

Towards predicting liquid fuel physicochemical properties using molecular dynamics guided machine learning models

Rodolfo S. M. Freitas^a, Ágatha P. F. Lima^a, Cheng Chen^b, Fernando A. Rochinha^a, Daniel Mira^c and Xi Jiang^{b,*}

^aCOPPE, Federal University of Rio de Janeiro, Rio de Janeiro, Rio de Janeiro 21941-598, Brazil

^bSchool of Engineering & Materials Science, Queen Mary University of London, Mile End Road, London E1 4NS, UK

^cBarcelona Supercomputing Center (BSC), Barcelona, Spain

ARTICLE INFO

Keywords:

Fuel properties
Molecular dynamics
Deep Learning
Machine learning models

ABSTRACT

Accurate determination of fuel properties of complex mixtures over a wide range of pressure and temperature conditions is essential to utilizing alternative fuels. The present work aims to construct cheap-to-compute machine learning (ML) models to act as closure equations for predicting the physical properties of alternative fuels. Those models can be trained using the database from MD simulations and/or experimental measurements in a data-fusion-fidelity approach. Here, Gaussian Process (GP) and probabilistic generative models are adopted. GP is a popular non-parametric Bayesian approach to build surrogate models mainly due to its capacity to handle the aleatory and epistemic uncertainties. Generative models have shown the ability of deep neural networks employed with the same intent. In this work, ML analysis is focused on two particular properties, the fuel density and diffusion, but it can also be extended to other physicochemical properties. This study explores the versatility of the ML models to handle multi-fidelity data. The results show that ML models can predict accurately the fuel properties of a wide range of pressure and temperature conditions.

1. Introduction

Fossil fuels have been playing a major role in energy supply and liquid fossil fuels have dominated the energy use in transport, which will continue to be so for many decades to come, especially for sectors that are difficult to decarbonise [1, 2]. With the pressing needs of decarbonisation and sustainable energy utilisation, renewable fuels and biofuels are becoming

*Corresponding author. Email: xi.jiang@qmul.ac.uk

ORCID(s): 0000-0001-6036-8534 (R.S.M. Freitas); 0000-0002-2155-6185 (Á.P.F. Lima); 0000-0001-7292-9490 (C. Chen); 0000-0001-8035-9651 (F.A. Rochinha); 0000-0001-9901-7942 (D. Mira); 0000-0003-2408-8812 (X. Jiang)

Nomenclature**Abbreviations**

ANN	Artificial neural network
CFD	Computational Fluid Dynamics
CN	Cetane number
EMD	Equilibrium Molecular Dynamics
EoS	Equation of State
FAME	Fatty acid methyl ester
GANs	Generative Adversarial Networks
GP	Gaussian Process
MD	Molecular dynamics
ML	Machine learning
MLPNNs	Multilayer Perceptron Neural Networks
NARGP	Nonlinear autoregressive multifidelity Gaussian Process
NEMD	Nonequilibrium Molecular Dynamics
NIST	National Institute of Standards and Technology
OMEs	Oxymethylene Dimethyl Ethers
TraPPE	Transferable Potential for Phase Equilibria
VAE	Variational auto-encoders

Greek letters

β	Residual penalty parameter
θ	A vector of hyper-parameters
γ	A generic property

λ	Entropy regularization parameter
μ	Expected value
ϕ	A vector of parameters
ρ	Density
σ	Standard deviation
ξ	A potential noisy

Latin letters

\mathbf{x}, \mathbf{y}	Input and output vectors
cv	Coefficient of variation
C	Number of atoms of carbon
D	Diffusion coefficient
f	Gaussian function
g	Mapping function
K	Covariance matrix
k	A kernel function
l	Correlation length
n	Dimension of the input and output
N_s	Number of samples
P	Pressure
p	Probability distribution
P_c	Critical pressure
T	Temperature
t	Time
T_c	Critical temperature
z	Latent variable

6 increasingly important [3, 4]. For instance, synthetic fuels like Oxymethylene Dimethyl Ethers
7 (OMEs) have shown high potential for low-carbon transport applications due to their capacity to
8 avoid soot formation [5]. However, the physicochemical properties of these fuels must be known
9 for their rapid integration into current infrastructures for storage, transport and direct injection

R. S. M. Freitas, Á. P. F. Lima, C. Chen, F. A. Rochinha, D. Mira & X. Jiang / Preprint submitted to
10 in combustion engines. This represents a significant challenge, due to the fact that practical
11 fuels are often composed by complex mixtures and vary widely in their chemical compositions
12 depending on the production source and process [3]. For example, petroleum diesel is a complex
13 mixture involving molecules with carbon chains that typically contain between 9 and 25 carbon
14 atoms per molecule. To simplify the complex chemical compositions of these fuels, surrogate
15 models have been used to represent the chemical composition and combustion characteristics in
16 practical applications [6, 7]. In addition, modern combustion engines have to operate at high
17 pressure conditions in order to improve the energy conversion efficiency. Fuel properties at extreme
18 conditions such as high pressure and high temperature conditions, are very difficult to measure and
19 predict [5], leading to an additional challenge.

20 Accurate determination of fuel properties of complex mixtures over a wide range of pressure
21 and temperature conditions is essential to adapt the system operation to alternative fuels. In
22 recent years, molecular dynamics (MD) simulations have been used to predict the physicochemical
23 properties of practical fuels including transport properties at supercritical conditions [8]. By using
24 equilibrium molecular dynamics (EMD) and nonequilibrium molecular dynamics (NEMD), Yang
25 et al [9, 10] predicted the viscosity and thermal conductivity of alkanes (n-decane, n-undecane
26 and n-dodecane). Kondratyuk et al [11, 12, 13] performed a serial of MD simulation to study the
27 viscosity of hydrocarbons (1-methylnaphthalene, methylcyclohexane and 2,2,4-trimethylhexane)
28 in high pressure conditions up to 1000 MPa. Coleman et al [14] tested the capacity of existing
29 force fields on prediction of properties (density, enthalpy of vaporization, surface tension and
30 heat capacity etc) of organic liquids. Although MD simulations provide molecular details that can
31 be potentially used to accurately predict fuel properties, they are generally expensive in terms
32 of computational costs (CPU time and memory). In addition, MD predictions also need to be
33 validated against experimental measurements, which can be even more costly especially at extreme
34 conditions. Accordingly, it is not feasible to establish complete and detailed fuel property databases
35 consisting of a wide range of pressure and temperature conditions using either MD simulations or
36 experiments.

37 Machine learning has great potentials to discover the relation between inputs and outputs in
38 a thermodynamic system directly from the data of complex systems [15] and for predicting the
39 properties of materials based on their composition [16]. ML can be a powerful tool to predict
40 fuel properties from chemical compositions of the fuel mixture and/or chemical structures of the
41 fuel molecules. Several works have been devoted to designing ML models capable of predicting
42 complex fuels properties from experimental data. In this regard, ML models obtained accurate
43 predictions of cetane number (CN) compared to experimental data [17, 18, 19]. A satisfactory
44 ML approach for modeling the CN of biodiesel based on four operating conditions given by
45 iodine volume (IV), carbon number, double bounds, and saponification value was proposed [20].
46 Recently, an artificial neural network (ANN) was applied to predict and identify the underlying
47 links between the fuel properties and the octane number (ON) [21]. Moreover, ML models were
48 tuned with evolutionary algorithms to predict the CN of biodiesel as a function of its fatty acid
49 methyl ester (FAME) profile [22, 23]. The predictability, i.e. the ability to predict, of the ML
50 approaches also can be improved by using different optimization algorithms for the training and/or
51 hyperparameter search such as teaching-learning based optimization (TLBO), backpropagation,
52 Quasi-Newton and particle swarm optimization (PSO) [24, 25, 26]. Also, ML models have been
53 used for modeling the kinematic viscosity of diesel-derived fuels as a function of their FAMEs
54 profiles [27, 28, 29]. In the last years, Multilayer Perceptron Neural Networks (MLPNNs) have
55 been successfully built to estimate the physicochemical characteristics of biodiesel [30, 31, 32, 33]
56 combining different parameters of model inputs. Furthermore, ML models based on state variables
57 such as temperature and pressure showed high potential to obtain physicochemical properties of
58 biodiesel/diesel fuels more accurately [34, 35, 36]. In particular, ML models have been developed
59 to predict thermodynamic properties such as critical pressure and temperature, vapor pressures,
60 and densities of pure fluids [37]. Moreover, approaches combining MD simulations and ML have
61 been applied to modeling the diffusion of pure liquids [38, 39]. Following the same context, a ML
62 approach based on support vector regression (SVR) was proposed by [40] for predicting the PVT
63 properties of pure fluids (H_2O , CO_2 , and H_2) and their mixtures, where the training database is

R. S. M. Freitas, Á. P. F. Lima, C. Chen, F. A. Rochinha, D. Mira & X. Jiang / Preprint submitted to
64 provided by the National Institute of Standards and Technology (NIST) and MD simulations. Also,
65 an ML approach was proposed to assess the macroscopic Engine Combustion Network (ECN)
66 Spray-A characteristics and predictions of fluid properties for the thermodynamic states found in
67 such conditions [41]. Yet, from our knowledge, little work has been dedicated towards exploring the
68 thermodynamic properties of practical fuels combining MD simulations and ML models. ML can
69 be a powerful tool to predict fundamental fuel properties directly from the chemical compositions of
70 the fuel mixture by using databases from MD simulations or available experimental measurements.

71 The aim of the study was to demonstrate and validate a ML-MD methodology to predict
72 fundamental properties of liquid fuels. In this approach, the ML models are built from data provided
73 by MD simulations, while a combination of MD and NIST data is used for model assessment and
74 validation. This study is the first attempt of using ML models with Gaussian process regression
75 [42] and probabilistic conditional generative learning [43, 44] for the property predictions of
76 single-compounds. The ML analysis is focused on fuel density in this study as one of fundamental
77 properties of liquid fuels, though it can easily be extended to other physicochemical properties
78 of relevance for practical applications like diffusion coefficient, viscosity, conductivity or surface
79 tension.

80 The rest of the paper is organized as follows. Section 2 presents the ML models and the
81 molecular dynamics simulation methodology. Section 3 describes the ML results for typical fuel
82 surrogates of diesel. Finally, Section 4 concludes the study with recommendation for further
83 investigations.

84 **2. Methodology: Building Machine Learning Models to Describe Physicochemical Properties**

85 In order to reduce energy consumption and pollutant formation, supercritical combustion
86 has been increasingly explored in the context of high pressure internal combustion engines and
87 rocket engines [45]. Specifically, in supercritical conditions, the devices operate with pressures
88 and temperatures higher than the critical values, which implies that physicochemical properties
89 of fluids are quite different from those at liquid conditions [46]. In such scenarios, the design of

R. S. M. Freitas, Á. P. F. Lima, C. Chen, F. A. Rochinha, D. Mira & X. Jiang / Preprint submitted to
90 devices become more complex, specially due to limitations of replicating flow and combustion in
91 controlled laboratory environments. In order to cope with these challenges, computational models
92 can provide adequate tools for obtaining more accurate predictions of state variables and increase
93 cycle performance in transcritical conditions.

94 From a computational fluid dynamics (CFD) perspective, combustion models are built upon
95 the combination of solid and reliable physico/chemical principles with closure models, typically
96 describing physicochemical properties of the fuels and their mixtures using approaches that nor-
97 mally entail uncertainties. The use of numerical simulations for practical applications encompass a
98 wide range of conditions, resulting in different fundamental problems depending on the nozzle
99 geometry, engine architecture or thermodynamic conditions. A good example is the database
100 from the Engine Combustion Network [47] for which different sprays for diesel- and gasoline-like
101 conditions are investigated. For instance, pressure can go from sub-atmospheric to 2,000 bar, and
102 temperatures from cold to highly preheated conditions. In that context, having accurate values for
103 macroscopic fuel characteristics and properties over such wide variety of spatial and time scales
104 is one of the main challenges for physically-driven methods. That is particularly more dramatic
105 for modern compounds depicting complex chemical compositions, and simplified surrogate fuels
106 [48] are employed to estimate the properties of the original compounds. That allows the systematic
107 use of controlled experiments and, also, Molecular Dynamics simulations [49, 50]. Indeed, here
108 our focus lies on using ML models to leverage such type of simulations when obtaining liquid
109 fuel physicochemical properties. Those properties are generally expressed as functions of local
110 thermodynamic conditions like pressure and temperature, which motivate to refer to closure models
111 such as the Equations of State (EoS). In general, the EoS is embedded in complex CFD simulations
112 resulting in divergence or numerical oscillations when used with traditional methods based on
113 tabular and interpolation schemes [51]. It is worth to remark that we are seeking for models capable
114 of describing physicochemical properties over a wide range of flow conditions and we expected to
115 observe abrupt changes around critical conditions.

116 We built two different ML models, namely Gaussian Processes (GP) [42] and a probabilistic
117 conditional generative approach [44]. We train both in a supervised learning fashion using data
118 produced with expensive MD simulations. Therefore, we rely on their ability to learn from a
119 small amount of data and their capacity of extrapolation. Moreover, we also want to take into
120 consideration the unavoidable uncertainties arising from limited information (epistemic) and from
121 noisy data (aleatoric).

122 GPs have become popular due to its success on being a proxy for physics-based high-fidelity
123 models in different applications [52, 53, 54, 55, 56, 57]. Another well proved ML approach are the
124 so called generative models that explore existing low-dimensional structures capable of explaining
125 high-dimensional data introducing probabilistic latent variables.

126 In the remainder of this chapter, we present a brief description of both ML models for a generic
127 property $\gamma(P, T)$ function of pressure and temperature, along with the corresponding training
128 algorithms. For the training of the models, we assume the availability of, potentially expensive,
129 dataset comprising input/output pairs $\{(P, T)_i, \gamma_i \mid i = 1, \dots, n\}$ generated by an implicit mapping
130 g characterizing the macroscopic thermodynamic relation between the property and the state
131 variables:

$$\gamma = g(P, T; \xi). \quad (1)$$

132 The role of g here is played by upscaling MD simulations or, to a less extent, by experimental
133 available data. The vector ξ denotes potential noisy and is often considered a random. In order
134 to keep a compact notation, we refer to the above dataset as $\mathcal{D} = (\mathbf{x}, \mathbf{y})$, with $\mathbf{x} \in \mathbb{R}^{2n}$ and
135 $\mathbf{y} \in \mathbb{R}^n$ vectors containing inputs and outputs. We intentionally do not use the word surrogate to
136 designate any of the two ML models to avoid misleading. In the combustion technical literature,
137 it is employed to refer to compounds with simpler compositions to replace complex fuels in
138 experimental or numerical analysis.

139 2.1. Gaussian process regression

140 A GP is an infinite collection of random variables, in which any finite number of such variables
 141 depict a joint Gaussian distribution [42]. In line with Bayesian estimation, to approximate g we
 142 assign a GP zero mean prior $f(\mathbf{x})$, i.e., $f \sim GP(f|\mathbf{0}, k(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}))$, where k is a kernel parametrized
 143 by a vector of hyper-parameters $\boldsymbol{\theta}$ to be learned from \mathcal{D} and engenders a symmetric positive-definite
 144 $n \times n$ covariance matrix $K_{ij} = k(x_i, x_j; \boldsymbol{\theta})$. Instead of choosing the squared exponential form of the
 145 kernel as usual [42], here, we test some forms of covariance matrix belonging to the Matern family.
 146 More specifically, we employ the Matern 3/2 covariance matrix given as

$$k(\mathbf{r}) = \sigma^2 \left(1 + \sqrt{6} \frac{|\mathbf{r}|}{l} \right) \exp \left(-\sqrt{6} \frac{|\mathbf{r}|}{l} \right) \quad (2)$$

147 with $\mathbf{r} = \mathbf{x} - \mathbf{x}'$ denoting the distance between different inputs. The hyper-parameters are the
 148 standard deviation σ , and the correlation lengths $\mathbf{l} = \{l_1, l_2, \dots, l_{n_k}\}$, and n_k denotes the dimension
 149 of input \mathbf{r} . Hence, the hyper-parameters vector reduces to $\boldsymbol{\theta} = \{\mathbf{l}, \sigma\}$.

150 We do not follow a fully Bayesian approach, and obtain the vector of hyper-parameters $\boldsymbol{\theta}$ by
 151 maximizing the marginal log-likelihood of the model, i.e.

$$\log p(\gamma|\mathbf{x}, \boldsymbol{\theta}) = -\frac{1}{2} \log |\mathbf{K}| - \frac{1}{2} \gamma^T \mathbf{K}^{-1} \gamma - \frac{n}{2} \log 2\pi. \quad (3)$$

152 using a conjugate gradient descend method.

153 The final goal of the regression is obtaining a predictive model for γ , which means to compute
 154 its value for an untested state \mathbf{x}_* [53]

$$\mu_*(\mathbf{x}_*) = k_{*n} \mathbf{K}^{-1} \mathbf{y} \quad (4)$$

155 and

$$\sigma_*^2(\mathbf{x}_*) = k_{**} - k_{*n} \mathbf{K}^{-1} k_{*n}^T \quad (5)$$

156 where $k_{*n} = [k(\mathbf{x}_*, \mathbf{x}_1), \dots, k(\mathbf{x}_*, \mathbf{x}_n)]$ and $k_{**} = k(\mathbf{x}_*, \mathbf{x}_*)$. The predictions are computed using the
 157 posterior mean μ_* , and the uncertainty associated with that predictions is quantified through the
 158 posterior variance σ_*^2 . It is worth to mention that in absence of noisy in the training data, the later
 159 represents epistemic uncertainty due to lack of data.

160 2.2. Probabilistic conditional generative model

161 Now, we explore a probabilistic conditional generative approach [43, 44], that integrates
 162 variational auto-encoders (VAE) [58] and generative adversarial networks (GANs) [59]. Moreover,
 163 it employs a probabilistic perspective that enables to take into consideration noisy and limited data
 164 from the beginning. It is also capable of dealing with high-dimensionality of inputs and outputs,
 165 what is not explored here due to the specific aspects of our needs.

166 The final goal is to build probabilistic neural networks that follow a conditional probability
 167 density function $p(\gamma|(P, T), \mathcal{D})$ learnt from the data. So, the surrogate model can deploy accurate
 168 values for the property γ by estimating the expectation $\mathbb{E}(\gamma|(P, T), \mathcal{D})$, and also, to quantify the
 169 uncertainty associated with that prediction in CFD calculations.

170 The main ingredient for this approach is the introduction of a vector of latent random variables
 171 aiming at seeking for a hidden low dimensional structure for explaining the data structure. In a
 172 formal abstract perspective, such latent variables allow us to express the conditional probability
 173 associate to the data \mathcal{D} , not included in the expression to keep the notion clear, $p(\gamma|P, T)$, as an
 174 infinite mixture model through

$$p(\gamma|P, T) = \int p(\gamma, \mathbf{z}|P, T) d\mathbf{z} = \int p(\gamma|P, T, \mathbf{z}) p(\mathbf{z}|P, T) d\mathbf{z} \quad (6)$$

175 where $p(\mathbf{z}|P, T)$ is a prior distribution on the latent variables. The above hierarchical mathematical
 176 ansatz, despite being very elegant and rigorous, has to be approximated [44], where a regularized
 177 adversarial inference framework is proposed and detailed. The final result is a generator model
 178 $\gamma = f_\phi(p, T, \mathbf{z})$ parametrized by vector ϕ , like trained deep neural networks. In conjunction
 179 with $p(\mathbf{z})$, the statistics of γ can be characterized. More specifically, we can compute its low

R. S. M. Freitas, Á. P. F. Lima, C. Chen, F. A. Rochinha, D. Mira & X. Jiang / Preprint submitted to
 180 order statistics via Monte Carlo sampling. It is important to remark that the predictions with the
 181 identified probabilistic generator, that, in present context, plays the role of a proxy for obtaining
 182 macroscopic thermodynamic properties of mixtures for pressures and temperatures not contained
 183 in \mathcal{D} , is negligible when compared to MD simulations. The mean and variance of the predictive
 184 distribution at a new point (p^*, T^*) are computed as

$$\mu_\gamma(P^*, T^*) = \mathbb{E}[\gamma | P^*, T^*, \mathbf{z}] \approx \frac{1}{N_s} \sum_{i=1}^{N_s} [f_\phi(P^*, T^*, \mathbf{z}_i)] \quad (7)$$

$$\sigma_\gamma^2(P^*, T^*) = \text{Var}[\gamma | P^*, T^*, \mathbf{z}] \approx \frac{1}{N_s} \sum_{i=1}^{N_s} [f_\phi(P^*, T^*, \mathbf{z}_i) - \mu_\gamma(P^*, T^*)]^2, \quad (8)$$

185 where $\mathbf{z}_i \sim p(\mathbf{z})$, $i = 1, \dots, N_s$, and N_s corresponds to the total number of samples.

187 At this point, it is important to clarify that the predictive uncertainty encoded in \mathbf{z} is due to
 188 noise in the Molecular Dynamics computations originated by numerical approximations and to the
 189 potential small amount of data employed in the training process. Therefore, it encapsulates aleatoric
 190 and epistemic uncertainties.

191 Later, we explore the versatility of the probabilistic ML model employing the fusion of data
 192 produced by MD with experimental data obtained for supercritical behavior of the mixture.

193 2.3. Physicochemical properties prediction in EMD simulation

194 In this study, all MD simulations are performed in Gromacs package [60] with Transferable
 195 Potentials for Phase Equilibria (TraPPE) force field [61]. United-atom molecular description is used
 196 in order to reduce the computational cost. Before simulation, 1000 molecules are distributed in a
 197 box with relatively large edge length of 14 nm to avoid atom's overlap. After energy minimisation,
 198 a 2 ns simulation is performed with time setup of 1fs in isobaric-isothermal NPT (fix the number
 199 of atoms, pressure and temperature of the system) ensemble by using Parrinello-Rahman method
 200 [62] to maintain the pressure. Then 1ns NVT (fix the number of atoms, volume and temperature
 201 of the system) simulation is followed for production run. The temperature is controlled by velocity

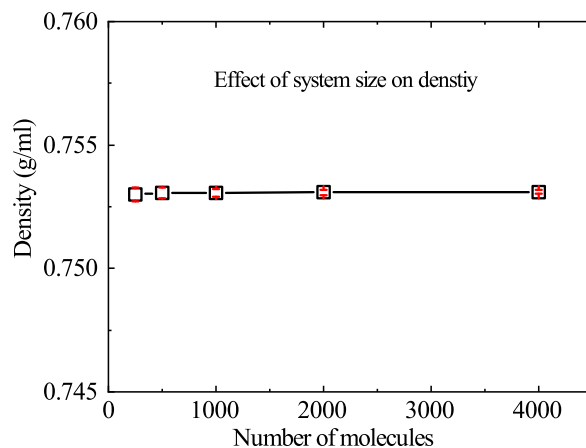


Figure 1: Effect of the system size on density prediction.

202 rescale. The fixed bond length in TraPPE force field is achieved by using LINCS algorithm [63].

203 The density and diffusion is calculated in NVT simulation.

204 The diffusion coefficient (D) can be obtained from the linear fittings of mean square displace-
 205 ment (MSD) of molecules:

$$MSD(t) = \langle |\mathbf{r}_i(t) - \mathbf{r}_i(0)|^2 \rangle \quad (9)$$

$$D(t) = \frac{1}{6} \frac{d}{dt} \langle |\mathbf{r}_i(t) - \mathbf{r}_i(0)|^2 \rangle \quad (10)$$

206 where $\mathbf{r}_i(t)$ is the position of the i^{th} particle at time t , angle bracket indicates the ensemble average
 207 over all the particles in the system.

208 The number of fuel molecules and simulation time in our simulation is setup according to
 209 previous studies. For example, Yang et al. [64] used 250 molecules with 2ns simulation time in
 210 transport property prediction of n-alkanes, and Kondratyuk et al. [65] used 1000 molecules in
 211 modelling branched alkanes running in EMD simulation of 1 ns. Figure 1 depicts the effect of the
 212 system size on the n-dodecane density prediction. As we can see 1000 molecules are sufficient to
 213 achieve convergence of the density prediction at an affordable computational cost.

214 **3. Results and discussion**

215 Here, we demonstrate the performance of the proposed methodology. Despite alternative
216 fuels can be very complex mixtures consisting of hundreds of compounds, we consider single-
217 component alkanes C_nH_{2n+2} , so reliable data for model assessment and validation can be used.
218 In general, realistic fuels are usually described by surrogate models [8] because of availability of
219 validated chemical mechanisms and experimental measurements. The data to train our ML models
220 consist of properties of a family of alkanes, ranging from normal to supercritical conditions. More
221 specifically, we construct ML models to characterize density dependency on some operational
222 conditions in which data is not available. As mentioned before, in order to take into consideration
223 unavoidable uncertainties, we approximate the conditional probability $p(\gamma|\mathbf{x}, \theta)$, with \mathbf{x} being the
224 input vector with components pressure p , temperature T and chemical composition. Moreover,
225 it is worth mentioning here that for simplicity we consider as the input that characterizes the
226 chemical compositions the number of atoms of carbon C in the molecule of the pure compounds, a
227 categorical variable. However, parameters from the EMD used to characterize the physicochemical
228 properties of the fuel molecule can be used. Also, for the GP learning model, the hyper-parameters
229 vector reduces to $\theta = \{\mathbf{l}, \sigma\}$, and for the generative model θ represents the vector of parameters of
230 the deep neural networks ϕ . The latent variable z is embedded in the input vector \mathbf{x} . We employ a
231 one-dimensional latent space with a standard normal prior, $p(z) \sim \mathcal{N}(0, 1)$.

232 The pure compounds considered are n-octane, n-nonane, n-decane, n-dodecane, and n-
233 hexadecane, operating from high-pressure nozzle to supercritical chamber environment conditions.
234 The dataset used to build the ML models consists of 1200 density values. Specifically, there are
235 240 values of the density for each compound, computed at a regular temperature grid within $T \in$
236 $[320, 900]$ K, varying by 20K, and at the specific pressures values: $P = \{3, 4, 6, 8, 10, 20, 100, 150\}$
237 MPa. It is worth remarking that in this dataset we included density values for supercritical regions,
238 more specifically values above the critical temperature (T_c) of the compounds, being the critical
239 values for n-octane ($T_c = 569.32K$), n-nonane ($T_c = 594.55K$), n-decane ($T_c = 617.7K$), n-
240 dodecane ($T_c = 658.1K$), and n-hexadecane ($T_c = 722K$), which replicate engine-like conditions

241 .

242 In the learning process, 80% of the data points are selected randomly to training the ML models.
 243 The remaining 20% are used to validating them. Moreover, the training data set is organized in three
 244 subsets with 10%, 50%, and 100% of data available to train the models. The aim here is to evaluate
 245 the convergence and impacts of constructing the ML models in a small data regime. Accuracy is
 246 measured using the distance between the expected values predicted with the ML models and the
 247 predictions computed with the MD simulations. We check this accuracy computing the L_2 mean
 248 relative error (L_{2-MRE})

$$L_{2-MRE} = \frac{1}{N} \sum_{i=1}^N \left(\frac{\rho_i - \hat{\rho}_i}{\rho_i} \right)^2 \quad (11)$$

249 where ρ_i is the density computed with MD simulations, $\hat{\rho}_i$ is the expected ML output and N is the
 250 number of test samples. Also, we compute the coefficient of determination (R^2 -score) metric [66]

$$R^2 = 1 - \frac{\sum_{i=1}^N \|\rho_i - \hat{\rho}_i\|_2^2}{\sum_{i=1}^N \|\rho_i - \bar{\rho}\|_2^2} \quad (12)$$

251 where $\bar{\rho} = \frac{1}{N} \sum_{i=1}^N \rho_i$ is the mean density of test samples. The R^2 -score metric represents the
 252 normalized error, allowing the comparison between ML models trained by different data sets, with
 253 values close to 1 corresponding to the ML models best accuracy, while L_{2-MRE} is a common metric
 254 used to check the accuracy of ML models during the optimization process.

255 We obtain the GP regression model of Eq. (1) via maximizing the marginal log-likelihood
 256 of Eq. (3) using the Mattern 3/2 kernel function, as that shown in Eq. (2). Also, we have used
 257 the gradient descend optimizer L-BFGS [67] using randomized restarts to ensure convergence
 258 to a global optimum. The GP learning model was implemented in GPy: Gaussian Process (GP)
 259 framework written in python [68].

260 On the other hand, to construct the generative learning model, we departed from the architecture
 261 proposed and validated by Yang and Perdikaris [44]. More specifically, the conditional generative
 262 model is constructed using fully connected feed-forward architectures for the encoder and generator
 263 networks with 4 hidden layers and 100 neurons per layer, while the discriminator architecture

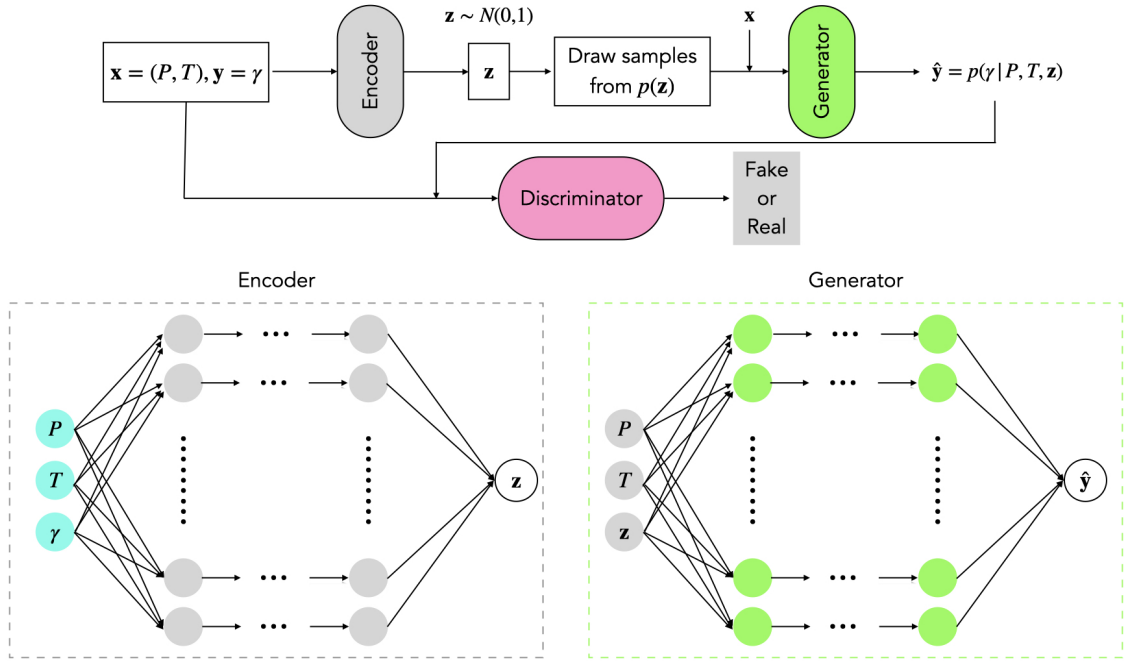


Figure 2: Schematic view of the conditional generative model.

264 has 2 hidden layers with 100 neurons per layer. A schematic view of the conditional generative
 265 model is depicted in Figure 2. The neural networks are constructed by combining try-and-error
 266 and Hyperopt algorithm [69] to search for the hyperparameters that give the lowest L_{2-MRE} . All
 267 activation uses a hyperbolic tangent non-linearity. The models are trained for 50,000 stochastic
 268 gradient descent steps using the Adam optimizer [70] with a learning rate of 10^{-4} , while fixing a
 269 two-to-one ratio for the discriminator versus generator updates. Furthermore, we have also fixed the
 270 entropy regularization and the residual penalty parameters to $\lambda = 1.5$ and $\beta = 0.5$, respectively. The
 271 proposed model was implemented in TensorFlow v2.1.0 [71], and computations were performed
 272 in single precision arithmetic on a single NVIDIA GeForce RTX 2060 GPU card.

273 We also explore some alternatives versions of the above described ML models by proposing
 274 fusion with experimental data and the use of multi-fidelity formulations.

275 3.1. ML results for typical fuel surrogates

Table 1

Gaussian Process training accuracy.

Train data	L_{2-MRE}	R ² -score
10 %	6.2805×10^{-2}	0.8538
50 %	4.7438×10^{-2}	0.9976
100 %	2.7272×10^{-2}	0.9991

Table 2

Generative model training accuracy.

Train data	L_{2-MRE}	R ² -score
10 %	4.9316×10^{-2}	0.9359
50 %	2.8989×10^{-3}	0.9983
100 %	2.1409×10^{-3}	0.9990

276 Tables 1 and 2 show the coefficient of determination (R²-score) and L_2 mean relative error,
 277 respectively, for GP and the probabilistic conditional generative models. The accuracy metrics are
 278 computed with the test samples. We observe that they are not satisfactory in the small training data
 279 scenario, with 10% of training data. R²-scores for the GP and conditional generative models in
 280 this specific training scenario are 0.8538 and 0.9359, respectively. For a data richer situation, with
 281 50% of training data, we observe that the models return good predictions with R²-score higher than
 282 0.99. Also, we observe that the conditional generative model returns better predictions than the
 283 GP model in a small data scenario, with an accuracy of $L_{2-MRE} = 2.8989 \times 10^{-3}$ while the GP
 284 accuracy is $L_{2-MRE} = 4.7428 \times 10^{-2}$. Finally, with 100% of the training data, we can see that the
 285 surrogate models return excellent predictions with R²-score very near 1.0 and mean relative errors
 286 lower than 0.03%.

287 As a further illustration of the performance of such approaches to predict the density, we plot its
 288 values for n-octane, n-dodecane, and n-hexadecane densities with respect to temperature for the ML
 289 models trained with 50% of the dataset, since this training scenario returns the best relation between
 290 accuracy and computational cost. Figure 3 shows the n-octane density predictions at the pressures
 291 equal to 3, 10, and 100 MPa. We can observe that at 3MPa the GP model fails to deliver good
 292 results around the transcritical region, while the generative model provides robust predictions with

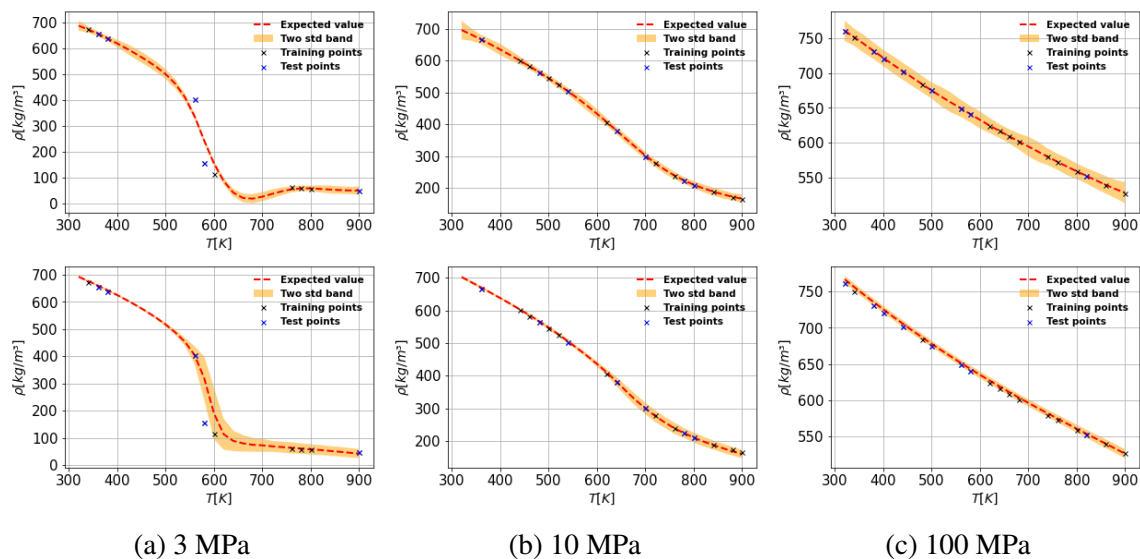


Figure 3: n-Octane predictions with the GP (top) and probabilistic conditional generative models (bottom) at the pressures 3, 10, and 100 MPa.

293 uncertainties bounds that capture the data. The predictive uncertainty of the proposed approaches
 294 reflects limited data for training the models, the epistemic uncertainty. We can also note that both
 295 models perform well at 10 and 100 MPa, wherein the density dependency on the temperature has
 296 a smooth behavior.

297 Also, the n-dodecane and n-hexadecane densities are depicted along with temperature in
 298 Figures 4 and 5. We observe that the ML models return robust predictions at three different
 299 pressures. Besides, it is noted that the GP model returns larger uncertainty bounds at high pressures,
 300 specifically at density points not used in the training process.

301 We also validate how the proposed ML technology perform in an extrapolation scenario. We
 302 validate them for the n-heptane, a fuel not used for building the models. In order to do that, instead
 303 of employing data provided by ML computations, we use an experimental database furnished by the
 304 National Institute of Standards and Technology (NIST). Figure 6 shows that at 3 MPa and liquid
 305 condition the ML model returns good predictions of the n-heptane density behavior, with small
 306 uncertainties. However, at supercritical conditions ($T_c = 540.13K$), the GP model returns density
 307 predictions far from satisfactory. Also, we note that the generative model has uncertainty bounds

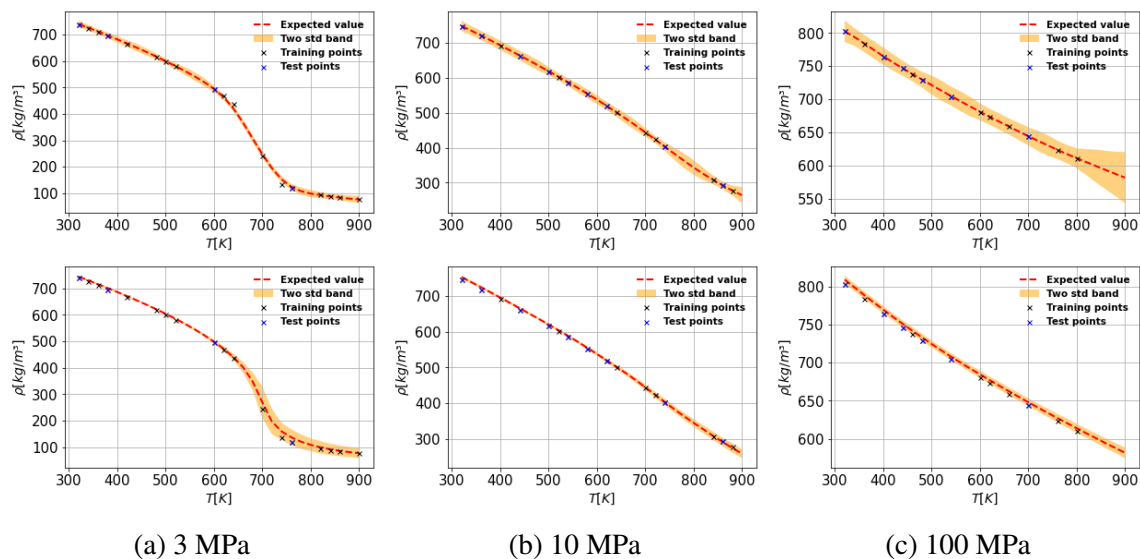


Figure 4: n-Dodecane predictions with the GP machine learning model (top) and conditional generative machine learning model (bottom) at the pressures 3, 10, and 100 MPa.

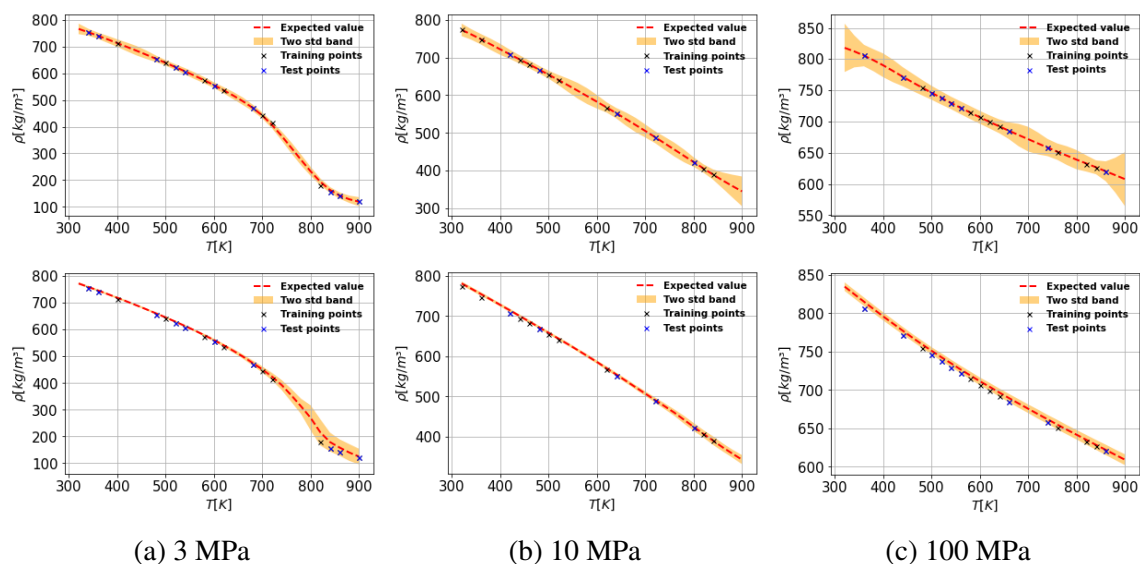


Figure 5: n-Hexadecane predictions with the GP machine learning model (top) and conditional generative machine learning model (bottom) at the pressures 3, 10, and 100 MPa.

308 able to capture the thermophysical property. The L_2 mean relative error between the NIST dataset
 309 and the expected values predicted by the GP and conditional generative models are 7.1697×10^{-2}
 310 and 2.0838×10^{-2} , respectively. We can also note that at higher pressure where the density behavior
 311 is smooth, the models present better predictions, with the GP model showing larger uncertainties
 312 bounds and the generative model returns smaller uncertainty bounds. Moreover, the L_2 mean

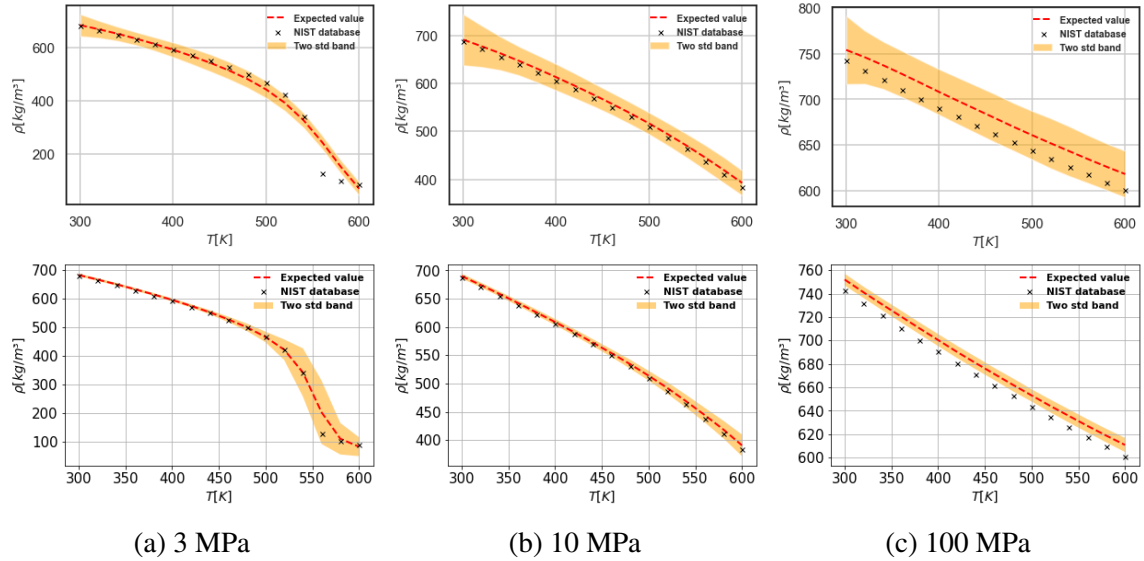


Figure 6: n-Heptane predictions with the GP machine learning model (top) and conditional generative machine learning model (bottom) at the pressures 3, 10, and 100 MPa.

313 relative errors of the GP model at 10 and 100 MPa are respectively 1.8152×10^{-4} and 6.3072×10^{-4} ,
314 and for the conditional generative model the L_2 mean relative errors at the same pressures are
315 8.4484×10^{-5} and 2.0322×10^{-4} .

316 Furthermore, we use a coefficient of variation to measure the degree of uncertainty of the
317 density predictions. It is defined as the ratio between the standard deviation σ_ρ and the mean μ_ρ of
318 the prediction

$$cv(p, T) = \frac{\sigma_\rho(p, T)}{\mu_\rho(p, T)} \quad (13)$$

319 Figure 7 gives an overall picture by displaying a mapping between the operating conditions
320 and the uncertainty on n-octane density predictions. We present an explicit quantification of the
321 epistemic uncertainty resulting from the lack of data, which helps to understand limits of the ML
322 models. More specifically, to make more accessible the visualization of the results, we plot this
323 mapping for $\log_{10} p \in [0.5, 2.5]$ MPa and $T \in [320, 900]$ K with regular intervals of 20K, allowing
324 us to make explicit the strong dependence of the epistemic uncertainties regarding different regions
325 of operating conditions. A critical aspect to be remarked is the higher values of cv in particular

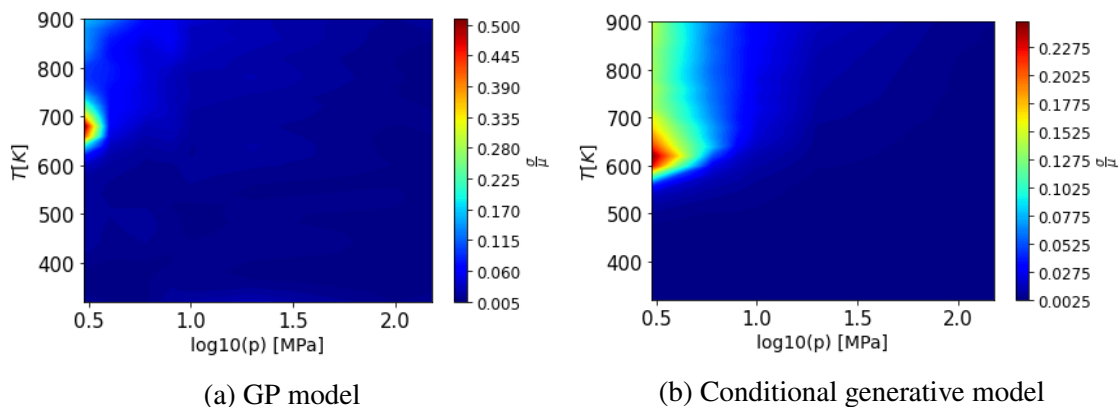


Figure 7: n-Octane density variability for a range of temperatures and pressures.

326 regions of the operating conditions space, especially at transcritical conditions displaying higher
 327 gradients of the property. We can note that the GP model returns a degree of uncertainty slightly
 328 large in this region. That can be mitigated by providing more training data for this specific region.
 329 Also, it is noted that variability of density provided by the conditional generative model is less
 330 pronounced at liquid regions and for high-pressure supercritical regions, which is due to the smooth
 331 density behavior resulting in a low degree of uncertainty in the predictions at these regions.

332 In addition, we explore the ability of ML models considering other physicochemical properties.
 333 Specifically, we extend the above approaches to predict the diffusion coefficient of the alkane
 334 compounds. The diffusion coefficient controls mass transport in combustion engines. Therefore,
 335 understanding diffusion is extremely important in order to optimize industrial processes and
 336 improve device efficiency, especially for supercritical combustion, where the physicochemical
 337 properties of fluids are quite different from those in liquid conditions. It is worth emphasizing that
 338 constructing accurate and simple predictive models overcoming costly simulations and expensive
 339 experimental procedures is crucial for describing physicochemical properties over a wide range of
 340 flow conditions.

341 The dataset used to build the ML models consists of 1240 values of the diffusion coefficient,
 342 computed within a regular temperature grid $T \in [300, 900]$ K, varying by 20K, and at specific
 343 pressures: $P = \{1, 2, 4, 10, 20, 40, 100, 150\}$ MPa. In the training process, 70% of the data points

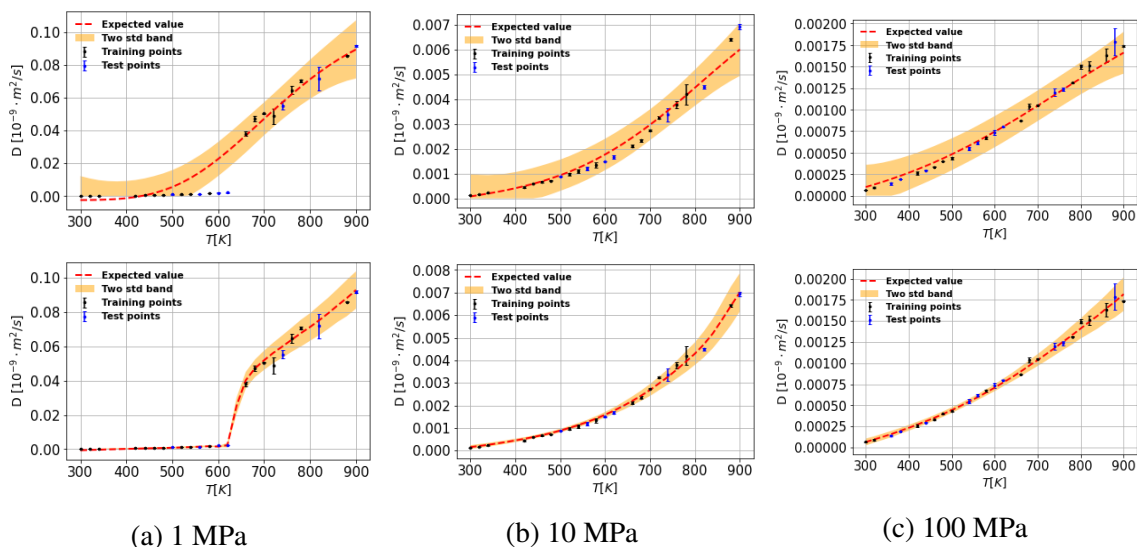


Figure 8: n-Dodecane predictions of the diffusion coefficient with the GP machine learning model (top) and conditional generative machine learning model (bottom) at the pressures 1, 10, and 100 MPa.

344 are selected randomly to training the ML models. The remaining 30% are used to testing them.
 345 Moreover, the training data set is organized in three subsets with 20%, 50%, and 70% of data.
 346 Figure 8 shows the n-dodecane diffusion coefficient predictions at the pressures equal to 1, 10, and
 347 100 MPa for the ML models trained with 50% of the dataset. We observe that the ML models
 348 return robust predictions at three different pressures with GP model returns larger uncertainty
 349 bounds. We can also note that similar to density the model perform better at higher pressures,
 350 wherein the diffusion coefficient dependency on the temperature has a smooth behavior. That
 351 is further confirmed by calculating the L_2 mean relative error, where for a pressure of 1 MPa
 352 the models return worse predictions, as shown in Figure 9. That might be explained by the fact
 353 that the physicochemical properties display higher gradients near transcritical regions at lower
 354 pressures, which decreases the predictability of the models under these conditions. Also, we note
 355 that the generative model has slightly better predictions than the GP model. These results show
 356 the robustness of the proposed approaches to construct predictive models for physicochemical
 357 properties of diesel fuels.

358 3.2. Data-fusion Machine Learning models

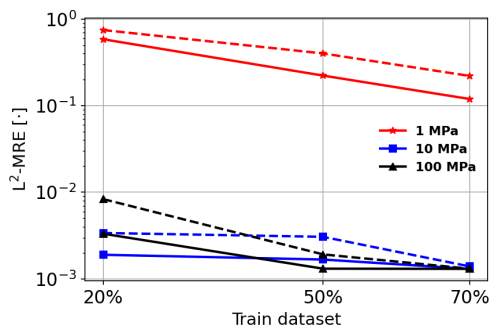


Figure 9: Comparison of the L_2 mean relative error in different data regimes for training. Gaussian process (dashed-line) and generative model (solid-line).

359 Although MD simulation is considered to be a robust tool to predict thermodynamic properties,
 360 it returns unsatisfactory values at critical points/transcritical regions. It was shown [8] that the
 361 transport properties predictions of diesel surrogate fuels are far from satisfactory near such critical
 362 points. That is also the case with n-dodecane in that study. The results depict that EMD simulation
 363 might be unsuitable for predicting the properties at regions near the critical point. Non-equilibrium
 364 molecular dynamics simulation may leverage the results near the critical points, which is beyond the
 365 scope of the present study. Density predictions with MD simulations and NIST data at transcritical
 366 regions present considerable discrepancies, as shown in Figure 10. More specifically, in operating
 367 conditions near the critical point of n-dodecane, critical pressure ($P_c = 1.8170$ MPa) and critical
 368 temperature ($T_c = 658.1$ K), our ML models based on the MD data fail to accurately predict the
 369 density. Figures 11 (a) and 12 (a) show the density predictions at 2 MPa for GP and conditional
 370 generative model, respectively. Moreover, Figures 11 (b) and 12 (b) also show that the main
 371 discrepancies between the expected values of ML models against the NIST database are into the
 372 transcritical regions.

373 Aiming at improving the predictability of our ML models at transcritical regions, we adopt two
 374 strategies, exploring the fusion of MD simulations with experimental data. The aim here is not to
 375 compare these different strategies but to evaluate their potential. Both are formulated with the same
 376 idea, promoting the fusion of data from MD simulations and experiments datasets. In the first one,
 377 we propose a data-fusion strategy in which density points of the transcritical region provided by the

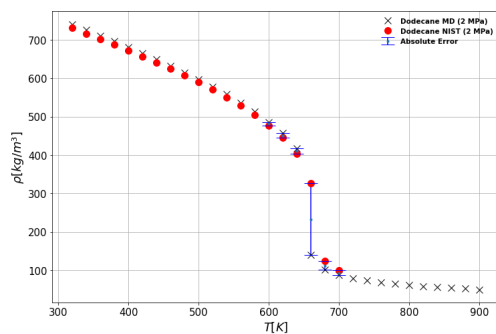
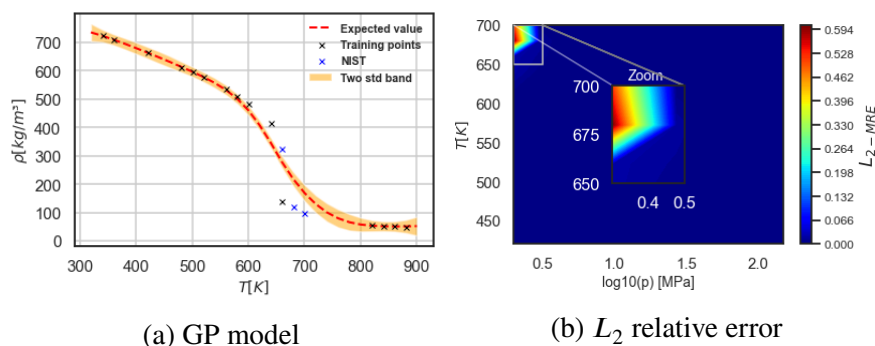


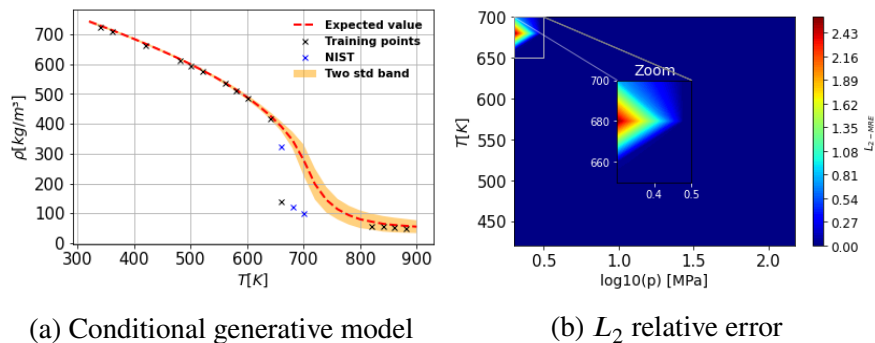
Figure 10: Comparison between n-dodecane density predictions along with temperature at 2 MPa between MD simulations against NIST dataset.



(a) GP model

(b) L_2 relative error

Figure 11: n-Dodecane density predictions GP model: (a) n-Dodecane density along with temperature at 2 MPa. (b) L_2 error between the expected value predicted by the ML model against NIST.



(a) Conditional generative model

(b) L_2 relative error

Figure 12: n-Dodecane density predictions conditional generative model: (a) n-Dodecane density along with temperature at 2 MPa. (b) L_2 error between the expected value predicted by the ML model against NIST.

378 NIST database are simply concatenated into the training dataset. The second differs as we propose
 379 a multi-fidelity arrangement of the data. A detailed description of both strategies is given further
 380 ahead.

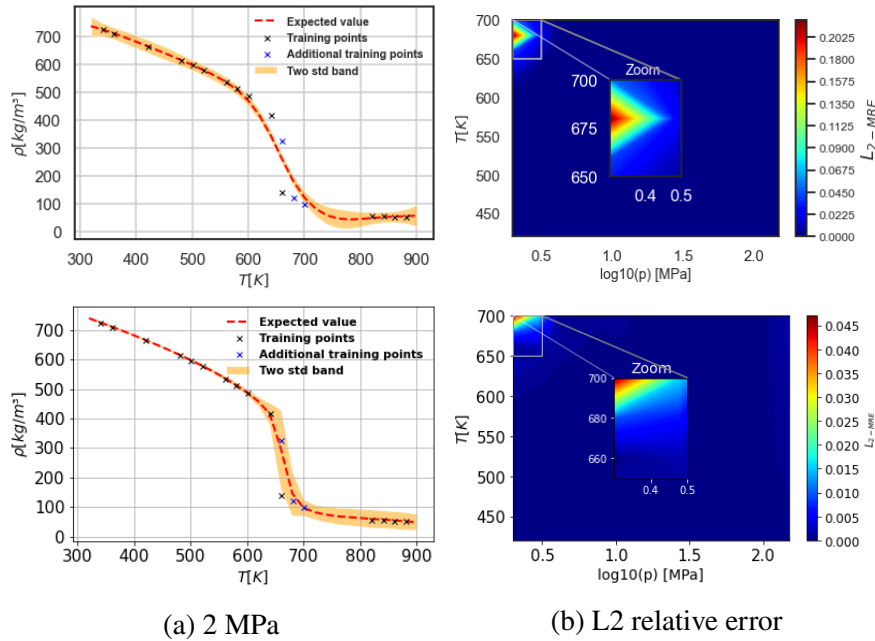


Figure 13: n-Dodecane density predictions with the GP model (top) and conditional generative model (bottom) using the data-fusion approach with three density points from the NIST database.

381 In the data-fusion approach, we add three density values from NIST to the original training
 382 dataset, as depicts in Fig 13 (a). Note that the fusion improves considerably the predictions of the
 383 conditional generative model with relative errors lower than 5%, while the GP model still returns
 384 relative errors not satisfactory. Further details about this data-fusion approach can be found in the
 385 Appendix A.

386 As discussed above, generating reliable data with MD simulations to be used in supervised
 387 learning might require a great computational effort. To tackle such a drawback, numerical
 388 formulations combining models displaying different levels of fidelity are frequently employed.
 389 Those multi-fidelity simulators employ, for instance, coarse grid discretizations, models based on
 390 simplified physics, or simplified iterative methods. Here, again we merge experimental data with
 391 MD simulations, restricting our approaches to two levels of fidelity.

392 In this new context, we propose extensions of the previous introduced ML models. We start
 393 by obtaining high-fidelity $\{\mathbf{x}_H, \gamma_H\}$ and low-fidelity $\{\mathbf{x}_L, \gamma_L\}$ input-output samples. Typically, the
 394 number of samples in the first case tends to be much smaller due to the related costs. We assign

R. S. M. Freitas, Á. P. F. Lima, C. Chen, F. A. Rochinha, D. Mira & X. Jiang / Preprint submitted to
395 the high-fidelity score to the experimental data, according to the considerations above about the
396 potential inaccuracy of the MD obtained computed properties for transcritical regions.

397 We start with our first, in this multi-fidelity context, ML model approximating the conditional
398 probability $p(\gamma_H|\mathbf{x}_H, \gamma_L, z)$, using the generative model $\gamma_H = f_\phi(\mathbf{x}_H, \gamma_L, z)$, $z \sim p(z)$. In another
399 words, the ML model is supposed to capture the correlation between the two level of fidelity
400 data. Once this is achieved, we have a predictive model computing outputs for a new point \mathbf{x}^* :
401 $\mathbf{y}_H^* = f_\phi(\mathbf{x}^*, f_L(\mathbf{x}^*), z)$. At this point, it is worth remarking that one of the inputs is the output of
402 the low-fidelity model, leading to a recursive scheme to obtain the predictions of the multifidelity
403 model. In fact, here the considered low fidelity data is produced with expensive MD simulations.
404 Therefore, in order to achieve a feasible scheme, we need to build an auxiliary, cheap to compute
405 and accurate, proxy for the low fidelity model using the available data.

406 As a second approach, the one based on GPs, we employ the nonlinear autoregressive multi-
407 fidelity GP (NARGP) regression model [53]. The main idea of the NARGP model is to extend
408 GP modeling to capture nonlinear correlations from data generated by sources of different fidelity
409 [72, 73]. It enables the construction of probabilistic models prone to encapsulate uncertainties, built
410 upon the recursive relation $y_H = g(x_H, f_L(x_H))$ involving low and high fidelity data, in which f_L
411 is a GP model for the former. Moreover, we put a GP prior on g . After the training, we obtain the
412 predictive model, which turns to be also a GP, $y_H = g(x^*, f_L(x^*))$.

413 To assess the above multi-fidelity ML approaches, we use an illustrative example involving
414 data from "low-fidelity" MD simulations and "high-fidelity" NIST experimental values. For
415 both approaches, the training dataset consists of 7 density values of n-dodecane $\rho_H(p, T_H)$
416 and $\rho_L(p, T_L)$, at the pressure of 2 MPa and a set of temperatures given by $T_H = T_L =$
417 $\{320, 440, 500, 620, 660, 680, 700\}$ K. Note that we prioritize points located in the transcritical part,
418 since this region presents larger discrepancies between the values predicted by MD simulations and
419 the NIST database.

420 The conditional generative model is constructed using fully connected feed-forward architec-
421 ture for the encoder and generator networks with 4 hidden layers and 100 neurons per layer, while

R. S. M. Freitas, Á. P. F. Lima, C. Chen, F. A. Rochinha, D. Mira & X. Jiang / Preprint submitted to
422 the discriminator architecture has 2 hidden layers with 100 neurons per layer. All activation uses
423 a hyperbolic tangent non-linearity. The models are trained for 20,000 stochastic gradient descent
424 steps using the Adam optimizer [70] with a learning rate of 10^{-4} , while fixing a one-to-five ratio for
425 the discriminator versus generator updates. Furthermore, we have fixed the entropy regularization
426 parameter to $\lambda = 1.5$, and we also employed a one-dimensional latent space with a standard normal
427 prior, $p(z) \sim \mathcal{N}(0, 1)$.

428 We train the NARGP model via maximizing the marginal log-likelihood using the Mattern 3/2
429 kernel function. The gradient descend optimizer L-BFGS is used considering randomized restarts
430 to ensure convergence to a global optimum. Once the high-fidelity recursive GP is trained, we can
431 compute the predictive posterior mean and variance at a given untested point \mathbf{x}^* by sampling the
432 probabilistic predictive model.

433 The main results are summarized in Fig 14. More specifically, the results indicate that
434 the NARGP model was able to satisfactorily reconstruct the high-fidelity data. To make this
435 comparison quantitative, we compute the mean L_2 relative error between the expected values
436 predicted by the generative model and the NIST data. It shows predictions with accuracy of
437 $L_{2-MRE} = 1.4524 \times 10^{-2}$. Moreover, it returns good uncertainty bounds able to capture the
438 high-fidelity response at the transcritical region. Also, we note a perfect agreement between the
439 expected value provided by the probabilistic conditional generative model and the high-fidelity
440 data, resulting in an accuracy of $L_{2-MRE} = 4.4782 \times 10^{-5}$. Finally, we observe that the multi-
441 fidelity model returns small uncertainty bounds despite the small amount of data employed in the
442 training process.

443 4. Conclusions

444 In this work, we propose a computational methodology based on the use of ML with Molecular
445 Dynamics simulations to compute physicochemical properties of single compound fuels at engine-
446 relevant conditions. The ML models have been revealed to be a powerful tool to predict accurately
447 the fuel properties of pure compounds. Moreover, this study explores the versatility of the ML

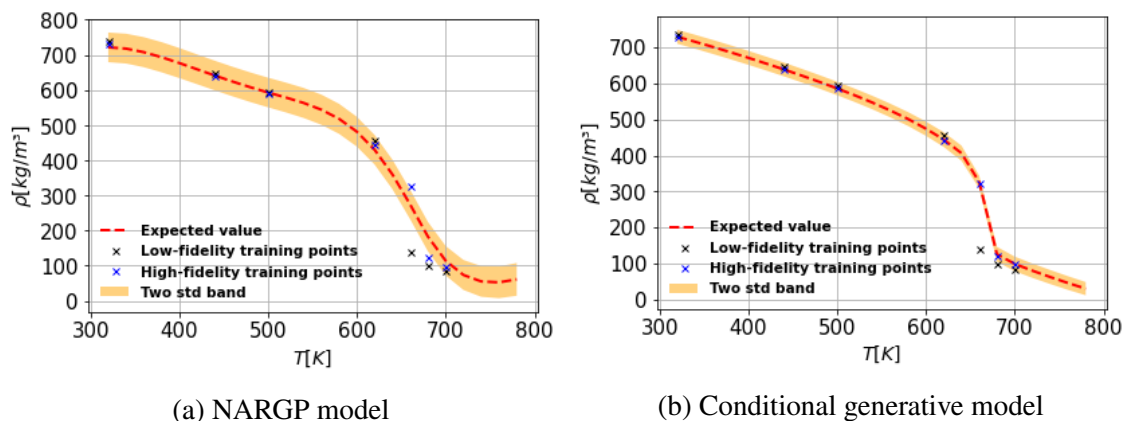


Figure 14: Multi-fidelity modeling of n-dodecane diesel surrogate fuel density.

448 models to handle data from different sources, which can then be integrated efficiently in the context
 449 of UQ workflows with many-query tasks.

450 We place our contribution in the emerging area of physics-aware ML, where the final model,
 451 in many different ways, blends two main components: availability of experimental data and/or
 452 often expensive computational models relying on first principles and phenomenological closure
 453 equations, and deep learning data-driven models. Such combination allows describing physico-
 454 chemical properties over a wide range of flow conditions at relatively low cost, and also offers a
 455 broad spectrum of opportunities to enhance CFD codes.

456 This study has shown a successful prediction of fuel physical quantities, in this case density
 457 and diffusion coefficient, that can also be extended to other physicochemical properties as well as
 458 more complex fuel molecules or multicomponent mixtures like dimethyl ethers or oxymethylene
 459 dimethyl ethers. The generation of reliable physicochemical properties of renewable fuels is an
 460 important step forward towards the generation of digital tools that can assist on the decarbonization
 461 by the use of renewable fuels.

462 **Acknowledgements.** The research leading to these results had received funding from the Brazilian
 463 National Agency for Petroleum, Natural Gas and Biofuels (ANP) through Programa de Recursos
 464 Humanos (PRH) under the PRH 8 - Mechanical Engineering for the Efficient Use of Biofuels, grant
 465 agreement numbers F0A5.EDDE.B5C0.3BCB and 2B61.4F5C.A83B.A713.

References

- 466
- 467 [1] S. Naimoli, S. Ladislav, Climate solutions series: Decarbonizing heavy industry, Tech. rep., Center for Strategic and International Studies
468 (CSIS) (2020).
469 URL <http://www.jstor.org/stable/resrep26402>
- 470 [2] H. Mandová, T. Vass, A. F. Pales, P. Levi, T. Gül, The challenge of reaching zero emissions in heavy industry, Tech. rep., International Energy
471 Agency (IEA) (2020).
472 URL <https://www.iea.org/articles/the-challenge-of-reaching-zero-emissions-in-heavy-industry>
- 473 [3] A. Omari, B. Heuser, S. Pischinger, Potential of oxymethylenether-diesel blends for ultra-low emission engines, *Fuel* 209 (July) (2017) 232–
474 237. doi:10.1016/j.fuel.2017.07.107.
475 URL <http://dx.doi.org/10.1016/j.fuel.2017.07.107>
- 476 [4] D. Pélerin, K. Gaukel, M. Härtl, E. Jacob, G. Wachtmeister, Potentials to simplify the engine system using the alternative diesel fuels
477 oxymethylene ether OME1 and OME3&LŠ-6 on a heavy-duty engine, *Fuel* 259 (2020) 116231. doi:10.1016/j.fuel.2019.116231.
478 URL <https://linkinghub.elsevier.com/retrieve/pii/S0016236119315856>
- 479 [5] J. V. Pastor, J. M. García-Oliver, C. Micó, A. A. García-Carrero, A. Gómez, Experimental study of the effect of hydrotreated vegetable oil and
480 oxymethylene ethers on main spray and combustion characteristics under engine combustion network spray a conditions, *Applied Sciences*
481 10 (16) (2020) 5460.
- 482 [6] W. J. Pitz, C. J. Mueller, Recent progress in the development of diesel surrogate fuels, *Progress in Energy and Combustion Science* 37 (3)
483 (2011) 330–350. doi:<https://doi.org/10.1016/j.pecs.2010.06.004>.
484 URL <https://www.sciencedirect.com/science/article/pii/S0360128510000535>
- 485 [7] J. Y. Lai, K. C. Lin, A. Violi, Biodiesel combustion: Advances in chemical kinetic modeling, *Progress in Energy and Combustion Science*
486 37 (1) (2011) 1–14. doi:<https://doi.org/10.1016/j.pecs.2010.03.001>.
487 URL <https://www.sciencedirect.com/science/article/pii/S036012851000033X>
- 488 [8] C. Chen, X. Jiang, Transport property prediction and inhomogeneity analysis of supercritical n-dodecane by molecular dynamics simulation,
489 *Fuel* 244 (2019) 48–60. doi:<https://doi.org/10.1016/j.fuel.2019.01.181>.
490 URL <https://www.sciencedirect.com/science/article/pii/S0016236119301826>
- 491 [9] X. Yang, M. Zhang, Y. Gao, J. Cui, B. Cao, Molecular dynamics study on viscosities of sub/supercritical n-decane, n-undecane and n-dodecane,
492 *Journal of Molecular Liquids* 335 (2021) 116180.
- 493 [10] X. Yang, Y. Gao, M. Zhang, W. Jiang, B. Cao, Comparison of atomic simulation methods for computing thermal conductivity of n-decane at
494 sub/supercritical pressure, *Journal of Molecular Liquids* (2021) 117478.
- 495 [11] N. D. Kondratyuk, V. V. Pisarev, J. P. Ewen, Probing the high-pressure viscosity of hydrocarbon mixtures using molecular dynamics
496 simulations, *The Journal of Chemical Physics* 153 (15) (2020) 154502.
- 497 [12] N. Kondratyuk, D. Lenev, V. Pisarev, Transport coefficients of model lubricants up to 400 mpa from molecular dynamics, *The Journal of*
498 *Chemical Physics* 152 (19) (2020) 191104.
- 499 [13] N. D. Kondratyuk, V. V. Pisarev, Calculation of viscosities of branched alkanes from 0.1 to 1000 mpa by molecular dynamics methods using
500 compass force field, *Fluid Phase Equilibria* 498 (2019) 151–159.
- 501 [14] C. Caleman, P. J. Van Maaren, M. Hong, J. S. Hub, L. T. Costa, D. Van Der Spoel, Force field benchmark of organic liquids: density, enthalpy
502 of vaporization, heat capacities, surface tension, isothermal compressibility, volumetric expansion coefficient, and dielectric constant, *Journal*
503 *of chemical theory and computation* 8 (1) (2012) 61–74.

- 504 [15] R. S. M. Freitas, F. A. Rochinha, D. Mira, X. Jiang, Parametric and model uncertainties induced by reduced order chemical mechanisms for
505 biogas combustion, *Chemical Engineering Science* 227 (2020) 115949. doi:<https://doi.org/10.1016/j.ces.2020.115949>.
506 URL <https://www.sciencedirect.com/science/article/pii/S0009250920304814>
- 507 [16] L. Ward, A. Agrawal, A. Choudhary, C. Wolverton, A general-purpose machine learning framework for predicting properties of inorganic
508 materials, *npj Computational Materials* 2 (16028) (2016). doi:<https://doi.org/10.1038/npjcompumats.2016.28>.
- 509 [17] A. Ramadhas, S. Jayaraj, C. Muraleedharan, K. Padmakumari, Artificial neural networks used for the prediction of the cetane number of
510 biodiesel, *Renewable Energy* 31 (15) (2006) 2524–2533. doi:<https://doi.org/10.1016/j.renene.2006.01.009>.
511 URL <https://www.sciencedirect.com/science/article/pii/S0960148106000395>
- 512 [18] R. Piloto-Rodríguez, Y. Sánchez-Borroto, M. Lapuerta, L. Goyos-Pérez, S. Verhelst, Prediction of the cetane number of biodiesel using
513 artificial neural networks and multiple linear regression, *Energy Conversion and Management* 65 (2013) 255–261, global Conference on
514 Renewable energy and Energy Efficiency for Desert Regions 2011 "GCREEDER 2011". doi:<https://doi.org/10.1016/j.enconman.2012.07.023>.
515 URL <https://www.sciencedirect.com/science/article/pii/S0196890412003093>
- 516 [19] S. M. Miraboutalebi, P. Kazemi, P. Bahrami, Fatty acid methyl ester (fame) composition used for estimation of biodiesel cetane number
517 employing random forest and artificial neural networks: A new approach, *Fuel* 166 (2016) 143–151. doi:<https://doi.org/10.1016/j.fuel.2015.10.118>.
518 URL <https://www.sciencedirect.com/science/article/pii/S0016236115011321>
- 519 [20] S. Faizollahzadeh Ardabili, B. Najafi, S. Shamshirband, Fuzzy logic method for the prediction of cetane number using carbon number,
520 double bounds, iodine, and saponification values of biodiesel fuels, *Environmental Progress & Sustainable Energy* 38 (2) (2019) 584–599.
521 arXiv:<https://aiche.onlinelibrary.wiley.com/doi/pdf/10.1002/ep.12960>, doi:<https://doi.org/10.1002/ep.12960>.
522 URL <https://aiche.onlinelibrary.wiley.com/doi/abs/10.1002/ep.12960>
- 523 [21] S. Tipler, G. D'Álessio, Q. Van Haute, A. Parente, F. Contino, A. Coussement, Predicting octane numbers relying on principal component
524 analysis and artificial neural network, *Computers & Chemical Engineering* 161 (2022) 107784. doi:<https://doi.org/10.1016/j.compchemeng.2022.107784>.
525 URL <https://www.sciencedirect.com/science/article/pii/S0098135422001259>
- 526 [22] M. Mostafaei, Anfis models for prediction of biodiesel fuels cetane number using desirability function, *Fuel* 216 (2018) 665–672. doi:
527 <https://doi.org/10.1016/j.fuel.2017.12.025>.
528 URL <https://www.sciencedirect.com/science/article/pii/S0016236117315958>
- 529 [23] A. Bemani, Q. Xiong, A. Baghban, S. Habibzadeh, A. H. Mohammadi, M. H. Doranehgard, Modeling of cetane number of biodiesel
530 from fatty acid methyl ester (fame) information using ga-, pso-, and hgapso- lssvm models, *Renewable Energy* 150 (2020) 924–934.
531 doi:<https://doi.org/10.1016/j.renene.2019.12.086>.
532 URL <https://www.sciencedirect.com/science/article/pii/S0960148119319585>
- 533 [24] A. Baghban, M. N. Kardani, A. H. Mohammadi, Improved estimation of cetane number of fatty acid methyl esters (fames) based biodiesels
534 using tlbo-nn and pso-nn models, *Fuel* 232 (2018) 620–631. doi:<https://doi.org/10.1016/j.fuel.2018.05.166>.
535 URL <https://www.sciencedirect.com/science/article/pii/S0016236118310111>
- 536 [25] A. S. Noushabadi, A. Dashti, M. Raji, A. Zarei, A. H. Mohammadi, Estimation of cetane numbers of biodiesel and diesel oils using regression
537 and pso-anfis models, *Renewable Energy* 158 (2020) 465–473. doi:<https://doi.org/10.1016/j.renene.2020.04.146>.
538 URL <https://www.sciencedirect.com/science/article/pii/S0960148120306844>

- 542 [26] Y. SÁnchez-Borroto, R. Piloto-Rodríguez, M. Errasti, R. Sierens, S. Verhelst, Prediction of cetane number and ignition delay of biodiesel
543 using artificial neural networks, *Energy Procedia* 57 (2014) 877–885, 2013 ISES Solar World Congress. doi:[https://doi.org/10.1016/](https://doi.org/10.1016/j.egypro.2014.10.297)
544 [j.egypro.2014.10.297](https://doi.org/10.1016/j.egypro.2014.10.297).
545 URL <https://www.sciencedirect.com/science/article/pii/S1876610214016646>
- 546 [27] W. Yu, F. Zhao, Prediction of critical properties of biodiesel fuels from fames compositions using intelligent genetic algorithm-based back
547 propagation neural network, *Energy Sources, Part A: Recovery, Utilization, and Environmental Effects* 0 (0) (2019) 1–14. arXiv:<https://doi.org/10.1080/15567036.2019.1641575>, doi:10.1080/15567036.2019.1641575.
548 [//doi.org/10.1080/15567036.2019.1641575](https://doi.org/10.1080/15567036.2019.1641575).
549 URL <https://doi.org/10.1080/15567036.2019.1641575>
- 550 [28] D. Alviso, G. Artana, T. Duriez, Prediction of biodiesel physico-chemical properties from its fatty acid composition using genetic
551 programming, *Fuel* 264 (2020) 116844. doi:<https://doi.org/10.1016/j.fuel.2019.116844>.
552 URL <https://www.sciencedirect.com/science/article/pii/S0016236119321982>
- 553 [29] X. Meng, M. Jia, T. Wang, Neural network prediction of biodiesel kinematic viscosity at 313K, *Fuel* 121 (2014) 133–140. doi:<https://doi.org/10.1016/j.fuel.2013.12.029>.
554 [//doi.org/10.1016/j.fuel.2013.12.029](https://doi.org/10.1016/j.fuel.2013.12.029).
555 URL <https://www.sciencedirect.com/science/article/pii/S0016236113011733>
- 556 [30] K. Cheenkachorn, Predicting properties of biodiesels using statistical models and artificial neural networks, 2006.
- 557 [31] S. O. Giwa, S. O. Adekomaya, K. O. Adama, M. O. Mukaila, Prediction of selected biodiesel fuel properties using artificial neural network,
558 *Frontiers in Energy* 9 (2015) 433–445. doi:<https://doi.org/10.1007/s11708-015-0383-5>.
- 559 [32] C. Rocabrúno-ValdÁs, L. RamÁnrez-Verduzco, J. HernÁndez, Artificial neural network models to predict density, dynamic viscosity, and
560 cetane number of biodiesel, *Fuel* 147 (2015) 9–17. doi:<https://doi.org/10.1016/j.fuel.2015.01.024>.
561 URL <https://www.sciencedirect.com/science/article/pii/S0016236115000381>
- 562 [33] F. M. de Oliveira, L. S. de Carvalho, L. S. G. Teixeira, C. H. Fontes, K. M. G. Lima, A. B. F. CÁmara, H. O. M. AraÁzjo, R. V. Sales, Predicting
563 cetane index, flash point, and content sulfur of dieselÁbiodiesel blend using an artificial neural network model, *Energy & Fuels* 31 (4) (2017)
564 3913–3920. arXiv:<https://doi.org/10.1021/acs.energyfuels.7b00282>, doi:10.1021/acs.energyfuels.7b00282.
565 URL <https://doi.org/10.1021/acs.energyfuels.7b00282>
- 566 [34] L. Zhou, Toward prediction of kinematic viscosity of biodiesel using a robust approach, *Energy Sources, Part A: Recovery, Utilization,*
567 *and Environmental Effects* 40 (23) (2018) 2895–2902. arXiv:<https://doi.org/10.1080/15567036.2018.1513099>, doi:10.1080/
568 [15567036.2018.1513099](https://doi.org/10.1080/15567036.2018.1513099).
569 URL <https://doi.org/10.1080/15567036.2018.1513099>
- 570 [35] T. EryÁslmaz, M. Yesilyurt, A. Taner, Á. A. ÁĖelik, Prediction of kinematic viscosities of biodiesels derived from edible and non-edible
571 vegetable oils by using artificial neural networks, *Arabian Journal for Science and Engineering* 40 (2015) 3745–3758.
- 572 [36] T. Eryilmaz, M. Arslan, M. K. Yesilyurt, A. Taner, Comparison of empirical equations and artificial neural network results in terms of kinematic
573 viscosity prediction of fuels based on hazelnut oil methyl ester, *Environmental Progress & Sustainable Energy* 35 (6) (2016) 1827–1841.
574 arXiv:<https://aiche.onlinelibrary.wiley.com/doi/pdf/10.1002/ep.12410>, doi:<https://doi.org/10.1002/ep.12410>.
575 URL <https://aiche.onlinelibrary.wiley.com/doi/abs/10.1002/ep.12410>
- 576 [37] K. Zhu, E. A. Müller, Generating a machine-learned equation of state for fluid properties, *The Journal of Physical Chemistry B* 124 (39)
577 (2020) 8628–8639, pMID: 32870675. arXiv:<https://doi.org/10.1021/acs.jpcc.0c05806>, doi:10.1021/acs.jpcc.0c05806.
578 URL <https://doi.org/10.1021/acs.jpcc.0c05806>

- 579 [38] C. J. Leverant, J. A. Harvey, T. M. Alam, J. A. Greathouse, Machine learning self-diffusion prediction for lennard-jones fluids in pores,
580 The Journal of Physical Chemistry C 125 (46) (2021) 25898–25906. arXiv:<https://doi.org/10.1021/acs.jpcc.1c08297>, doi:
581 10.1021/acs.jpcc.1c08297.
582 URL <https://doi.org/10.1021/acs.jpcc.1c08297>
- 583 [39] J. P. Allers, C. W. Priest, J. A. Greathouse, T. M. Alam, Using computationally-determined properties for machine learning prediction of
584 self-diffusion coefficients in pure liquids, The Journal of Physical Chemistry B 125 (47) (2021) 12990–13002, pMID: 34793167. arXiv:
585 <https://doi.org/10.1021/acs.jpcc.1c07092>, doi:10.1021/acs.jpcc.1c07092.
586 URL <https://doi.org/10.1021/acs.jpcc.1c07092>
- 587 [40] Y. Liu, W. Hong, B. Cao, Machine learning for predicting thermodynamic properties of pure fluids and their mixtures, Energy 188 (2019)
588 116091. doi:<https://doi.org/10.1016/j.energy.2019.116091>.
589 URL <https://www.sciencedirect.com/science/article/pii/S0360544219317864>
- 590 [41] P. Koukouvinis, C. Rodriguez, J. Hwang, I. Karathanassis, M. Gavaises, L. Pickett, Machine learning and transcritical sprays: A demonstration
591 study of their potential in ecn spray-a, International Journal of Engine Research (2021) 14680874211020292 arXiv:<https://doi.org/10.1177/14680874211020292>, doi:10.1177/14680874211020292.
592 URL <https://doi.org/10.1177/14680874211020292>
- 593 [42] C. E. Rasmussen, C. K. I. Williams, Gaussian Processes for Machine Learning, The MIT Press, 2006.
- 594 [43] Y. Yang, P. Perdikaris, Adversarial uncertainty quantification in physics-informed neural networks, Journal of Computational Physics 394
595 (2019) 136–152. doi:<https://doi.org/10.1016/j.jcp.2019.05.027>.
596 URL <https://www.sciencedirect.com/science/article/pii/S0021999119303584>
- 597 [44] Y. Yang, P. Perdikaris, Conditional deep surrogate models for stochastic, high-dimensional, and multi-fidelity systems, Computational
598 Mechanics 64 (2019) 417–434. doi:<https://doi.org/10.1007/s00466-019-01718-y>.
599
- 600 [45] J. Safarov, U. Ashurova, B. Ahmadov, E. Abdullayev, A. Shahverdiyev, E. Hassel, Thermophysical properties of diesel fuel over a wide range
601 of temperatures and pressures, Fuel 216 (2018) 870–889. doi:<https://doi.org/10.1016/j.fuel.2017.11.125>.
602 URL <https://www.sciencedirect.com/science/article/pii/S0016236117315399>
- 603 [46] I. Pioro, S. Mokry, S. Draper, Specifics of thermophysical properties and forced-convective heat transfer at critical and supercritical pressures,
604 Rev. Chem. Eng. 27 (3-4) (2011) 191–214. doi:doi:10.1515/REVCE.2011.501.
605 URL <https://doi.org/10.1515/REVCE.2011.501>
- 606 [47] L. M. Pickett, C. L. Genzale, G. Bruneaux, L.-M. Malbec, L. Hermant, C. Christiansen, J. Schramm, Comparison of diesel spray combustion
607 in different high-temperature, high-pressure facilities (oct 2010). doi:<https://doi.org/10.4271/2010-01-2106>.
608 URL <https://doi.org/10.4271/2010-01-2106>
- 609 [48] Y. Shen, Y.-B. Liu, B.-Y. Cao, C4+ surrogate models for thermophysical properties of aviation kerosene rp-3 at supercritical pressures,
610 Energy & Fuels 35 (9) (2021) 7858–7865. arXiv:<https://doi.org/10.1021/acs.energyfuels.1c00326>, doi:10.1021/acs.
611 energyfuels.1c00326.
612 URL <https://doi.org/10.1021/acs.energyfuels.1c00326>
- 613 [49] M. Razi, A. Narayan, R. Kirby, D. Bedrov, Fast predictive models based on multi-fidelity sampling of properties in molecular dynamics
614 simulations, Computational Materials Science 152 (2018) 125–133. doi:<https://doi.org/10.1016/j.commatsci.2018.05.029>.
615 URL <https://www.sciencedirect.com/science/article/pii/S0927025618303367>

- 616 [50] W. Xing, A. Shah, P. Wang, S. Zhe, Q. Fu, R. Kirby, Residual gaussian process: A tractable nonparametric bayesian emulator for multi-fidelity
617 simulations, *Applied Mathematical Modelling* 97 (2021) 36–56. doi:<https://doi.org/10.1016/j.apm.2021.03.041>.
618 URL <https://www.sciencedirect.com/science/article/pii/S0307904X21001724>
- 619 [51] P. Koukouvinis, A. Vidal-Roncero, C. Rodriguez, M. Gavaises, L. Pickett, High pressure/high temperature multiphase simulations of dodecane
620 injection to nitrogen: Application on ecn spray-a, *Fuel* 275 (2020) 117871. doi:<https://doi.org/10.1016/j.fuel.2020.117871>.
621 URL <https://www.sciencedirect.com/science/article/pii/S001623612030867X>
- 622 [52] V. Alves, V. Gazzaneo, F. V. Lima, A machine learning-based process operability framework using gaussian processes, *Computers & Chemical
623 Engineering* 163 (2022) 107835. doi:<https://doi.org/10.1016/j.compchemeng.2022.107835>.
624 URL <https://www.sciencedirect.com/science/article/pii/S0098135422001739>
- 625 [53] P. Perdikaris, M. Raissi, A. Damianou, N. D. Lawrence, G. E. Karniadakis, Nonlinear information fusion algorithms for data-efficient multi-
626 fidelity modelling, *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences* 473 (2198) (2017) 20160751.
627 arXiv:<https://royalsocietypublishing.org/doi/pdf/10.1098/rspa.2016.0751>, doi:10.1098/rspa.2016.0751.
628 URL <https://royalsocietypublishing.org/doi/abs/10.1098/rspa.2016.0751>
- 629 [54] G. Su, L. Peng, L. Hu, A gaussian process-based dynamic surrogate model for complex engineering structural reliability analysis, *Structural
630 Safety* 68 (2017) 97–109. doi:<https://doi.org/10.1016/j.strusafe.2017.06.003>.
631 URL <https://www.sciencedirect.com/science/article/pii/S016747301730214X>
- 632 [55] T. Chen, J. Morris, E. Martin, Gaussian process regression for multivariate spectroscopic calibration, *Chemometrics and Intelligent Laboratory
633 Systems* 87 (1) (2007) 59–71, selected papers presented at the Conferentia Chemometrica 2005 HajdÅzszoboszlÅs, Hungary 28-31 August
634 2005. doi:<https://doi.org/10.1016/j.chemolab.2006.09.004>.
635 URL <https://www.sciencedirect.com/science/article/pii/S0169743906001900>
- 636 [56] J. Yuan, K. Wang, T. Yu, M. Fang, Reliable multi-objective optimization of high-speed wedm process based on gaussian process regression,
637 *International Journal of Machine Tools and Manufacture* 48 (1) (2008) 47–60. doi:<https://doi.org/10.1016/j.ijmachtools.2007.07.011>.
638
639 URL <https://www.sciencedirect.com/science/article/pii/S0890695507001265>
- 640 [57] G. M. Guerra, R. Freitas, F. A. Rochinha, Constructing accurate phenomenological surrogate for fluid structure interaction models, in:
641 K. L. Cavalca, H. I. Weber (Eds.), *Proceedings of the 10th International Conference on Rotor Dynamics – IFToMM*, Springer International
642 Publishing, Cham, 2019, pp. 295–305.
- 643 [58] D. P. Kingma, M. Welling, Auto-encoding variational bayes (2014). arXiv:1312.6114.
- 644 [59] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial networks
645 (2014). arXiv:1406.2661.
- 646 [60] D. Van Der Spoel, E. Lindahl, B. Hess, G. Groenhof, A. E. Mark, H. J. Berendsen, Gromacs: fast, flexible, and free, *Journal of computational
647 chemistry* 26 (16) (2005) 1701–1718.
- 648 [61] M. G. Martin, J. I. Siepmann, Transferable potentials for phase equilibria. 1. united-atom description of n-alkanes, *The Journal of Physical
649 Chemistry B* 102 (14) (1998) 2569–2577.
- 650 [62] M. Parrinello, A. Rahman, Polymorphic transitions in single crystals: A new molecular dynamics method, *Journal of Applied physics* 52 (12)
651 (1981) 7182–7190.
- 652 [63] B. Hess, H. Bekker, H. J. Berendsen, J. G. Fraaije, Lincs: a linear constraint solver for molecular simulations, *Journal of computational
653 chemistry* 18 (12) (1997) 1463–1472.

- 654 [64] X. Yang, M. Zhang, Y. Gao, J. Cui, B. Cao, Molecular dynamics study on viscosities of sub/supercritical n-decane, n-undecane and n-dodecane,
655 Journal of Molecular Liquids 335 (2021) 116180. doi:<https://doi.org/10.1016/j.molliq.2021.116180>.
656 URL <https://www.sciencedirect.com/science/article/pii/S0167732221009077>
- 657 [65] N. Kondratyuk, D. Lenev, V. Pisarev, Transport coefficients of model lubricants up to 400 mpa from molecular dynamics, The Journal of
658 Chemical Physics 152 (19) (2020) 191104. doi:10.1063/5.0008907.
659 URL <https://doi.org/10.1063/5.0008907>
- 660 [66] S. Weisberg, Applied Linear Regression, John Wiley & Sons, Inc., 2005.
- 661 [67] D. C. Liu, J. Nocedal, On the limited memory bfgs method for large scale optimization, Math. Program. 45 (1-3) (1989) 503–528.
- 662 [68] GPy, GPy: A gaussian process framework in python, <http://github.com/SheffieldML/GPy> (since 2012).
- 663 [69] J. Bergstra, D. Yamins, D. D. Cox, Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision
664 architectures, in: Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28, ICML'13,
665 JMLR.org, 2013, p. 115–123.
- 666 [70] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization (2017). arXiv:1412.6980.
- 667 [71] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow,
668 A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah,
669 M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden,
670 M. Wattenberg, M. Wicke, Y. Yu, X. Zheng, TensorFlow: Large-scale machine learning on heterogeneous systems, software available from
671 tensorflow.org (2015).
672 URL <http://tensorflow.org/>
- 673 [72] L. L. Gratiet, J. Garnier, Recursive co-kriging model for design of computer experiments with multiple levels of fidelity, International Journal
674 for Uncertainty Quantification 4 (5) (2014) 365–386.
- 675 [73] M. C. Kennedy, A. O'Hagan, Predicting the output from a complex computer code when fast approximations are available, Biometrika 87 (1)
676 (2000) 1–13.
677 URL <http://www.jstor.org/stable/2673557>

678 Appendix A Data-fusion studies

679 In order to enhance the predictability of the ML models at transcritical regions, here we propose
680 a data-fusion approach. Specifically, we concatenate density points of the transcritical region
681 provided by the NIST database into the training dataset. The aim here is to improve the density
682 predictability of our ML models, by supplying reliable information about this state variable in the
683 specific region where MD data is scarce. Following this purpose, the first attempt is to add one
684 density point from the NIST database. Here, we concatenate the n-dodecane density at pressure
685 2 MPa and temperature 660 K to the training data. By adding this point to the training set, it is
686 verified that the ML models can recover the density at 660 K, as shown in Figure 15. However,

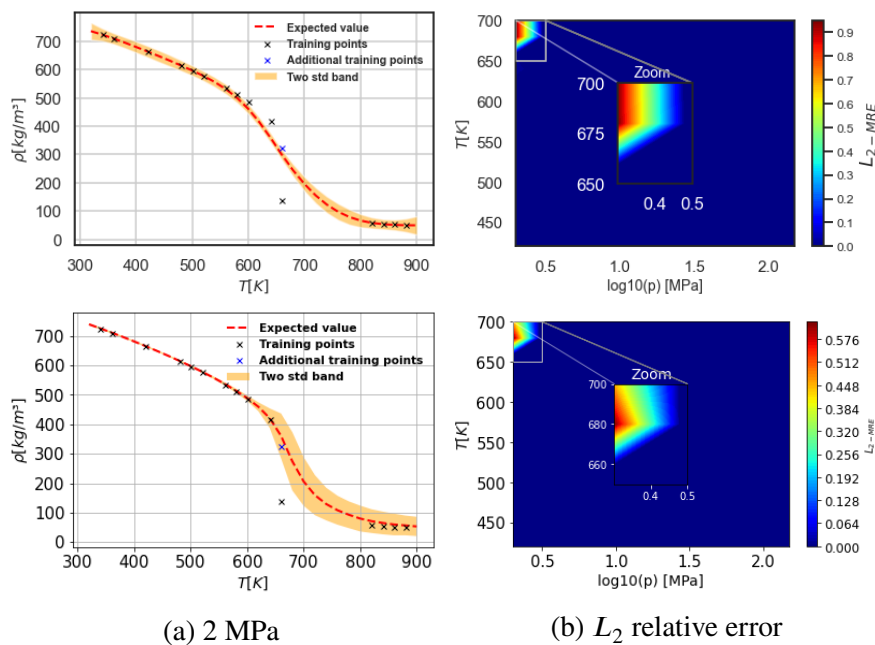


Figure 15: n-Dodecane density predictions with the GP model (top) and conditional generative model (bottom) using the data-fusion-fidelity approach with one density point from the NIST database.

687 the L_2 relative errors between the expected values predicted by ML models and the NIST data are
 688 still considerable in transcritical regions. Also, we can note that the conditional generative model
 689 has larger uncertainty bounds at the transcritical region trying to recover density behavior due to
 690 the lack of data in this region. Furthermore, Figure 16 shows that adding density points from NIST
 691 into the training data does not change the degree of uncertainty at other operating conditions.

692 As a further attempt to enhance the density predictions at the transcritical region, we now
 693 concatenate one more density point from the NIST database. More specifically, in addition to
 694 concatenating the n-dodecane density at pressure 2 MPa and temperature 660 K to the training
 695 data, we also add the n-dodecane density at 680 K. Figure 17 shows that adding two density points
 696 from NIST data in the transcritical region slightly improves the predictions of the GP model, while
 697 the relative error remains considerable. However, we can verify that the generative model returns
 698 satisfactory predictions with L_2 relative error lower than 10% in the transcritical region. This shows
 699 the capability of the conditional generative model to enhance the predictability of the density when

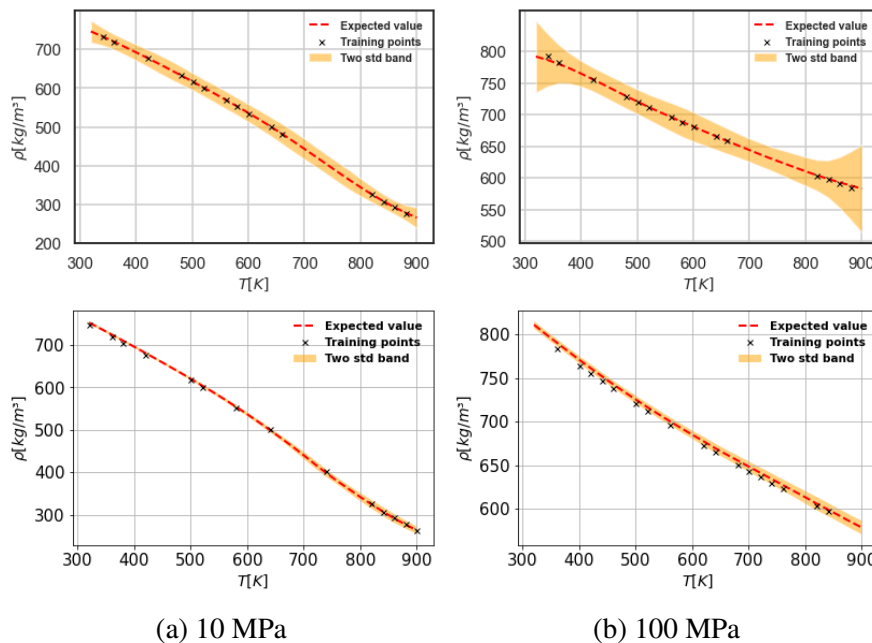


Figure 16: n-Dodecane density predictions with the GP model (top) and conditional generative model (bottom) using the data-fusion-fidelity approach with one density point from the NIST database at the pressures 10 and 100 MPa.

700 some pieces of information about the correct behavior of the transport property are given to the
 701 model.

702 Finally, to further increase the predictability of our ML models, a third attempt is proposed
 703 based on adding three density points from NIST to the training data, those being the n-dodecane
 704 densities at 2 MPa and temperatures 660, 680, and 700K. Figure 13 depicts that in this training
 705 scenario the density predictions of the GP model have some improvements, but the L_2 relative
 706 error is still considerable. Furthermore, we can verify that the conditional generative model returns
 707 accurate predictions, with relative errors lower than 5% in transcritical regions. Finally, we note
 708 that the generative model has uncertainty bounds able to recover the density predictions near the
 709 critical point.

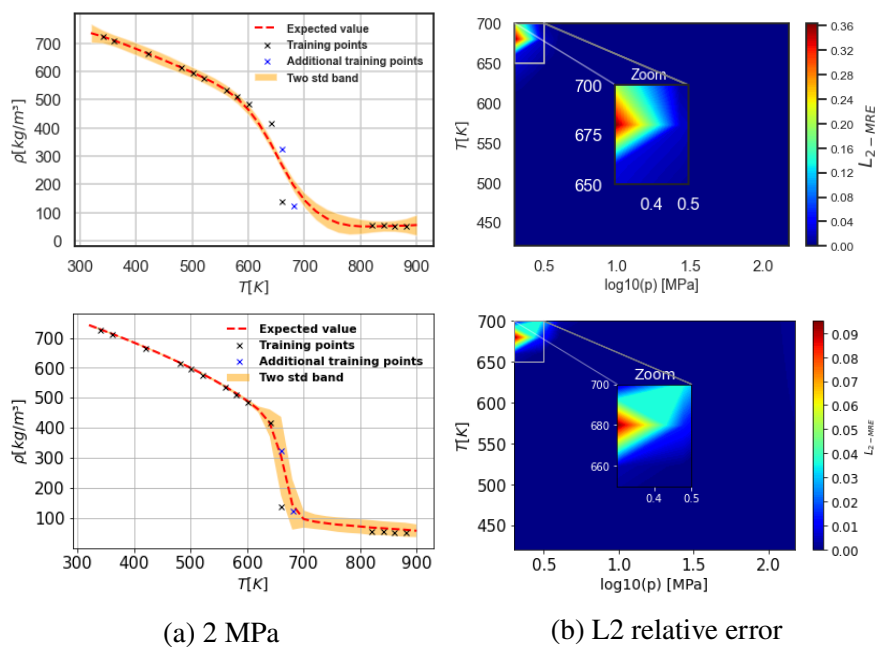


Figure 17: n-Dodecane density predictions with the GP model (top) and conditional generative model (bottom) using the data-fusion-fidelity approach with two density points from the NIST database.