

# Advancing Sustainable Aviation Fuel Design: Machine Learning for High Energy Density Liquid Polycyclic Hydrocarbons

*Dilip Rijal<sup>1</sup>, Vladislav Vasilyev<sup>2</sup> and Feng Wang<sup>1\*</sup>*

<sup>1</sup>Department of Chemistry and Biotechnology, School of Science, Computing and Engineering Technologies, Swinburne University of Technology, Melbourne, Victoria 3122, Australia

<sup>2</sup>National Computational Infrastructure, Australian National University, Canberra, ACT 0200, Australia

\*Corresponding author: [fwang@swin.edu.au](mailto:fwang@swin.edu.au)

## Abstract

Sustainable aviation fuels (SAFs) are crucial for addressing carbon emissions in the aviation industry. With a focus on SAFs, the research aims to establish a quantitative structure-property relationship for polycyclic hydrocarbons (PCHCs) and their net heat of combustion (NHOC) using the innovative approach of machine learning (ML). The model trained with support vector machine (SVM) algorithms in ML is selected as it demonstrates superior performance over other available algorithms with a high coefficient of determination ( $R^2$ ) and low mean absolute error (MAE) of 27.821 KJ/mol for 20% test data. Using the optimum SVM model, thirty-five potential PCHCs are identified as SAF candidates from C6 to C15 sourced from reputable scientific literature and databases. Furthermore, structural analysis revealed that high-performance PCHCs typically consist of saturated alkanes with multiple 3, 4, and 5-

1 membered rings, suggesting that strained energy plays a role in their high energy density. The  
2 model obtained from ML can be employed to screen new hydrocarbons for their suitability as  
3 SAF candidates before costly experiments and ASTM evaluations.

4

5 **Keywords:** Sustainable aviation fuel; Machine learning; Polycyclic hydrocarbons; Net heat of  
6 combustion; Fuel efficiency

7

8

### Nomenclature

ICAO	International Civil Aviation Organization
ML	Machine Learning
LTO	Landing and Take-Off
$\rho$	Density
SAFs	Sustainable Aviation Fuels
GUI	Graphical User Interface
PtL	Power-to-Liquid
HED	High Energy Density
PCHCs	Polycyclic Hydrocarbons
SVM	Support Vector Machines
NHOC	Net Heat of Combustion
KNN	K-Nearest Neighbor
NHOC <sub>G</sub>	Gravimetric Net Heat of Combustion
RF	Random Forest
NHOC <sub>V</sub>	Volumetric Net Heat of Combustion
DFT	Density Functional Theory
QSPR	Quantitative Structure-Property Relationship
AF	Aviation Fuel
R <sup>2</sup>	Coefficient of Determination
MAE	Mean Absolute Error
RMSE	Root Mean Square Error
GC	Group Contribution

9

## 10 1. Introduction

11 The aviation industry is flourishing with the steady growth of air travel worldwide,  
12 driven by factors such as increasing global connectivity, limited time, and expanding tourism  
13 industries, all contributing to heightened demand for fuel to power commercial aircraft.

1 According to an International Civil Aviation Organization (ICAO) report, it is predicted that  
2 aviation fuel (AF) consumption will increase by 1.9 to 2.6 times the value of 2018 in 2050 [1].  
3 Ultimately, heavy reliance on conventional AF such as Jet-A, Jet-A1, JP-4, and JP-5 poses a  
4 significant challenge in the fight against greenhouse gas (GHG) emissions. For instance, as per  
5 the forecast provided by ICAO, emissions from international aviation during both full-flight  
6 operations and landing and take-off (LTO) are anticipated to rise between 2 to 4 times by  
7 2050 compared to the level observed in 2018 [1], highlighting the urgent need for an  
8 alternative source of AF. Despite the pressing need for alternative, cleaner energy, options  
9 like batteries and hydrogen still need to be viable for immediate implementation in aviation  
10 due to technical and infrastructure limitations [2].

11 A promising progress in developing alternative source sustainable AF (SAFs) derived  
12 from Power-to-Liquid (PtL) pathways. In the PtL process, H<sub>2</sub> is typically obtained through  
13 water electrolysis. The electricity used for electrolysis is often sourced from renewable energy  
14 systems such as solar, wind, or hydropower. As for the carbon source, CO<sub>2</sub> is commonly  
15 sourced from various industrial processes, such as power plants, cement production, and  
16 steel manufacturing, where it is emitted as a byproduct, as well as directly from the  
17 atmosphere. The PtL process can help mitigate GHG emissions by converting CO<sub>2</sub> into valuable  
18 liquid fuels, thereby contributing to carbon neutrality or even carbon negativity [3]. It was  
19 revealed that the system's electrical efficiency is higher when the solid oxide electrolyzer  
20 operates in co-electrolysis mode compared to the steam mode for many hydrocarbon-based  
21 fuel production systems [3]. The PtL pathways produce synthetic fuels known as eFuel, which  
22 offer a potential solution to decarbonize the aviation industry.

23 Aviation fuel is a complex mixture of hydrocarbons like alkanes, cycloalkanes, and  
24 aromatics, further complicating the transition to greener alternatives. To enhance the

1 efficiency of eFuel and to increase the current 50% blending constraint of SAFs [4], there is a  
2 crucial need to optimize polycyclic hydrocarbons (PCHCs), the building blocks of these  
3 synthetic eFuels. Focusing on the composition and characteristics of PCHCs, one aims to  
4 develop more efficient and environmentally friendly alternatives to traditional AF, paving the  
5 way for a greener future in air travel [2, 5, 6]. eFuel may play a pivotal role in the aviation  
6 industry's imperative to reduce GHG emissions, offering an environmentally responsible  
7 alternative that holds the key to achieving long-term carbon-neutral growth and realizing net-  
8 zero targets.

9       Aviation fuel is subject to strict specifications. Evaluation of the physicochemical  
10 properties of each component in the AF mixture and the blended AF is a complex and  
11 expensive but essential process. First, fuels often comprise diverse components with unique  
12 chemical properties, requiring a comprehensive analysis of the mixture (blend) and individual  
13 components. Next, the physicochemical properties encompass a wide array of characteristics,  
14 including net heat of combustion (NHOC), density at 15°C, viscosity at -20°C and -40°C, flash  
15 point, surface tension at 22°C, cetane number, and octane number [4, 7]. Evaluating these  
16 properties for individual components in the blended fuel adds complexity. In addition, the  
17 compounds in a blended fuel may interact, leading to synergistic or antagonistic effects on  
18 the properties of the fuel. Third, aviation industries have stringent standards and  
19 specifications for fuel quality and performance. To meet these standards, a blended AF  
20 requires thoroughly evaluating its physicochemical properties. Finally, obtaining precise and  
21 accurate measurements and analysis of the fuel properties needs advanced testing  
22 equipment and methodologies to ensure the reliability of results. All contribute to the overall  
23 complexity and cost of the process.

1           The reliability and availability of AF depend on the quantity and composition of the  
2 components in the fuel. SAFs can be achieved by reducing or replacing unwanted components  
3 such as aromatics in conventional AF with strained PCHCs [5, 6] from PtL pathway [2] as  
4 blended fuel, and finally achieve 100% SAF. In the development phase of SAF, data obtained  
5 from measurements, testing, and computer modeling are employed to achieve insight into  
6 the desired properties from known structures of the compounds. For example, the net heat  
7 of combustion (NHOC) of AF is crucial for understanding a fuel's energy content and  
8 performance characteristics for energy efficiency, flight range, and payload capacity of jets  
9 [8]. The NHOC is categorized into gravimetric (NHOC<sub>G</sub>) and volumetric (NHOC<sub>V</sub>), and the  
10 former (NHOC<sub>G</sub>) is suitable for weight-limited aircraft, such as rockets and spacecraft.  
11 Likewise, the latter (NHOC<sub>V</sub>) helps to reduce the aircraft fuel tank volume. Therefore, for  
12 volume-limited aircraft such as missiles and military aircraft, high NHOC<sub>V</sub> fuel helps increase  
13 the payload without changing the tank size [9]. AF, such as RJ-4, RJ-5, and JP-10, have already  
14 been developed from PCHCs [10].

15           In this study, we developed a machine learning (ML) model from a training set of  
16 diverse hydrocarbons with known properties of NHOC and density. The model is then applied  
17 to pre-screen the chemical structures of PCHCs with high-performance for SAF applications.  
18 The present study is an initial pre-screen of a multilevel study aiming at rational design for  
19 high-performance SAF from the PtL pathway [2]. The selected high-performance PCHCs will  
20 be studied quantum mechanically using DFT calculations for their quantitative structure-  
21 property relationship (QSPR) to identify suitable SAF candidates, followed by synthesis  
22 reaction pathway and catalyst development, and finalized with a techno-economic  
23 assessment (TEA) for feasibility and costs associated.

## 24   **2. Methods development**

## 1 **2.1 General process of ML model**

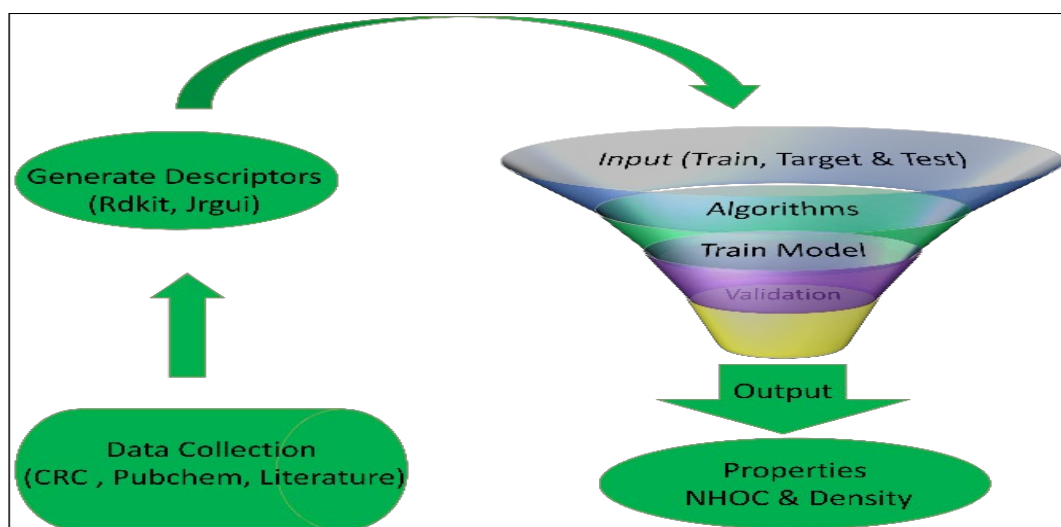
2           Developing predictive models for fuel properties, particularly focusing on net heat of  
3 combustion (NHOC) and density in alternative fuels, involves systematically integrating  
4 experimental or computational data and machine learning algorithms. Training data collection  
5 is a crucial first step, which depends on experimentally measured values. Experimental  
6 measurements utilize different devices. For instance, calorimeters are the major  
7 measurement devices for accurately and reliably determining the NHOC of alternative fuels.  
8 A bomb calorimeter determines the NHOC by quantifying the energy released as the fuel  
9 undergoes combustion, a straightforward and popular method [11]. Other calorimeters  
10 include adiabatic flame calorimeters and oxygen bomb calorimeters [12]. However, practical  
11 constraints such as cost, fuel volume, and time impede the acquisition of experimental data  
12 [4].

13           In addition to experimental data, one can also use density functional theory (DFT)  
14 [13, 14] computed data for training models, as DFT is a reliable and accurate computational  
15 method. For example, Alibakhshi recently estimated the NHOC of up to 40 organic molecules  
16 using the more expensive quantum mechanical CCSD-F12b and DSD-PBEB86 methods,  
17 achieving correlation coefficients of 0.9999 and 0.9998, respectively [13]. Similarly, another  
18 recent study calculated the NHOC for as many as 295 sesquiterpenoid high-energy density  
19 (HED) fuels and reported an average absolute error as small as 2.6% [14]. While these methods  
20 ensure high accuracy, they can be resource-intensive and time-consuming.

21           To address the limitations of experimental and computationally demanding methods,  
22 predictive methodologies such as group contribution (GC) methods [15, 16] and machine  
23 learning (ML) algorithms [17, 18] offer efficient alternatives. Albahri [19] introduced a more  
24 accurate GC method for computing net heat of combustion (NHOC). This method computes

1 32 atom-type structural groups for NHOC in up to 452 hydrocarbons, yielding NHOC  
2 predictions with an average absolute error of 0.71% and a correlation coefficient (R) of 0.9982  
3 [19]. However, this method may be time-consuming for many hydrocarbons and may be  
4 unsuitable for newly synthesized complex molecules. GC methods are empirical [15, 16],  
5 relying on empirical data and expert knowledge to assign contributions to functional groups.

6 Machine learning (ML), in particular, emerges as a cost-effective option for analyzing  
7 large datasets and predicting complex chemical properties by following the general steps in  
8 Fig. 1. By training models on molecular structures and corresponding descriptors, ML extracts  
9 patterns and relationships, facilitating the rapid screening of numerous fuel candidates. This  
10 predictive capability complements experimental measurements and computationally  
11 demanding quantum chemical calculations, enabling the virtual construction of desired  
12 molecular structures and significantly reducing the time and resources required for fuel  
13 screening. This rapid screening identifies the most promising candidates, saving significant  
14 time and resources before future investigation can be done for NHOC (ASTM D4809) and  
15 density ( $\rho$ ) at 15° C (ASTM D4052) and other alternative fuel properties [4]. The general  
16 process of ML is presented in Fig. 1.



17  
18 Fig. 1. Flowchart of model development in ML process in the present study.

## 1 2.2. Data collection and feature engineering

2 A comprehensive dataset, AFProp(N, M), is utilized in this study as a training dataset,  
3 where N represents the number of fuel properties, and M denotes the number of organic  
4 compounds with known properties. Focusing on two fundamental properties (N=2), NHOC  
5 and density ( $\rho$ ). For NHOC properties, the dataset AFProp(1, M=452) compiles data on up to  
6 452 pure hydrocarbons, encompassing paraffins, olefins, naphthenes, and aromatics, sourced  
7 from reference [19]. This dataset incorporates NHOC values obtained from experimental  
8 measurements and calculations drawn from the American Petroleum Institute - Technical  
9 Data Book (API-TDB) [20]. Additionally, the dataset AFProp(2, M=486) comprises density ( $\rho$ )  
10 properties at temperatures ranging from 15°C to 30°C for up to 486 distinct hydrocarbons,  
11 sourced from the CRC Handbook of Chemistry and Physics [21] and [18]. Up to 17  
12 hydrocarbons were excluded from the NHOC dataset AFProp(1, 452) because they were  
13 identified as duplicates due to sharing identical Simplified Molecular Input Line Entry System  
14 ([SMILES](#)) notation [22] for being isomeric hydrocarbons. Furthermore, certain data points  
15 were computed values that displayed significant deviations, as highlighted by Albahri [19] and  
16 were excluded. Hence, the dataset for NHOC becomes AFProp(1, 435), which is the final  
17 training dataset for NHOC.

18 The chemical structures of the candidate hydrocarbons undergoing screening are also  
19 compiled into a dataset, SAFCan(P, Q), where they are systematically represented using  
20 [SMILES](#) notation. The [SMILES](#) notations of dataset AFProp(N, M) and 30 existing PCHCs were  
21 directly sourced from Pubchem [23]. Additionally, the [SMILES](#) notation for five novel PCHCs  
22 circled in Fig. 5 was generated using the Open Babel toolbox [24], ensuring a compact and  
23 unambiguous representation of each compound's structure. We utilize a Python program



1 called GUIDEMOL [25, 26] to derive molecular descriptors from the [SMILES](#) representations.  
2 GUIDEMOL leverages the RDKit toolkit for cheminformatics [25, 26], offering a range of  
3 functionalities, including molecular structure handling, substructure searching, molecular  
4 similarity calculation, chemical reaction handling, and descriptor calculation.

### 5 **2.3. ML configuration and molecular descriptors**

6 Machine learning (ML) relies on hyperparameters, which are external configuration  
7 settings that cannot be learned directly from the data. These parameters are predetermined  
8 and remain fixed throughout training, influencing the model's behavior and performance.  
9 Precisely adjusting hyperparameters is essential for accurate model predictions and  
10 optimizing predictive precision. To achieve this, hyperparameter tuning techniques such as  
11 GridSearchCV [27] are employed. The optimal parameter obtained in different algorithms for  
12 NHOC and density ( $\rho$ ) using GridSearchCV refer to Table S1 in supplementary. GridSearchCV  
13 systematically explores a subset of hyperparameters, evaluating the model's performance  
14 through cross-validation of the training data. This approach simplifies the tuning process and  
15 improves predictive accuracy and optimized model performance [27].

16 Molecular descriptors or features representing measurable data points' properties  
17 were generated using the JRgui graphical user interface (GUI) [25, 28] powered by the Tkinter  
18 package. This tool computes descriptors integrated into RDKit [26] and generates grid  
19 representations of 3D molecular structures [29]. Leveraging the JRgui and RDKit toolkit [25,  
20 28, 29], we extracted a comprehensive set of approximately 200 descriptors. However, after a  
21 thorough examination, only 6 descriptors for NHOC and 40 for density ( $\rho$ ) were meticulously  
22 selected from this set to train the models. In other words, the hydrocarbons' properties can  
23 be expressed as a function of these descriptors,

$$AFProp(j, i) = \sum_{i=1}^{n_j} C_{ji} x_i + PConst_j \quad (1)$$

Where  $AFProp(j, i)$  presents the aviation fuel compound property  $j$  and its descriptor  $i$ . In the present study,  $j=1, 2$  as only two properties, NHOC and density ( $\rho$ ), are investigated. As a result,  $AFProp(1, i=1,2,\dots,6)$  for NHOC and  $AFProp(2, i=1,2,\dots,40)$  for density ( $\rho$ ).  $PConst_j$  represents the intercept (constant term) of property  $j$ . While  $x_i$  is descriptors, the index  $i$  runs from 1 to  $n_j$  (the number of descriptors) for the property  $j$  of the compound under study. For example, for NHOC ( $AFProp(1, i=1,2,\dots,6)$ ) property of a compound contains six descriptors ( $n_1=6$ ), and 40 descriptors ( $n_2=40$ ) for density ( $\rho$ ) ( $Prop_2$ ). That is,  $AFProp(1, 6)$  for NHOC and  $AFProp(2, 40)$  for density. Here,  $C_{ji}$  ( $i=1, 2,\dots,n_j$ ) are the obtained coefficients for property  $j$  obtained from the ML model through the training dataset. Table 1 reports the information of the six descriptors  $x_i$  ( $i=1, 2,\dots,6$ ) of NHOC, whereas the 40 descriptors and corresponding coefficients obtained for the density ( $\rho$ ) property ( $j=2, n_2=40$ ) in the model are given in Table S2 in the supplementary materials.

Table 1: The descriptors of property NHOC ( $AFProp(1, i=1, 2,\dots,6)$ ) in the ML model.\*

Symbol	$AFProp(1, i=1, 2,\dots,6)$	Description
$X_1$	NC	Number of carbons
$X_2$	NH	Number of hydrogens
$X_3$	Num_of Atoms	Total number of atoms
$X_4$	BalabanJ	Topological index
$X_5$	NumAromaticRings	Number of aromatic rings
$X_6$	Kappa3	Coefficient of characteristics

\*Descriptors with a strong correlation to NHOC are highlighted in grey.

As can be seen in Table 1, the three descriptors, NC, NH, and Num\_of \_Atoms, are dependent as  $NC+NH=Num\_of\_Atoms$  for hydrocarbons. Theoretically, orthogonal (independent) descriptors in Equ (1) are preferred because they simplify the interpretation of the model and reduce multicollinearity, which can lead to instability in the estimates of model coefficients. However, in practice, it's not always possible to have completely orthogonal

1 descriptors. Sometimes, including non-orthogonal variables can improve the fitting for  
2 several reasons. For example, non-orthogonal variables might provide additional information  
3 that improves the model's predictive performance and reduces bias; they can also handle  
4 nonlinear relationships between the descriptors and the target variable, leading to more  
5 accurate predictions. Sometimes, domain-specific tasks need non-orthogonal variables that  
6 are known to be relevant to the prediction task, even if they are correlated with other  
7 variables. Moreover, Equ (1) can also be extended to other organic compounds rather than  
8 hydrocarbons. We are working on developing a new set of orthogonal descriptors for AF  
9 candidates, which are closely related to molecular structures in 3D and energies.

#### 10 **2.4. Machine learning model training**

11 After collecting and performing feature engineering on the dataset AFProp(N, M) and  
12 SAFCand(P, Q), the subsequent step entails partitioning the dataset AFProp(N, M) into (80%)  
13 for training and (20%) for testing. Following this partitioning, the next step involves selecting  
14 a suitable algorithm and training the model for analysis. This study chooses a supervised  
15 machine learning (ML) training model, as it can accurately predict target properties [30].  
16 Several commonly used algorithms within the supervised ML training model are considered,  
17 including the support vector machines (SVM) [31], random forest (RF) [32], and k-nearest  
18 neighbors (KNN) [33]. These algorithms are evaluated to determine the most appropriate one  
19 for the task. The model training process is conducted using Google Collaboratory [34], based  
20 on Python 3.10 [35]. This platform offers data exploration and visualization flexibility through  
21 libraries such as Pandas, NumPy, and Matplotlib. Additionally, it seamlessly integrates with  
22 ML libraries like Scikit-Learn, TensorFlow, and PyTorch, providing a comprehensive toolkit for  
23 model development [36].

1           In summary, as indicated earlier, the ML training process depicted in Fig. 1 of the  
2 present study begins with data collection from multiple databases and resources. Descriptor  
3 generation follows, utilizing tools like Rdkit and Jugui [25, 28, 29]. Subsequently, a supervised  
4 training model is chosen, employing an appropriate algorithm such as SVM. The output  
5 properties are then validated against a set of properties with available experimental values.  
6 This iterative process continues until the output aligns with the target fuel properties.

### 7 **3. Results and discussion**

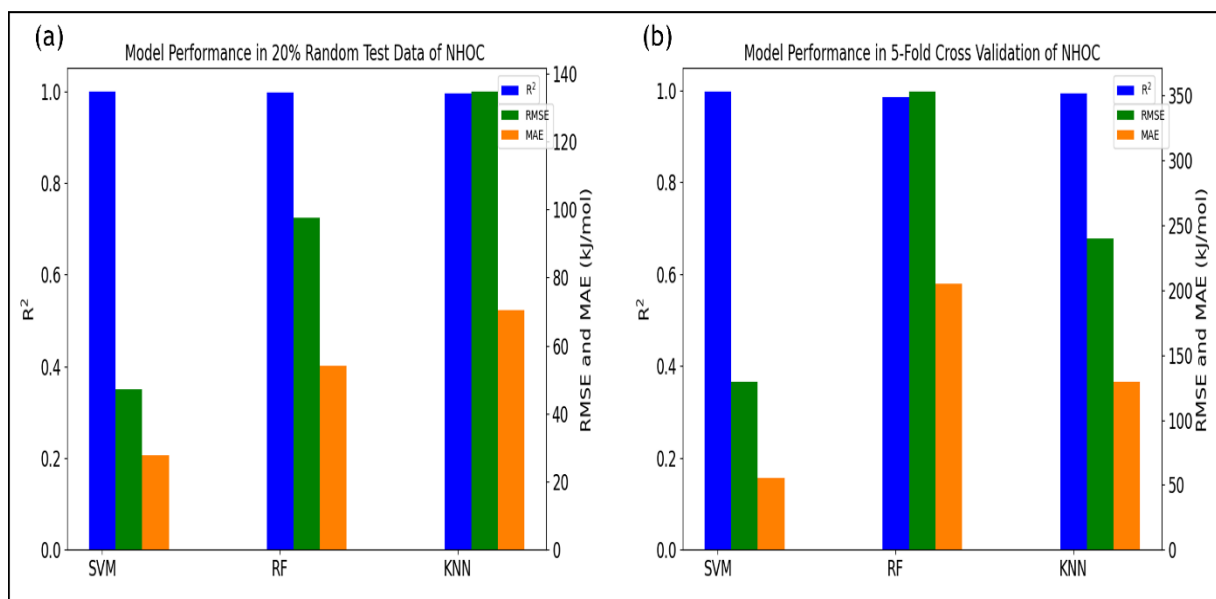
#### 8 **3.1. Performance of the algorithms**

9           The performance of three major supervised ML algorithms is considered for training.  
10 That is, the support vector machines (SVM) [31], random forest (RF) [32], and k-nearest  
11 neighbors (KNN) [33]. Evaluation metrics such as mean absolute error (MAE), root mean  
12 square error (RMSE), and coefficient of determination ( $R^2$ ) are employed for the performance  
13 of these algorithms. MEA and RMSE are both metrics used to evaluate the accuracy of  
14 predictions made by an algorithm or model, whereas  $R^2$  quantifies the extent to which the  
15 model's predicted values align with the observed ones [18, 37]. The optimal algorithm  
16 performance is characterized by achieving the maximum value of  $R^2$  (0-1) while minimizing  
17 the values of MAE and RMSE.

18           The 5-fold cross validation (5-fold CV) is a widely used approach for assessing  
19 prediction accuracy and validating machine learning models applied to evaluate their efficacy  
20 [38]. In a 5-fold CV, the data are randomly divided into 5 folds or groups, and the model's  
21 ability is summarized using the sample of model evaluation scores. Moreover, as mentioned  
22 earlier, the dataset was randomly divided into training (80%) and testing (20%) subsets to  
23 ensure robust model development. Additionally, descriptors were standardized to bring

1 different units onto a common scale without altering their original units, thus enhancing the  
2 model's performance. Combining these techniques can avoid issues such as overfitting and  
3 underfitting and obtain a sense of how the model will transfer to a different dataset [39]. This  
4 rigorous approach aimed to improve the accuracy of the training data, resulting in more  
5 precise predictions. Fig. 2 compares the performance of the three algorithms, SVM, RF, and  
6 KNN, in the prediction of NHOC at 20% random test data (Fig. 2a) and at 5-fold cross-  
7 validation (Fig. 2b). Detailed results can be found in Tables S3 and S4 in the supplementary  
8 materials.

9 Fig. 2 illustrates the coefficient of determination ( $R^2$ ) of NHOC produced by three  
10 algorithms, SVM, RF, and KNN, all of which are close to 1.0, indicating a high level of  
11 agreement between predicted and observed values. In terms of RMSE values for 20% random  
12 test data (green in Fig. 2a), SVM performs the best NHOC value of 47.237 kJ/mol, followed by  
13 RF with the NHOC value of 97.493 kJ/mol, and KNN with 134.753 kJ/mol. Similarly, the MAE  
14 values (orange in Fig. 2a) are the smallest for SVM at 27.821 kJ/mol again, while RF and KNN  
15 have higher MAE values of 54.058 kJ/mol and 70.472 kJ/mol, respectively. These trends are  
16 consistent in Fig. 2b as well. Overall, the SVM algorithm demonstrates superior performance  
17 compared to RF and KNN, making it the preferred choice for further calculations. The SVM  
18 algorithm is known for its computational efficiency and robust predictive capabilities,  
19 particularly when dealing with limited data availability [40].



1  
2 Fig. 2. Comparison of the performance parameters R<sup>2</sup> (Blue), RMSE (green), and MAE (orange)  
3 of three algorithms, SVM, RF, and KNN prediction of NHOC. (a) 20% random test data, and (b)  
4 5-fold cross-validation.

### 5 3.2. Performance of the SVM trained model on NHOC

6 The SVM algorithm is applied to the training dataset so that the model (i.e., the  
7 coefficients in Equ (1)) is trained and obtained. Table 2 represents the coefficients obtained  
8 from the SVM algorithm used in the machine learning model (see Equ (1) and Table 1).

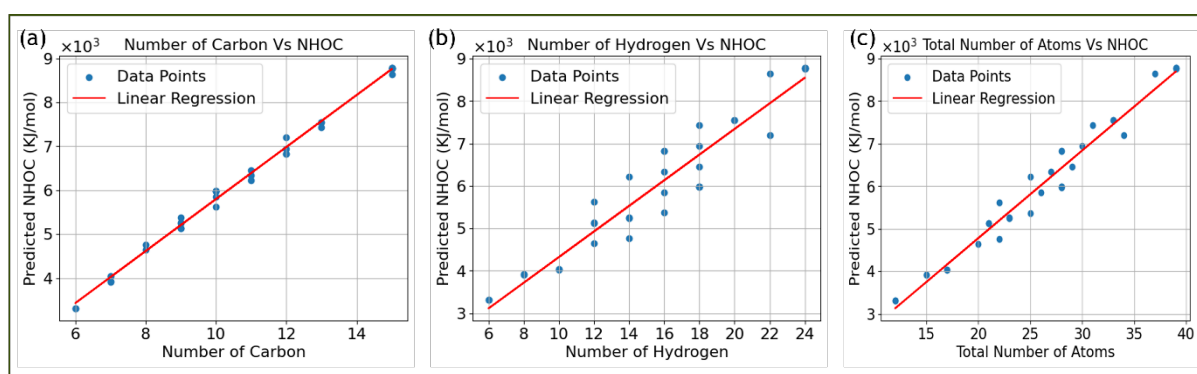
9 Table 2: The model (coefficients of Equ (1)) obtained for NHOC property.

Descriptors	Coefficients symbol	Coefficients Values
NC	C <sub>11</sub>	509.9569
NH	C <sub>12</sub>	78.8326
Num_of Atoms	C <sub>13</sub>	-21.0102
BalabanJ	C <sub>14</sub>	13.3489
NumAromaticRings	C <sub>15</sub>	-205.7653
Kappa3	C <sub>16</sub>	5.9531
Intercept	PConst <sub>1</sub>	0.6981

10

11 As reported in Table 2, the 6+1 coefficients of the multiple descriptor linear equation  
12 (Equ (1)) are positive except for the total number of atoms (Num\_of\_Atoms) and the number

1 of aromatic rings (NumAromaticRings), which are negative. A positive coefficient suggests  
 2 that an increase in the descriptor is associated with an increase in the target property (NHOC),  
 3 whereas a negative coefficient suggests that a decrease in the descriptor is associated with  
 4 an increase in the target property (NHOC). As a result, if one wishes to enhance the target  
 5 property NHOC, the descriptors with positive coefficients in Table 2, including NC, NH,  
 6 BalabanJ, and Kappa3, need to be enhanced. For example, the largest coefficient of the  
 7 multiple descriptor linear equation (Equ (1)) is NC, the number of carbons, with a coefficient  
 8 value as large as 509.9569. Fig. 3 displays the relationship of the NHOC property with the  
 9 number of atoms.



10

11 Fig. 3. Impact of positive descriptors on NHOC. (a) NC (Number of carbons), (b) NH (Number  
 12 of hydrogens), and (c) Total number of atoms in PCHCs (NC+NH).

13 However, for AF fuel candidate hydrocarbons, the NC of hydrocarbons can only  
 14 increase within a boundary of approximately  $6 < NC < 17$ , although this boundary varies due  
 15 to structures and is possible for higher NC hydrocarbons for liquid. If the NC number of a  
 16 hydrocarbon compound is up to 17, the maximum number of hydrogens is no more than  $NH$   
 17  $< NC \times 2 + 2$  (36) (for n-alkanes, any unsaturated carbons and rings will reduce the number of  
 18 hydrogens). As a result, to design novel hydrocarbons with high NHOC, one needs to increase  
 19 the structure and topological descriptors, BalabanJ and Kappa3, by designing novel

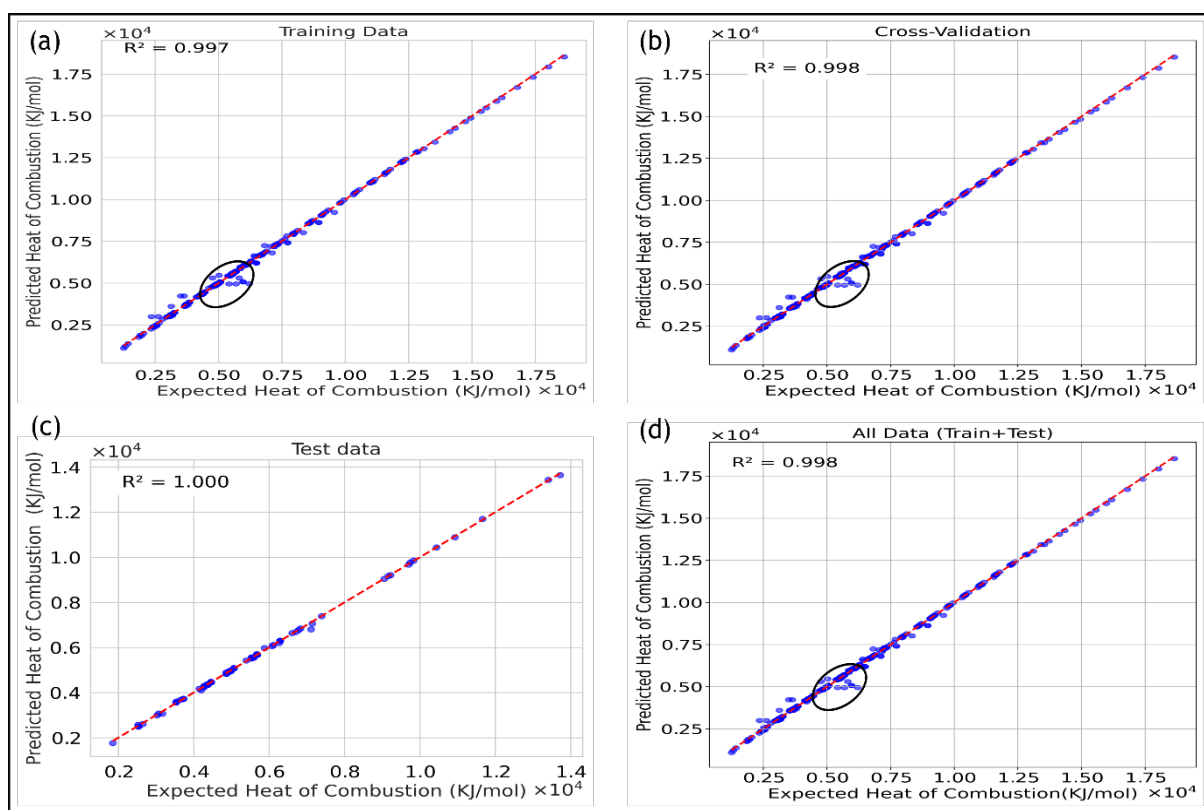
1 hydrocarbon structures. Here, BalabanJ provides essential topological insights [41], and  
2 Kappa3 is a third-order molecular connectivity index that captures crucial information  
3 regarding a molecule's structural topology and connectivity [41].

4 In addition, one of these negative coefficients in Table 2, Num\_of Atoms, is not  
5 orthogonal, as indicated earlier, which restricts the increase of NC and NH. The other negative  
6 coefficient, NumAromaticRings, is unwanted, as aromatic compounds produced 88% more  
7 soot formation than cycloalkanes due to incomplete combustion [42]. Consequently, the  
8 coefficients reported in Table 2 of the ML model provide rich information for the future  
9 development of the ML model (Equ (1)) with more efficient descriptors and molecular  
10 structures for new candidates.

11 The accuracy of the SVM-trained model was evaluated by applying it to compute the  
12 NHOC of hydrocarbons in various datasets, including training, testing, cross-validation (5-  
13 fold), and combined datasets (training + testing). Fig. 4 illustrates the correlation between the  
14 predicted NHOC values and the measured NHOC values in different datasets, such as the  
15 training dataset 80% of (AFProp(N, M)), the cross-validation (5-fold) dataset, the 20% test  
16 dataset of (AFProp(N, M)), and the combined dataset (train+test or AFProp(N, M)). The  
17 consistently high  $R^2$  values near unity with not less than 0.997 for NHOC indicate the model's  
18 high accuracy in NHOC property prediction. In addition, the SVM algorithm also demonstrates  
19 comparable proficiency in density ( $\rho$ ) prediction. For more detailed information on the  
20 performance of the SVM model in density ( $\rho$ ) estimation, please refer to Fig. S1 and S2 in the  
21 Supplementary Materials. These findings highlight the effectiveness of the ML model in  
22 precisely predicting the NHOC and density ( $\rho$ ) of SAF candidates.



1           The agreement between the predicted and the literature NHOC values of the  
2 compounds is excellent. Most of the compounds are along a straight line except for a small  
3 number of compounds in Fig. 4 (a, b, and d), with minor discrepancies as indicated in the oval.  
4 Specifically, the NHOC of the hydrocarbons in the vicinity of  $0.50 - 0.70 \times 10^4$  KJ/mol exhibits  
5 larger discrepancies, which are not seen in the test dataset (Fig. 4c). Further examination of  
6 the datasets reveals that the molecules in training data sets such as 2,3 pentadiene, 2-Methyl-  
7 2,4-Hexadiene, 2-Methyl-1,5-Hexadiene, 2,3-Hexadiene, 2,4-dimethylhexane, and 2,3-  
8 dimethyl-1-hexene were either with the computed NHOC or large errors in reference data as  
9 highlighted by Albahri [19]. Despite discrepancies from reference data, our SVM-trained  
10 model consistently produces accurate results across various datasets.

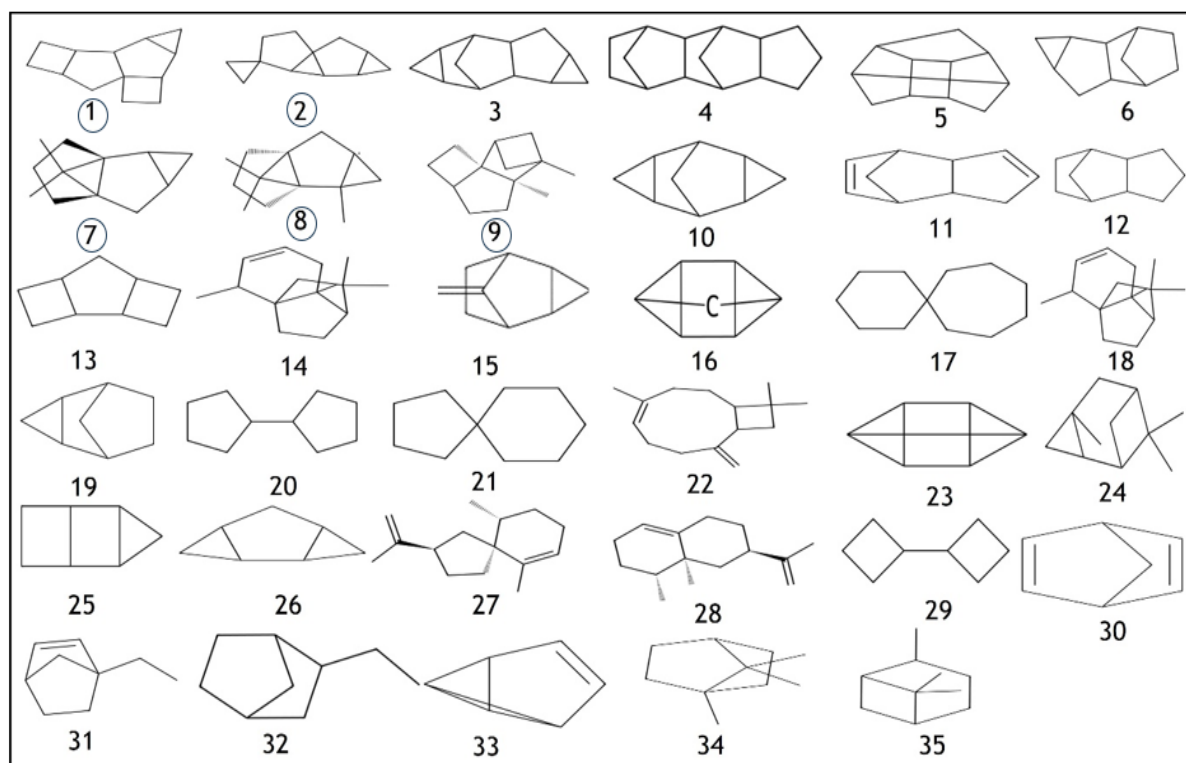


11  
12 Fig. 4. Predicted versus expected NHOC using the SVM model. (a) Training data (b) Cross-  
13 validation (5-fold) data (c) 20% test data (d) All data (Train+Test)

### 14 3.3 Screen of liquid hydrocarbons with high NHOC

1           The coefficient of the SVM-trained model given in Table 2 is employed to examine the  
2 hydrocarbons in SAFcand(P, N), which shows the structural properties of various  
3 hydrocarbons in Supplementary Fig. S3. The majority of PCHCs in SAFcand(P, Q) align with  
4 currently available hydrotreated esters and fatty acids (HEFA) SAFs, which usually contain C9-  
5 C16 carbons [43], and SAFs hydrocarbons produced from biomass or other waste in the range  
6 of C8-C18 [44]. Although squalane (C<sub>30</sub>H<sub>62</sub>), a saturated hydrocarbon with an IUPAC name  
7 of 2,6,10,15,19,23-hexamethyltetracosane, is a liquid hydrocarbon with up to 30 carbon  
8 atoms, the majority of liquid hydrocarbons under ambient temperature do not exceed 17  
9 carbon atoms, whereas the number of hydrogens of the potential PCHCs ranges from H6 to  
10 H24, as shown in Fig. 3 (b). Such the numbers of carbons and hydrogens in the compounds in  
11 the SAFcand(P, Q) dataset suggest these hydrocarbons are likely polycyclic with saturated C-  
12 C bonds and polycyclic hydrocarbons (PCHCs). Leveraging the structural information will help  
13 to design novel PCHCs with preferred fuel properties, such as high NHOC and high energy  
14 density [45]. Utilizing the SVM-trained model (Equ (1) and Table 2), up to 35 PCHCs are  
15 selected as suitable SAF candidates. Fig. 5 reports the chemical structures of these identified  
16 candidates using Google Collaboratory [34]. Most selected hydrocarbons correspond to  
17 existing compounds, except five PCHCs structures numbered 1 (6377), 2 (250609), 7 (268141),  
18 8 (82630), and 9 (33744) (circled) in Fig. 5, which were designed for HED fuel applications in a  
19 previous study [45]. The remaining 30 polycyclic hydrocarbons (PCHCs) were sourced from the  
20 PubChem database [23]. As shown in Fig. 5, these PCHCs predominantly comprise saturated  
21 cycloalkanes except for ten compounds (28.571%) containing unsaturated C=C bonds. This  
22 observation agrees with the fact that saturated hydrocarbons are often preferred for SAFs  
23 compared to unsaturated ones [4]. The majority (77.143%) of compounds in Fig. 5 exhibit

1 pentagon ring configuration, and nearly half (42.857%) contain triangular rings, consistent  
2 with the outcome reported earlier [45].



3  
4 Fig. 5. Chemical structures of 35 PCHCs obtained from ML screen in the present study.  
5 Structure 12 (exo-Tetrahydrodicyclopentadiene) is the dominant component of aviation fuel  
6 JP-10.

7 The fuel properties such as gravimetric NHOC ( $NHOC_G$ ), volumetric NHOC ( $NHOC_V$ ),  
8 H/C ratio, and density ( $\rho$ ) of these liquid PCHCs were obtained using the present ML are  
9 summarised in Table 3. Note that gravimetric NHOC ( $NHOC_G$ ) and density ( $\rho$ ) are obtained  
10 from ML, and other properties/descriptors such as volumetric NHOC ( $NHOC_V$ ), H/C ratio, and  
11 the total number of rings ( $N_{ring}$ ) are derived. The total number of rings ( $N_{ring}$ ) in Table 3 can be  
12 obtained from RDKit [28, 29] or counted manually from the structures. As seen in the table,  
13 almost all these compounds contain either triangular rings or rectangular rings with acute  
14 angles or pentagon rings, indicating that they are strained with possibly higher internal energy.

- 1 The PCHCs in the table exhibit required ranges of  $NHOC_G$  of 42.366-43.277 MJ/kg and  $NHOC_V$
- 2 of 35.849-52.039 MJ/L.
- 3 Table 3. Properties of selected PCHCs for SAF using ML. #

No	Hydrocarbon	Formula CAS No	H/C	$\rho$ (g/ml)	$NHOC_G$ (MJ/Kg)	$NHOC_V$ (MJ/L)	$N_{ring}$
1	6377*	C13H18	1.385	1.221	42.620	52.039	5
2	250609*	C12H16	1.333	1.221	42.565	51.972	5
3	Pentacyclo (6.3.1.0(2,7).0(3,5).0(9,11)) dodecane	C12H16 82110-70-1	1.333	1.200	42.573	51.087	5
4	THTCPD pentacyclo (6.5.1.13,6.02,7.09,13) pentadecane	C15H22 75172-85-9	1.467	1.192	42.701	50.900	5
5	Pentacyclo (5.4.0.02,6.03,10.05,9) undecane	C11H14 4421-32-3	1.273	1.165	42.529	49.547	5
6	Tetracyclo (6.2.1.0(2,7).0(3,5)) undecane	C11H16 1 777-44-2	1.455	1.093	42.750	46.726	4
7	268141*	C12H18	1.500	1.083	42.780	46.331	4
8	82630*	C13H20	1.538	1.075	42.813	46.024	4
9	33744*	C13H20	1.538	1.058	42.819	45.303	5
10	Tetracyclo (3.3.1.02,4.06,8) nonane	C9H12 187-49-5	1.333	1.062	42.646	45.290	4
11	Dicyclopentadiene	C10H12 77-73-6	1.200	1.018	42.517	43.282	3
12	Exo-THDCPD (JP-10)	C10H16 2825-82-3	1.600	0.990	42.968	42.539	3
13	Tricyclo (5.2.0.02,5) nonane	C9H14	1.555	0.982	42.928	42.156	3
14	Alpha neoclovene	C15H24 45-45-68-0	1.600	0.972	42.913	41.712	3
15	Tricyclo (3.2.1.0(2,4)) octane	C9H12 38310-48-4	1.333	0.970	42.900	41.613	3
16	Quadricyclane (QC)	C7H8 278-06-8	1.143	0.978	42.492	41.557	5
17	Spiro (5,6) dodecane	C12H22 181-15-7	1.833	0.957	43.277	41.416	2
18	Gamma neoclovene	C15H24	1.600	0.957	42.913	41.068	3
19	Tricyclo (3.2.1.02,4) octane	C8H12 13377-46-3	1.500	0.952	42.900	40.841	3
20	Bicyclopentane	C10H18 1636-39-1	1.800	0.933	43.252	40.354	3
21	Spiro (4,5) decane	C10H18 176-63-6	1.800	0.931	43.259	40.274	2
22	Caryophyllene	C15H24	1.600	0.928	42.976	39.882	2

23	Prismane	87-44-5 C <sub>6</sub> H <sub>6</sub> 650-42-0	1.000	0.941	42.366	39.866	5
24	4,7,7-Trimethyltricyclo (4.1.1.0 <sub>2,4</sub> ) octane	C <sub>11</sub> H <sub>18</sub>	1.636	0.927	42.981	39.843	3
25	Tricyclo (3.2.0.0 <sub>2,4</sub> ) heptane	C <sub>7</sub> H <sub>10</sub> 28102-61-6	1.429	0.929	42.834	39.793	3
26	Tricyclo (4.1.0.0 <sub>2,4</sub> ) heptane	C <sub>7</sub> H <sub>10</sub> 187-26-8	1.429	0.927	42.845	39.717	3
27	Premnaspirodien	C <sub>15</sub> H <sub>24</sub> 82189-85-3	1.600	0.915	42.970	39.318	2
28	Valencene	C <sub>15</sub> H <sub>24</sub> 3-07-4630	1.600	0.912	42.975	39.193	2
29	Bicyclobutyl	C <sub>8</sub> H <sub>14</sub> 7051-52-7	1.750	0.897	43.214	38.763	2
30	Norbornadiene	C <sub>7</sub> H <sub>8</sub> 16422-76-7	1.143	0.895	42.608	38.134	2
31	Ethylnorbornene	C <sub>9</sub> H <sub>14</sub> 2146-41-0	1.556	0.884	43.005	38.016	2
32	5-Ethylnorbornane	C <sub>9</sub> H <sub>16</sub> 2146-41-0	1.778	0.860	43.254	37.191	2
33	Benzvalene	C <sub>6</sub> H <sub>6</sub> 659-85-8	1.000	0.875	42.404	37.104	4
34	Camphane	C <sub>10</sub> H <sub>18</sub> 464-15-3	1.800	0.851	43.256	36.811	2
35	Pinane	C <sub>10</sub> H <sub>18</sub> 473-55-2	1.800	0.829	43.243	35.849	3

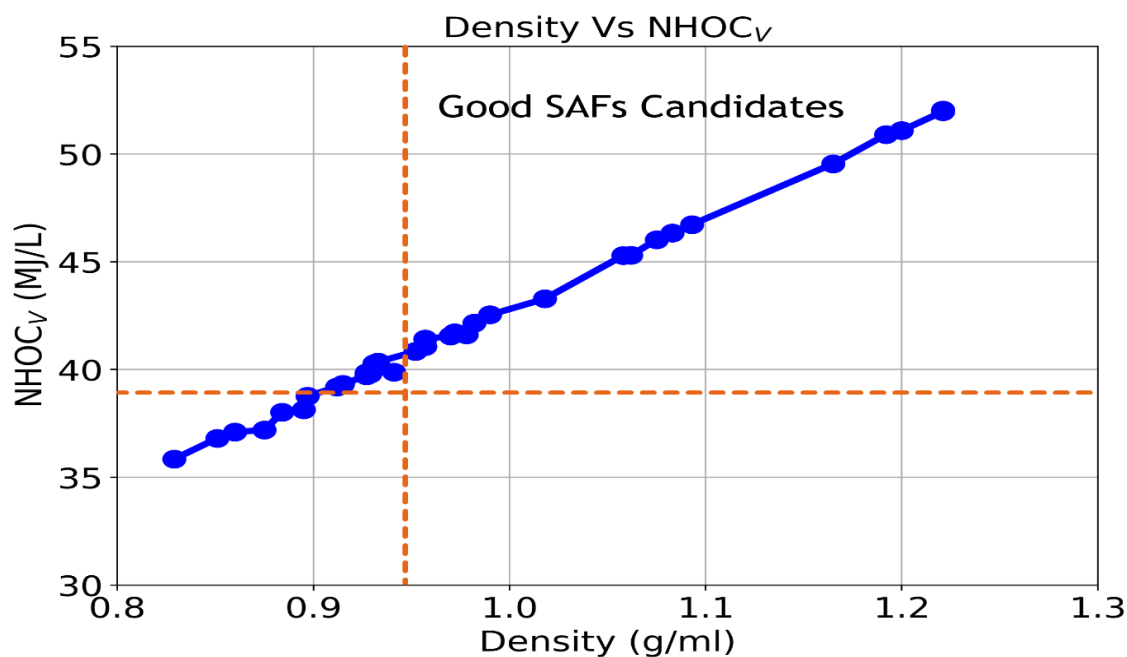
1 #The  $NHOC_G$  of jet fuels is 42.20-43.98 MJ/kg and  $NHOC_v$  is 32.26-39.64 MJ/L [10]. \*Ref [45]

2 (no CAS numbers available).

3       There exists a linear positive relationship between  $NHOC_G$  and the H/C ratio of the  
4 PCHCs. That is, the larger the H/C ratio, the larger the  $NHOC_G$ . For example, Spiro (5,6)  
5 dodecane (No 17) exhibits the largest  $NHOC_G$  of 43.277 MJ/Kg among the set of PCHCs with  
6 an H/C ratio of 1.833, whereas Prismane (No 22) displays the smallest  $NHOC_G$  of 42.366 KJ/kg,  
7 as detailed in Table 3, is characterized by an H/C ratio of 1. A larger H/C ratio suggests a  
8 preference for saturated hydrocarbons with more C-H bonds, aligning with the fact that  
9 composition standards of SAFs are characterized by a higher H/C ratio [46]. However, the H/C  
10 ratio should not be overstated, as molecular electronic configuration and chemical bonding  
11 play an essential role in molecular properties.

1 SAF candidates also need to exhibit balanced properties. The optimal SAF candidates  
2 in Table 3 are not necessarily those showing the largest  $\text{NHOC}_G$  nor the largest  $\text{NHOC}_V$ , as it is  
3 often unlikely that the hydrocarbons with large  $\text{NHOC}_G$  also have large  $\text{NHOC}_V$  or vice versa.  
4 Achieving the optimal balance between  $\text{NHOC}_G$  and  $\text{NHOC}_V$  requires a holistic approach  
5 considering specific aircraft requirements, such as operational conditions and technological  
6 advancements. The present study employs a high energy density (HED) aviation fuel JP-10 as  
7 the fuel reference. The density of JP-10 is 0.940 g/ml, with a high  $\text{NHOC}_G$  of 42.200 MJ/kg and  
8 an  $\text{NHOC}_V$  of 39.640 MJ/L [10].

9 The  $\text{NHOC}_G$  of SAF candidates in Table 3 is above the reference (i.e., the  $\text{NHOC}_G$  of JP-  
10 10 42.200 MJ/kg). As a result, it is important to examine the density and  $\text{NHOC}_V$  properties of  
11 these SAF candidates. Fig. 6 plots  $\text{NHOC}_V$  and density of these PCHCs, which varies significantly  
12 from as low as 35.849 MJ/L to as high as 52.039 MJ/L due to the density variation. For new  
13 SAF candidates with higher density ( $>0.94$  g/ml) and larger  $\text{NHOC}_V$  (39.64 MJ/L), the preferred  
14 PCHCs need to be on the right-hand side of the vertical orange dash line and above the  
15 horizontal orange dash line (north-east or phase I). Up to 20 PCHCs in this region in Fig. 6 fit  
16 these criteria and, therefore, can be excellent candidates with superior density and  $\text{NHOC}_V$   
17 than JP-10 to proceed with the development.



1

2 Fig. 6. Relationship between NHOC<sub>v</sub> and density of SAF candidates in Table 3. The 20 preferred  
 3 candidates are No 1-No 19 and No 23 PCHCs.

4 The results show that properties such as NHOC<sub>v</sub> and density of PCHCs are more  
 5 sensitive to the molecular structure of the PCHCs than NHOC<sub>G</sub>. Further structure analysis of  
 6 these PCHCs reveals that the preferred PCHCs (25 out of 35) possess 3-5 rings. NHOC<sub>v</sub> and  
 7 density are likely related to the number of rings in the structures. For example, compound (No  
 8 1, C<sub>13</sub>H<sub>18</sub>) in Table 3 with 5 rings has the largest NHOC<sub>v</sub> of 52.039 MJ/L, whereas the majority  
 9 of the PCHCs possessing 2 rings are at the bottom of the list. JP-10 fuel (dominated by exo-  
 10 THDCPD C<sub>10</sub>H<sub>16</sub>) possesses 3 pentagon rings (See Structure 12 in Fig. 5) with a density of  
 11 0.990 g/ml, NHOC<sub>G</sub> 42.968 MJ/kg, and NHOC<sub>v</sub> 42.539 MJ/L, likewise THTCPD pentacyclic  
 12 (6.5.1.13,6.02,7.09,13) pentadecane (C<sub>15</sub>H<sub>22</sub>) possesses 5 pentagon rings analogy to the  
 13 structure of JP-10 (See Structure 4 in Fig. 5). The predicted properties for this compound 4 are  
 14 density 1.192 g/ml, NHOC<sub>G</sub> 42.701 MJ/kg, and NHOC<sub>v</sub> 50.900 MJ/L, which indicate that  
 15 THTCPD is a highly promising compound for SAF with a potential as an HED fuel for military  
 16 aircraft as well as the substitute for aromatic components in conventional aviation fuels.

### 1 **3.4. Structure-property relationship for hydrocarbons in SAF**

2 Understanding the structural characteristics of these PCHCs provides insight for the  
3 QSPR for future design and development of SAFs. Among the 35 PCHCs screened from the ML  
4 model, Structures 1- 11 in Table 3 exhibit superior  $NHOC_v$  and density properties than the  
5 HED JP-10 aviation fuel. Examination of the characteristics of these hydrocarbons in Table 3  
6 reveals that they share the following features:

- 7 1. Compact molecular structure: the obtained PCHCs have a compact molecular structure  
8 with multiple fused rings. This compactness allows for efficient packing of molecules,  
9 leading to higher energy density per unit volume.
- 10 2. Ratio of H/C:  $NHOC_v$  of PCHCs is predominantly determined by the ratio of H/C. Selection  
11 and design of new PCHCs for SAFs can be advanced by prioritizing these influential factors.  
12 However, the number of total atoms descriptors, which is not independent, may be  
13 removed for new descriptor development for SAF.
- 14 3. Multiple C-C bonds: PCHCs contain several C-C bonds within their fused ring structures.  
15 These bonds store large amounts of energy during combustion reactions, producing high  
16 heat release rates and enhanced energy output.
- 17 4. Ring strain: The presence of fused rings in PCHCs often leads to significant ring strain,  
18 which arises from the forced bending or distortion of C-C bonds to accommodate ring  
19 fusion. This strain imparts high reactivity to the molecule, facilitating rapid combustion  
20 and efficient energy release.
- 21 5. Saturation vs unsaturation: PCHCs may be both saturated and unsaturated. Saturation  
22 dominates the PCHCs and contributes to the stability and thermal resistance of the  
23 molecule, while unsaturation enhances reactivity and combustion efficiency.



- 1 6. Substituents: Functional groups or substituents attached to the polycyclic ring system can  
2 further modulate the properties of a hydrocarbon, such as polarity, solubility, and  
3 reactivity. For example, alkyl groups may increase the hydrophobicity and stability of the  
4 molecule, while polar functional groups may enhance interactions with other molecules  
5 or surfaces.
- 6 7. Steric hindrance: The three-dimensional (3D) arrangement of atoms of the compounds  
7 can introduce steric hindrance, affecting the molecule's interactions with surrounding  
8 molecules, surfaces, or catalysts. This can influence factors such as combustion kinetics,  
9 reaction rates, and product distributions.

#### 10 **4. Conclusions**

11 The study focused on developing a machine learning (ML) model that efficiently  
12 estimates critical fuel properties like net heat of combustion (NHOC) and hydrocarbon  
13 density. Using the supervised support vector machines (SVM) algorithm, models with six and  
14 forty descriptors were trained for NHOC and density, respectively, ensuring accuracy and  
15 reliability. These models were then applied to screen molecules from literature and database,  
16 identifying 35 high-energy density polycyclic hydrocarbon (PCHCs) molecules suitable for  
17 sustainable aviation fuel (SAF) applications. Interestingly, around 70% of these PCHCs  
18 exhibited NHOC<sub>v</sub> values comparable to or better than JP-10 jet fuel. Notably, the optimal  
19 PCHCs favored multiple rings with C-C single bonds and a high H/C ratio. However, this pre-  
20 screening step is just the beginning, as further steps involve developing quantitative  
21 structure-property relationships (QSPR), selecting or developing suitable catalysts for PCHCs  
22 synthesis, blending the PCHCs into aviation fuel, and assessing the impact on fuel properties  
23 according to ASTM specifications. Feasibility studies, techno-economic analyses, and  
24 environmental impact assessments are also crucial aspects of fuel development.

1 **CRedit authorship contribution statement**

2 **Dilip Rijal:** Writing – original draft, Data collection, Methodology, Editing. **Feng Wang:**  
3 Conceptualization, Supervision, Visualization, Writing – review, Editing. **Vladislav Vasilyev:**  
4 Review, Supervision, Methodology

5 **Acknowledgements**

6 DR acknowledges the Tuition Fee Scholarship of the Swinburne University of Technology. This  
7 research did not receive any specific grant from funding agencies in the public, commercial,  
8 or not-for-profit sectors.

9 **Declaration of competing interest**

10 The authors declare that they have no known competing financial interests or personal  
11 relationships that could have appeared to influence the work reported in this paper.

12 **Appendix A. Supplementary materials**

13 **References**

14 [1] International Civil Aviation Organization (ICAO). *Environmental Trends in Aviation to*  
15 *2050*. [https://www.icao.int/environmental-](https://www.icao.int/environmental-protection/Documents/EnvironmentalReports/2022/ENVReport2022_Art7.pdf)  
16 [protection/Documents/EnvironmentalReports/2022/ENVReport2022\\_Art7.pdf](https://www.icao.int/environmental-protection/Documents/EnvironmentalReports/2022/ENVReport2022_Art7.pdf) (accessed  
17 March 14, 2024).

18 [2] Wang, F.; Rijal, D. Sustainable Aviation Fuels for Clean Skies: Exploring the Potential and  
19 Perspectives of Strained Hydrocarbons. *Energy & Fuels* **2024**. DOI:  
20 <https://doi.org/10.1021/acs.energyfuels.3c04935>

21 [3] Ince, A. C.; Colpan, C. O.; Hagen, A.; Serincan, M. F. Modeling and simulation of Power-to-  
22 X systems: A review. *Fuel* **2021**, *304*, 121354. DOI:  
23 <https://doi.org/10.1016/j.fuel.2021.121354>

- 1 [4] Heyne, J.; Rauch, B.; Le Clercq, P.; Colket, M. Sustainable aviation fuel prescreening tools  
2 and procedures. *Fuel* **2021**, *290*, 120004. DOI: <https://doi.org/10.1016/j.fuel.2020.120004>
- 3 [5] Landera, A.; Bambha, R. P.; Hao, N.; Desai, S. P.; Moore, C. M.; Sutton, A. D.; George, A.  
4 Building structure-property relationships of cycloalkanes in support of their use in sustainable  
5 aviation fuels. *Frontiers in Energy Research* **2022**, *9*, 771697. DOI:  
6 <https://doi.org/10.3389/fenrg.2021.771697>
- 7 [6] Kosir, S.; Heyne, J.; Graham, J. A machine learning framework for drop-in volume swell  
8 characteristics of sustainable aviation fuel. *Fuel* **2020**, *274*, 117832. DOI:  
9 <https://doi.org/10.1016/j.fuel.2020.117832>
- 10 [7] TB, G. H. Aviation Fuels Technical Review. **2007**. [https://www.chevron.com/-](https://www.chevron.com/-/media/chevron/operations/documents/aviation-tech-review.pdf)  
11 [/media/chevron/operations/documents/aviation-tech-review.pdf](https://www.chevron.com/-/media/chevron/operations/documents/aviation-tech-review.pdf) (accessed Jan 5, 2024).
- 12 [8] Zou, J.-J.; Zhang, X.; Pan, L. *High-energy-density fuels for advanced propulsion: Design and*  
13 *synthesis*; John Wiley & Sons, 2020.
- 14 [9] Nie, J.; Jia, T.; Pan, L.; Zhang, X.; Zou, J.-J. Development of high-energy-density liquid  
15 aerospace fuel: a perspective. *Transactions of Tianjin University* **2022**, *28* (1), 1-5. DOI:  
16 <https://doi.org/https://doi.org/10.1007/s12209-021-00302-x>
- 17 [10] Woodroffe, J.-D.; Harvey, B. G. High-performance, biobased, jet fuel blends containing  
18 hydrogenated monoterpenes and synthetic paraffinic kerosene's. *Energy & Fuels* **2020**, *34* (5),  
19 5929-5937. DOI: <https://doi.org/10.1021/acs.energyfuels.0c00274>
- 20 [11] Jessup, R. S. *Precise measurement of heat of combustion with a bomb calorimeter*; US  
21 Department of Commerce, National Bureau of Standards, 1960.
- 22 [12] Rojas-Aguilar, A. An isoperibol micro-bomb combustion calorimeter for measurement of  
23 the enthalpy of combustion. Application to the study of fullerene C60. *The Journal of Chemical*

1 *Thermodynamics* **2002**, *34* (10), 1729-1743. DOI: <https://doi.org/10.1016/S0021->  
2 [9614\(02\)00257-4](https://doi.org/10.1016/S0021-9614(02)00257-4)

3 [13] Alibakhshi, A. High precision evaluation of the combustion enthalpy by ab-intio  
4 computations. **2021**. DOI: <https://doi.org/10.26434/chemrxiv-2021-fvcph>

5 [14] Yang, H.; Yang, Z.-J.; Yang, Q.-F.; Wei, X.-M.; Yuan, Y.-Q.; Wang, L.-L.; Hu, Y.-F.; Ding, J.-J.  
6 Simple and high-precision DFT-QSPR prediction of enthalpy of combustion for  
7 sesquiterpenoid high-energy-density fuels. *Fuel* **2023**, *332*, 126157. DOI:  
8 <https://doi.org/10.1016/j.fuel.2022.126157>

9 [15] Frutiger, J.; Marcarie, C.; Abildskov, J.; Sin, G. r. A Comprehensive Methodology for  
10 Development, Parameter Estimation, and Uncertainty Analysis of Group Contribution Based  
11 Property Models□An Application to the Heat of Combustion. *Journal of Chemical &*  
12 *Engineering Data* **2016**, *61* (1), 602-613. DOI: <https://doi.org/10.1021/acs.jced.5b00750>

13 [16] Marrero, J.; Gani, R. Group-contribution based estimation of pure component properties.  
14 *Fluid phase equilibria* **2001**, *183*, 183-208. DOI: <https://doi.org/10.1016/S0378->  
15 [3812\(01\)00431-9](https://doi.org/10.1016/S0378-3812(01)00431-9)

16 [17] Abdul Jameel, A. G.; Al-Muslem, A.; Ahmad, N.; Alquaity, A. B.; Zahid, U.; Ahmed, U.  
17 Predicting enthalpy of combustion using machine learning. *Processes* **2022**, *10* (11), 2384.  
18 DOI: <https://doi.org/10.3390/pr10112384>

19 [18] Li, R.; Herreros, J. M.; Tsolakis, A.; Yang, W. Machine learning-quantitative structure  
20 property relationship (ML-QSPR) method for fuel physicochemical properties prediction of  
21 multiple fuel types. *Fuel* **2021**, *304*, 121437. DOI: <https://doi.org/10.1016/j.fuel.2021.121437>

22 [19] Albahri, T. A. Method for predicting the standard net heat of combustion for pure  
23 hydrocarbons from their molecular structure. *Energy conversion and management* **2013**, *76*,  
24 1143-1149. DOI: <https://doi.org/10.1016/j.enconman.2013.09.019>

- 1 [20] API. Technical data book—petroleum refining. *Metric Edition, Vols 1997, 1* (2).
- 2 [21] Rumble, J. CRC handbook of chemistry and physics. **2017**.
- 3 [22] DAYLIGHT. SMILES. <https://www.daylight.com/dayhtml/doc/theory/theory.smiles.html>
- 4 (accessed May 24, 2024).
- 5 [23] PubChem. *Explore Chemistry*. <https://pubchem.ncbi.nlm.nih.gov/> (accessed Jan 26,
- 6 2024).
- 7 [24] O'Boyle, N. M.; Banck, M.; James, C. A.; Morley, C.; Vandermeersch, T.; Hutchison, G. R.
- 8 Open Babel: An open chemical toolbox. *Journal of cheminformatics* **2011**, *3* (1), 1-14. DOI:
- 9 <https://doi.org/10.1186/1758-2946-3-33>
- 10 [25] Aires de Sousa, J. GUIDEMOL: a Python graphical user interface for molecular descriptors
- 11 based on RDKit. *Molecular Informatics* **2023**. DOI: <https://doi.org/10.1002/minf.202300190>
- 12 [26] The RDKit documentation. <https://www.rdkit.org/docs/index.html> (accessed Feb 25,
- 13 2024).
- 14 [27] Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.;
- 15 Prettenhofer, P.; Weiss, R.; Dubourg, V. Scikit-learn: Machine learning in Python. *the Journal*
- 16 *of machine Learning research* **2011**, *12*, 2825-2830.
- 17 [28] Shi, C.; Borchardt, T. B. JRgui: A Python program of Joback and Reid method. *ACS omega*
- 18 **2017**, *2* (12), 8682-8688. DOI: <https://doi.org/10.1021/acsomega.7b01464>
- 19 [29] Duarte Ramos Matos, G.; Pak, S.; Rizzo, R. C. Descriptor-Driven de Novo Design
- 20 Algorithms for DOCK6 Using RDKit. *Journal of Chemical Information and Modeling* **2023**, *63*
- 21 (18), 5803-5822. DOI: <https://doi.org/10.1021/acs.jcim.3c01031>
- 22 [30] Moubayed, A.; Injadat, M.; Nassif, A. B.; Lutfiyya, H.; Shami, A. E-learning: Challenges and
- 23 research opportunities using machine learning & data analytics. *IEEE Access* **2018**, *6*, 39117-
- 24 39138. DOI: <https://doi.org/10.1109/ACCESS.2018.2851790>

- 1 [31] Soman, K.; Loganathan, R.; Ajay, V. *Machine learning with SVM and other kernel methods*;  
2 PHI Learning Pvt. Ltd., 2009.
- 3 [32] Shaik, A. B.; Srinivasan, S. A brief survey on random forest ensembles in classification  
4 model. In *International Conference on Innovative Computing and Communications: Proceedings of ICICC 2018, Volume 2*, 2019; Springer: pp 253-260. DOI:  
5 [https://doi.org/10.1007/978-981-13-2354-6\\_27](https://doi.org/10.1007/978-981-13-2354-6_27)
- 6  
7 [33] Steinbach, M.; Tan, P.-N. *kNN: k-nearest neighbors*; 2009.
- 8 [34] Sukhdeve, D. S. R.; Sukhdeve, S. S. Google Colaboratory. In *Google Cloud Platform for*  
9 *Data Science: A Crash Course on Big Data, Machine Learning, and Data Analytics Services*,  
10 Springer, 2023; pp 11-34.
- 11 [35] Martelli, A.; Ravenscroft, A. M.; Holden, S.; McGuire, P. *Python in a Nutshell*; " O'Reilly  
12 Media, Inc.", 2023.
- 13 [36] 10 Essential Data Science Package for Python-Kite Blog.  
14 <https://www.kite.com/blog/python/data-science-packages-python/> (accessed Feb 25, 2024).
- 15 [37] Montgomery, D. C. *Design and analysis of experiments*; John Wiley & sons, 2017.
- 16 [38] Baykan, N. A.; Yılmaz, N. A mineral classification system with multiple artificial neural  
17 network using k-fold cross validation. *Mathematical and Computational Applications* **2011**, *16*  
18 (1), 22-30. DOI: <https://doi.org/10.3390/mca16010022>
- 19 [39] Montesinos López, O. A.; Montesinos López, A.; Crossa, J. Overfitting, model tuning, and  
20 evaluation of prediction performance. In *Multivariate statistical machine learning methods*  
21 *for genomic prediction*, Springer, 2022; pp 109-139.
- 22 [40] Otchere, D. A.; Ganat, T. O. A.; Gholami, R.; Ridha, S. Application of supervised machine  
23 learning paradigms in the prediction of petroleum reservoir properties: Comparative analysis

1 of ANN and SVM models. *Journal of Petroleum Science and Engineering* **2021**, *200*, 108182.  
2 DOI: <https://doi.org/10.1016/j.petrol.2020.108182>

3 [41] CHEMICAL COMPUTING GROUP. *QuaSAR-Descriptor*  
4 <https://cadaster.eu/sites/cadaster.eu/files/challenge/descr.htm> (accessed Feb 28, 2024).

5 [42] Ladommatos, N.; Rubenstein, P.; Bennett, P. Some effects of molecular structure of single  
6 hydrocarbons on sooting tendency. *Fuel* **1996**, *75* (2), 114-124. DOI:  
7 [https://doi.org/10.1016/0016-2361\(94\)00251-7](https://doi.org/10.1016/0016-2361(94)00251-7)

8 [43] Kittel, H.; Horský, J.; Šimáček, P. Properties of Selected Alternative Petroleum Fractions  
9 and Sustainable Aviation Fuels. *Processes* **2023**, *11* (3), 935. DOI:  
10 <https://doi.org/10.3390/pr11030935>

11 [44] Huq, N. A.; Hafenstine, G. R.; Huo, X.; Nguyen, H.; Tifft, S. M.; Conklin, D. R.; Stück, D.;  
12 Stunkel, J.; Yang, Z.; Heyne, J. S. Toward net-zero sustainable aviation fuel with wet waste-  
13 derived volatile fatty acids. *Proceedings of the National Academy of Sciences* **2021**, *118* (13),  
14 e2023008118. DOI: <https://doi.org/10.1073/pnas.2023008118>

15 [45] Li, G.; Hu, Z.; Hou, F.; Li, X.; Wang, L.; Zhang, X. Machine learning enabled high-throughput  
16 screening of hydrocarbon molecules for the design of next generation fuels. *Fuel* **2020**, *265*,  
17 116968. DOI: <https://doi.org/10.1016/j.fuel.2019.116968>

18 [46] Lin, J.-K.; Abi Nurazaq, W.; Wang, W.-C. The properties of sustainable aviation fuel I: Spray  
19 characteristics. *Energy* **2023**, *283*, 129125. DOI:  
20 <https://doi.org/10.1016/j.energy.2023.129125>

21  
22  
23  
24

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15