley, *Computer simulation of liquids* (Oxford University Press, Lon-
476 don, England, 2017).
- 477 ¹⁶P. R. Sassi, P. Mourier, M. H. Caude, and R. H. Rosset, *Anal. Chem.* **59**, 1164 (1987).
- 478 ¹⁷D. Chandler, *J. Chem. Phys.* **62**, 1358 (1975).
- 479 ¹⁸R. J. Speedy, *Mol. Phys.* **62**, 509 (1987).
- 480 ¹⁹R. Speedy, F. Prielmeier, T. Vardag, E. Lang, and H.-D. Lüdemann, *Mol. Phys.* **66**, 577 (1989).
- 481 ²⁰S. Chapman and T. G. Cowling, *The mathematical theory of non-uniform gases: An account of*
482 *the kinetic theory of viscosity, Thermal Conduction and Diffusion in Gases* (Cambridge Univer-
483 sity Press, 1990).
- 484 ²¹F. Demmel, D. Szubrin, W.-C. Pilgrim, and C. Morkel, *Phys. Rev. B* **84**, 014307 (2011).
- 485 ²²H. Liu and C. M. Silva, *Ind. Eng. Chem. Res.* **36**, 246 (1997).
- 486 ²³N. M. Blagoveshchenskii, A. G. Novikov, and V. V. Savostin, *Phys. Solid State* **56**, 120 (2014).
- 487 ²⁴N. Blagoveshchenskii, A. Novikov, A. Puchkov, V. Savostin, and O. Sobolev, *EPJ Web of*
488 *Conferences* **83**, 02018 (2015).

- 489 ²⁵G. V. Kharlamov and S. V. Zhilkin, *J. Phys.: Conf. Ser.* **899**, 052009 (2017).
- 490 ²⁶Y. A. Aljeshi, M. B. M. Taib, and J. P. M. Trusler, *Int. J. Thermophys.* **42**, 140 (2021).
- 491 ²⁷C. Corral-Casas, L. Gibelli, M. K. Borg, J. Li, S. F. K. Al-Afnan, and Y. Zhang, *Phys. Fluids*
492 **33**, 082009 (2021).
- 493 ²⁸R. L. Rowley and M. M. Painter, *Int J Thermophys* **18**, 1109 (1997).
- 494 ²⁹R. Laghaei, A. E. Nasrabad, and B. C. Eu, *J. Phys. Chem. B* **109**, 5873 (2005).
- 495 ³⁰L. V. Woodcock, *AIChE Journal* **52**, 438 (2006).
- 496 ³¹X. Wang, S. Ramírez-Hinestrosa, J. Dobnikar, and D. Frenkel, *Phys. Chem. Chem. Phys.* **22**,
497 **10624** (2020).
- 498 ³²S. E. Quiñones-Cisneros, C. K. Zéberg-Mikkelsen, J. Fernández, and J. García, *AIChE J.* **52**,
499 **1600** (2006).
- 500 ³³A. S. de Wijn, V. Vesovic, G. Jackson, and J. P. M. Trusler, *J. Chem. Phys.* **128**, 204901 (2008).
- 501 ³⁴F. Llovell, R. M. Marcos, and L. F. Vega, *J. Phys. Chem. B* **117**, 8159 (2013).
- 502 ³⁵S. Becker, E. Devijver, R. Molinier, and N. Jakse, *Phys. Rev. E* **105**, 045304 (2022).
- 503 ³⁶W. F. Reinhart, A. W. Long, M. P. Howard, A. L. Ferguson, and A. Z. Panagiotopoulos, *Soft*
504 *Matter* **13**, 4733 (2017).
- 505 ³⁷S. Dassetty, I. Coropceanu, J. Portner, J. Li, J. J. de Pablo, D. Talapin, and A. L. Ferguson, *Mol.*
506 *Syst. Des. Eng.* **7**, 350 (2022).
- 507 ³⁸J. O’Leary, R. Mao, E. J. Pretti, J. A. Paulson, J. Mittal, and A. Mesbah, *Soft Matter* **17**, 989
508 (2021).
- 509 ³⁹J. L. McDonagh, W. C. Swope, R. L. Anderson, M. A. Johnston, and D. J. Bray, *Polym. Int.* **70**,
510 **248** (2021).
- 511 ⁴⁰E. Boattini, M. Dijkstra, and L. Filion, *J. Chem. Phys* **151**, 154901 (2019).
- 512 ⁴¹M. Spellings and S. C. Glotzer, *AIChE J.* **64**, 2198 (2018).
- 513 ⁴²A. Fabrizio, B. Meyer, and C. Corminboeuf, *J. Chem. Phys.* **152**, 154103 (2020).
- 514 ⁴³C. Dietz, T. Kretz, and M. H. Thoma, *Phys. Rev. E* **96**, 011301 (2017).
- 515 ⁴⁴J. L. McDonagh, N. Nath, L. De Ferrari, T. van Mourik, and J. B. O. Mitchell, *J. Chem. Inf.*
516 *Model.* **54**, 844 (2014).
- 517 ⁴⁵J. L. McDonagh, D. S. Palmer, T. v. Mourik, and J. B. O. Mitchell, *J. Chem. Inf. Model.* **56**,
518 **2162** (2016).
- 519 ⁴⁶J. L. McDonagh, T. van Mourik, and J. B. O. Mitchell, *Mol. Inform.* **34**, 715 (2015).
- 520 ⁴⁷L. Joss and E. A. Müller, *J. Chem. Educ.* **96**, 697 (2019).

521 ⁴⁸K. Zhu and E. A. Müller, *J. Phys. Chem. B* **124**, 8628 (2020).

522 ⁴⁹A. S. Alshehri, A. K. Tula, F. You, and R. Gani, *AIChE J.* **68**, e17469 (2022).

523 ⁵⁰K. K. Yalamanchi, V. C. O. van Oudenhoven, F. Tutino, M. Monge-Palacios, A. Alshehri,
524 X. Gao, and S. M. Sarathy, *J. Phys Chem. A* **123**, 8305 (2019).

525 ⁵¹G. B. Goh, C. Siegel, A. Vishnu, N. O. Hodas, and N. Baker, *IEEE*, 1340 (2018).

526 ⁵²V. Mann, K. Brito, R. Gani, and V. Venkatasubramanian, *Fluid Ph. Equilibria* **561**, 113531
527 (2022).

528 ⁵³A. Tihic, G. M. Kontogeorgis, N. von Solms, M. L. Michelsen, and L. Constantinou, *Ind. Eng.*
529 *Chem. Res.* **47**, 5092 (2007).

530 ⁵⁴K. Hansen, F. Biegler, R. Ramakrishnan, W. Pronobis, O. A. von Lilienfeld, K.-R. Müller, and
531 A. Tkatchenko, *J. Phys. Chem. Lett.* **6**, 2326 (2015).

532 ⁵⁵J. Schmidt, M. R. G. Marques, S. Botti, and M. A. L. Marques, *Npj Comput. Mater.* **5**, 83
533 (2019).

534 ⁵⁶A. Abbasi Moud, *Colloids Interface Sci. Commun.* **47**, 100595 (2022).

535 ⁵⁷M. Ihme, W. T. Chung, and A. A. Mishra, *Prog. Energy Combust. Sci.* **91**, 101010 (2022).

536 ⁵⁸Y. Elbaz, D. Furman, and M. Caspary Toroker, *Adv. Funct. Mater.* **30**, 1900778 (2020).

537 ⁵⁹A. Abbasi and R. Eslamloueyan, *Chemometr. Intell. Lab. Syst.* **132**, 39 (2013).

538 ⁶⁰X. Zhao, T. Luo, and H. Jin, *Ind. Eng. Chem. Res.* **61**, 8542 (2022).

539 ⁶¹N. Melzi, L. Khaouane, S. Hanini, M. Laidi, Y. Ammi, and H. Zentou, *J. Appl. Mech. Tech.*
540 *Phys.* **61**, 207 (2020).

541 ⁶²J. Zhou, S. Chupradit, K. Ershov, W. Suksatan, H. Abdulameer Marhoon, M. Alashwal, S. Ghaz-
542 ali, M. Algarni, and A. El-Shafay, *J. Mol. Liq.* **353**, 118808 (2022).

543 ⁶³J. P. Allers, F. H. Garzon, and T. M. Alam, *Phys. Chem. Chem. Phys.* **23**, 4615 (2021).

544 ⁶⁴J. P. S. Aniceto, B. Zêzere, and C. M. Silva, *J. Mol. Liq.* **326**, 115281 (2021).

545 ⁶⁵J. P. Allers, J. A. Harvey, F. H. Garzon, and T. M. Alam, *J. Chem. Phys.* **153**, 034102 (2020).

546 ⁶⁶K. Papastamatiou, F. Sofos, and T. E. Karakasidis, *AIP Advances* **12**, 025004 (2022).

547 ⁶⁷C. J. Leverant, J. A. Harvey, and T. M. Alam, *J. Phys. Chem. Lett.* **11**, 10375 (2020).

548 ⁶⁸T. M. Alam, J. P. Allers, C. J. Leverant, and J. A. Harvey, *J. Chem. Phys.* **157**, 014503 (2022).

549 ⁶⁹G. Mie, *Ann. Phys. (Berl.)* **11**, 657 (1903).

550 ⁷⁰X. He, W. Shinoda, R. DeVane, and M. L. Klein, *Mol. Phys.* **108**, 2007 (2010).

551 ⁷¹C. Avendaño, T. Lafitte, A. Galindo, C. S. Adjiman, G. Jackson, and E. A. Müller, *J. Phys.*
552 *Chem. B* **115**, 11154 (2011).

553 ⁷²C. Avendaño, T. Lafitte, C. S. Adjiman, A. Galindo, E. A. Müller, and G. Jackson, *J. Phys.*
554 *Chem. B* **117**, 2717 (2013).

555 ⁷³C. G. Aimoli, E. J. Maginn, and C. R. A. Abreu, *J. Chem. Phys.* **141**, 134101 (2014).

556 ⁷⁴J. C. Shelley, M. Y. Shelley, R. C. Reeder, S. Bandyopadhyay, and M. L. Klein, *J. Phys. Chem.*
557 *B* **105**, 4464 (2001).

558 ⁷⁵J. J. Potoff and D. A. Bernard-Brunel, *J. Phys. Chem. B* **113**, 14725 (2009).

559 ⁷⁶N. S. Ramrattan, C. Avendaño, E. A. Müller, and A. Galindo, *Mol. Phys.* **113**, 932 (2015).

560 ⁷⁷K. S. Pitzer, *J. Chem. Phys.* , 583 (1939).

561 ⁷⁸E. Helfand and S. A. Rice, *J. Chem. Phys.* , 1642 (1960).

562 ⁷⁹J. A. Barker and D. Henderson, *Rev. Mod. Phys.* **48**, 587 (1976).

563 ⁸⁰W. G. Hoover, *Phys. Rev. A* **31**, 1695 (1985).

564 ⁸¹A. P. Thompson, H. M. Aktulga, R. Berger, D. S. Bolintineanu, W. M. Brown, P. S. Crozier,
565 P. J. in 't Veld, A. Kohlmeyer, S. G. Moore, T. D. Nguyen, R. Shan, M. J. Stevens, J. Tranchida,
566 C. Trott, and S. J. Plimpton, *Comp. Phys. Comm.* **271**, 108171 (2022).

567 ⁸²R. Kubo, *J. Phys. Soc. Jpn.* **12**, 570 (1957).

568 ⁸³M. S. Green, *J. Chem. Phys.* **22**, 398 (1954).

569 ⁸⁴Y. Zhu, X. Lu, J. Zhou, Y. Wang, and J. Shi, *Fluid Phase Equilib.* **194-197**, 1141 (2002).

570 ⁸⁵F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pret-
571 tenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot,
572 and E. Duchesnay, *J. Mach. Learn. Res.* **12**, 2825 (2011).

573 ⁸⁶T. Stephens, “Gplearn: Genetic programming in python, with a scikit-learn inspired API,”
574 (2015).

575 ⁸⁷J. Šlepavičius, C. Avendaño, B. O. Conchúir, and A. Patti, *Phys. Rev. E* **106**, 014604 (2022).

576 ⁸⁸P. Linstrom, “NIST Chemistry WebBook, NIST Standard Reference Database 69,” (1997).