

Clemens Hall^{1,*}, Benoit Creton², Bastian Rauch¹, Uwe Bauder¹, Manfred Aigner¹, Probabilistic Mean Quantitative Structure-Property Relationship modelling of Jet Fuels Properties, *Energy Fuels* 2022, 36, 1, 463–47

¹DLR, German Aerospace Center, Institute of Combustion Technology, MAT, 70569 Stuttgart, Germany.

²IFP Energies nouvelles, 1 et 4 avenue de Bois-Préau, 92852 Rueil-Malmaison, France.

This document is the Accepted Manuscript version of a Published Work that appeared in final form in *Energy & Fuels* copyright © American Chemical Society after peer review and technical editing by the publisher. To access the final edited and published work see

<https://pubs.acs.org/doi/10.1021/acs.energyfuels.1c03334>

Probabilistic Mean Quantitative Structure-Property Relationship modelling of Jet Fuels Properties

Clemens Hall^{1,}, Benoit Creton², Bastian Rauch¹, Uwe Bauder¹, Manfred Aigner¹*

¹DLR, German Aerospace Center, Institute of Combustion Technology, MAT, 70569 Stuttgart, Germany.

²IFP Energies nouvelles, 1 et 4 avenue de Bois-Préau, 92852 Rueil-Malmaison, France.

*Corresponding Author Email: Clemens.hall@dlr.de

KEYWORDS

Aviation Fuel, Synthetic Fuel, Machine Learning, Modelling, Mean Quantitative Structure Property Relationship, Quantitative Structure Property Relationship, Uncertainty Quantification, Predictive Capability, Monte-Carlo Dropout Neural Network Regression

ABSTRACT

We present a novel probabilistic Mean Quantitative Structure-Property Relationship (M-QSPR) method for the prediction of jet fuel properties considering two-dimensional gas chromatography measurements. Fuels are represented as one mean pseudo-structure that is inferred by a weighted average over structures of 1866 molecules that could be present in the individual fuel. The method allows the training of models on both data of pure components and of fuels and does not require mixing rules for the calculation of the bulk property. This drastically increases the number of available training data and allows the direct learning of the mixing behavior. For the modelling we use a Monte-Carlo dropout neural network, a probabilistic Machine Learning algorithm, that estimates prediction uncertainties due to possible unidentified isomers and dissimilarity of training and test data. Models are developed to predict the freezing point, flash point, net heat of combustion, and temperature dependent properties such as density, viscosity, and surface tension. We investigate the effect of the presence of fuels in the training data on the predictions for up to 82 conventional fuels and 50 synthetic fuels. The results of the predictions are compared on three metrics that quantify accuracy, precision and reliability. These metrics allow a comprehensive estimation of the predictive capability of the models. For the prediction of the density, surface tension and net heat of combustion the M-QSPR method yield highly accurate results even without the presence of fuels in the training data. For properties with non-linear behavior over temperature and complex fuel component interactions, like viscosity and freezing point, the presence of fuels in the training data was found to be essential for the method.

1. Introduction

To facilitate the development and approval of new jet fuel candidates, model-based prediction methods were identified as key enabler for fuel screening by research projects like the EU project for Jet Fuel Screening and Optimization (JETSCREEN) ¹ and the National Jet Fuels Combustion Project (NJFCP) ². Reliable prediction models solely based on the fuel composition could significantly reduce time and cost for necessary measurements that are needed to ensure safe application ³. Especially at early stages of fuel production process development, the fuel volumes required for the approval process regulated by the ASTM 4054 ⁴ often exceed the production capabilities. The application of model-based prediction methods can greatly support the decision making at this stage of the development process. Modern analytical methods like two-dimensional gas chromatography (GCxGC) are able to provide detailed characterization of the chemical composition of a fuel with volumes below 1 mL, which can be utilized as input for the modelling ⁵. The GCxGC method has been adopted in the new ASTM D4054 fast track certification process and was chosen as data basis for this work. For modelling we selected the properties density, kinematic viscosity, surface tension, flash point and freezing point, which are critical parameters for fuel specifications, relevant for the screening and development of jet fuels ^{1,2}.

The demands for the predictive capability and robustness of the models over the extensive input domain of possible future fuels are high for safety relevant use cases like jet fuel approval or combustor design ³. Predictions with significant deviations from the true property value are unacceptable and have to be prevented by utilizing robust modelling techniques that are accurate enough to predict properties of fuel candidates that might significantly differ from the training and validation data. There exist various methods for the modelling of jet fuels on the basis of GCxGC data. The methods range from direct correlation of the GCxGC measurement as input with the property desired output ^{5,6} to methods that approximate the fuel as mixture of representative molecules like the Quantitative Structure-Property Relationship (QSPR) method ⁷. The QSPR approach has been successfully applied for the modelling of properties of possible jet fuel compounds like density and viscosity ⁸, freezing point and net heat of combustion ⁹ as well as the flash point ¹⁰. The method represents fuels as mixtures of pure compounds where property values of fuel's component are predicted as individual species. The bulk property of the fuel is calculated afterwards, based on the pure compound values and mass fractions, within a mixing rule ¹¹. The method requires measured property data of pure compounds, which is available in various databases like the NIST Standard Reference Database 103a ¹² and DIPPR 801 ¹³. Data of pure compounds is thereby available for almost every molecular family, with a substantial number of representative molecules, which allows the modelling of almost every fuel composition. Compared to the method of direct correlation of the GCxGC with property measurement,

this allows for great flexibility and extensibility of the QSPR method to compositions which differ significantly from known fuels. The method could therefore cover the necessary domain of new potential fuel candidates. The correct selection of the molecules present in the fuel as well as an adequate mixing rule however pose a serious challenge. A table listing advantages and disadvantages of the direct correlation method, the QSPR method and the method presented in this work is provided in the supplementary material S1.

Jet fuels are highly complex mixtures composed of hundreds of molecules, with varying composition based on the production location and pathway. The exact identification of each molecule in the fuel is presently not possible with the presently utilized GCxGC methods (mass spectroscopy and flame ionization detectors) ^{5,14}. Rather than identifying the exact molecules in the fuel the GCxGC method allows classification of detected species in bins with respect to their molecular family and the number of C atoms present. To calculate fuel properties on the basis of the GCxGC measurement, the fuel needs to be approximated. In the current literature there exist three different approaches for this approximation: (1) selecting one representative molecule for each bin of the GCxGC measurement ¹¹ (2) sampling from a selection of possible representative molecules ¹⁵ and (3) finding a surrogate composition with a limited number of compounds that match the specific properties of the fuel (structural or chemo-physical) ¹⁶⁻¹⁸. Depending on the use case and the modelled property, approach (1) might be an oversimplification. The reason for this is the high number of possible isomers in a GCxGC bin, and the variance in the molecular structures. For properties that are less affected by branching effects of the molecular structure approach (1) might yield sufficient results, whereas for other properties significant deviations can occur. Approach (2) considers multiple possible isomers in one GCxGC bin, but the sampling method makes the method computationally expensive and time consuming. Approach (3) either requires additional property measurements or a selection of representative molecules for each GCxGC bin similar to approach (1) if the surrogate is formulated by structural similarity, as presented by Ri et al. ^{16,17}. However, approaches (1), (2), and (3) predict properties for individual compounds of the fuel. To calculate the bulk property of the fuel an appropriate mixing rule is needed for all three approaches. None of the approaches allows the training of the property model on data from both pure compounds and fuel measurements to directly infer the mixing behavior.

We present an approach that solves these issues of the consideration of multiple possible isomers for GCxGC bins as well as the need for an adequate mixing rule. This is achieved by transferring the concept of Ajmani et. al. of a mean mole weighted average quantitative descriptor ¹⁹ to calculate a pseudo-structure of a fuel and correlating this representation by probabilistic Machine Learning models with the properties. This approach allows the training on both data from pure compounds and fuels. The

concept of weighted average descriptor has already been successfully applied for mixtures with only a few compounds e.g. for the prediction of flash point ^{10,20,21}. We transfer this approach to jet fuels, mixtures of hundreds of possible molecules, with isomers that are not further identified by the GCxGC composition measurement.

We investigate the ability of this method to predict properties of jet fuels for models trained solely on data from pure compounds as well as for models trained on both data from pure compounds and fuels. We thereby systematically investigate the effect of presence of fuels in the training data to infer the mixing behavior for the different properties. The effects of different training datasets on the predictive capability of the models for the prediction of up to 82 crude oil-based conventional fuels and 50 synthetic fuels, produced from alternative production pathways, are assessed using a recently introduced concept ⁶. This concept assesses model predictions with respect to quantitative metrics of their accuracy, their validity and their precision. To cover a broad range of the possible application domain, the models are tested on conventional crude oil-based jet fuels as well as synthetic jet fuels and blends, as well as research fuels.

2. Datasets

For this work, compositional data of both pure compounds and jet fuels are considered to derive correlations with property measurements. Pure compounds are encoded using the Simplified Molecular-Input Line-Entry System (SMILES) key, a string representation of the molecular structure. For the fuels, the composition is represented by GCxGC measurements. Both compositional representations are transferred to a quantitative structural representation as explained later on.

2.1 Jet fuel data

The GCxGC measurement describes jet fuel composition as a matrix of bins of hydrocarbon family and carbon number. For this work species in a range of 1-25 carbon atoms from seven different hydrocarbon families: n-alkanes, iso-alkanes, mono-cyclo-alkanes, bi-cyclo-alkanes, mono-aromatics, cyclo-aromatics and di-aromatics are considered. Tri-cyclo-alkanes are lumped into the group of bi-cyclo-alkanes. To consider different kinds of fuels, a dataset from the DLR jet fuel database is utilized containing measurements of up to 82 conventional fuels and 50 synthetic fuels. Crude oil-based fuels with the jet fuel types Jet A, Jet A-1, JP-5, JP-8 and TS-1 are grouped as conventional fuels. The conventional fuels mainly originate from the Coordinating Research Council (CRC) world fuel survey of 2006 ²². The synthetic fuels, produced by processes that are not based on crude oil, were systematically gathered in the DLR jet fuel database from research projects like JETSCREEN ¹, Emission and Climate Impact of Alternative Fuels (ECLIF) ²³ and NJFCP ² as well as technical, screening and certification

reports. These synthetic fuels contain fuels produced by the Alcohol To Jet process (ATJ), the Fischer-Tropsch process (FT), from Hydroprocessed Esters and Fatty Acids (HEFA), the Integrated Hydrolysis process of Shell with a high cyclo-alkanes content Cyclo Paraffinic Kerosene (CPK), the Catalytic Hydrothermal Conversion Jet fuel process (CHCJ) as well as fuel blends. To visualize the variety in fuels' compositions *Figure 1* shows scatter plots for the mass fraction of each chemical family for both the conventional fuels (Conv.) in blue and the synthetic fuels (Syn.) in green. Blue and green shaded areas illustrate the possible range of each group.

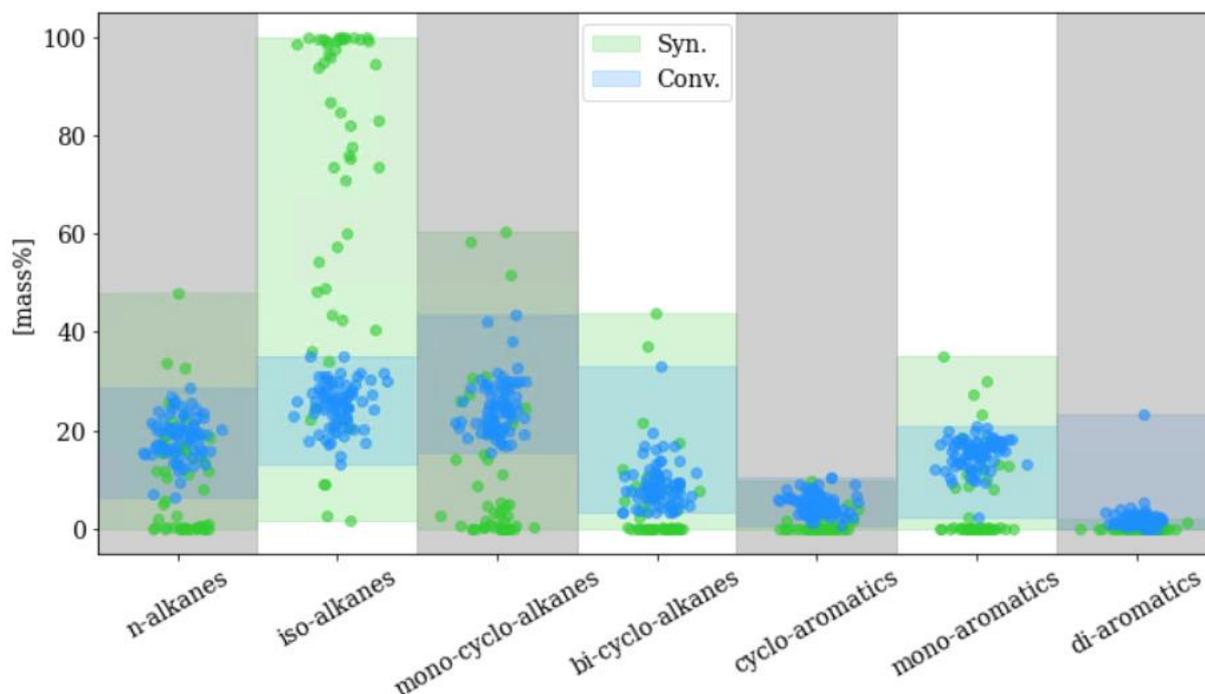


Figure 1: Scatter plot of GCxGC measurements of fuels summed up to hydrocarbon families. Blue: conventional fuels, green: synthetic fuels and blends. Blue and green shaded areas indicate the possible range.

The figure shows the variety of fuels' compositions and the possible range of hydrocarbon fractions of jet fuels. With the exception of di-aromatics, the synthetic fuels dataset (including blends) has a broader and more distributed range of compositions, whereas conventional fuels have a very distinctive compositional range. This is due to the large variety of different feedstock and production pathways of synthetic jet fuels and the lower number of hydrocarbon components compared to conventional ones. Predictive models developed mainly for conventional fuels, may have a challenge in correctly predicting over the vast space of possible fuel compositions.

2.2 Pure compounds data

To be able to approximate the fuels with representative molecules, a database with 1866 pure compounds was built collecting data available at the NIST ThermoData Engine with the NIST Standard Reference Database 103a¹², DIPPR 801¹³, Chemspider²⁴ and Pubchem²⁵. The classification of the

molecules in the different matrix bins is carried out based on the molecular formula and characteristic molecular substructures. Further details are given in the supplementary material S2.

Since the exact molecules in the fuels' compositions are unknown, we assume that the selection from the databases is representative for all considered fuel types and properties. In order to review the validity of this assumption, we compare the selection from the databases with all possible isomers for the different GCxGC bins, based on the molecular formula. To generate all possible isomers, we utilized the molecule generator MOLGEN v. 5²⁶ that generated all possible molecules for a given molecular formula of a hydrocarbon family (e.g. C_nH_{2n} for mono-cyclo-alkanes for n 1-25). For The generated molecules were subsequently classified using characteristic molecular substructures outlined in the supplementary material S2. This is a worst-case approximation since the range of possible molecules that could actually exist in conventional and synthetic fuels might consist of significantly smaller number of molecules. Molecules generated by MOLGEN that could theoretically exist based on the molecular formula might never occur in a fuel due to limitations given by the production process. *Table 1* shows the comparison of the number of possible molecules calculated by MOLGEN (column MG) and the number of molecules from our databases (column DB). Due to the exponential increase of possible isomers and computational limitations for the classification, only calculations up to molecules with 12 C atoms were executed. For iso-alkanes the numbers above C 12 were taken from the theoretical calculations²⁷.

We furthermore review the validity of our molecule selection based on a comparison of the chemical spaces of the QSPR features that we introduce in 3.2 Chemical space and application domain comparison. The percental difference of the area of the chemical spaces spanned by the molecules generated by MOLGEN and our selection for all families at a carbon number of C 12 are provided in the supplementary material S3. To give an exemplary visual comparison of the chemical spaces we also provide a parallel line plot for cyclo-aromatics in the supplementary material S4 with filled areas between the minimum and the maximum value for a QSPR feature.

C number	n-alkanes		iso-alkanes		mono-cyclo-alkanes		bi-cyclo-alkanes		cyclo-aromatics		mono-aromatics		di-aromatics	
	DB	MG	DB	MG	DB	MG	DB	MG	DB	MG	DB	MG	DB	MG
1	1	1												
2	1	1												
3	1	1			1	1								
4	1	1	1	1	2	2	0	1						
5	1	1	2	2	5	5	2	4						
6	1	1	4	4	12	12	5	14	0	4	1	1		
7	1	1	8	8	27	29	5	46	0	9	1	1		
8	1	1	17	17	43	73	8	150	0	31	4	4		
9	1	1	32	34	52	185	18	477	2	75	8	8		
10	1	1	49	74	58	475	18	1503	8	218	22	22	0	5
11	1	1	37	158	33	1231	12	4680	16	588	42	51	2	21
12	1	1	45	354	29	3232	18	14461	27	1657	60	136	9	103
13	1	1	34	801	14	>3232	15	>14461	32	>1657	51	>136	24	>103
14	1	1	29	1857	17	>3232	16	>14461	31	>1657	60	>136	38	>103
15	1	1	27	4346	18	>3232	14	>14461	29	>1657	34	>136	17	>103
16	1	1	32	10358	8	>3232	15	>14461	16	>1657	22	>136	9	>103
17	1	1	11	24893	6	>3232	8	>14461	7	>1657	15	>136	3	>103
18	1	1	12	60522	8	>3232	17	>14461	6	>1657	22	>136	13	>103
19	1	1	14	147283	13	>3232	6	>14461	12	>1657	11	>136	4	>103
20	1	1	20	366318	4	>3232	10	>14461	9	>1657	13	>136	9	>103
21	1	1	9	>366317	6	>3232	2	>14461	0	>1657	3	>136	1	>103
22	1	1	14	>366317	5	>3232	4	>14461	6	>1657	14	>136	4	>103
23	1	1	5	>366317	3	>3232	0	>14461	0	>1657	4	>136	0	>103
24	1	1	17	>366317	4	>3232	3	>14461	1	>1657	9	>136	0	>103
25	1	1	2	36797587	4	>3232	5	>14461	1	>1657	7	>136	1	>103

Table 1: Comparison of the number of representative molecules available in the database (DB) and the number of theoretically possible molecules, calculated by MOLGEN (MG)

Table 1 illustrates the exponential increase of possible isomers for jet fuels on the basis of GCxGC data. Compared to the possible molecules for C numbers above 12 only a fraction can be considered in this work, since no further data is available at the listed sources at the current time. This again illustrates the complexity of the property modelling based on GCxGC data. For properties, that are less affected by branching effects of isomers e.g. density, this might be less significant for the modelling since representative molecules have similar property values. For properties however, that are heavily affected by branching effects of isomers the variance of the property values of a GCxGC bin might result in significant deviation and uncertainty in the prediction. This could raise the demand for an even more detailed form of fuel compositional characterization and representation down to the molecular level.

2.3 Property data

For the evaluation and approval of jet fuels, values of the density are necessary at 15 °C, of the surface tension at 22 °C and of the kinematic viscosity at -40 °C and -20 °C at a pressure of 1 atm. Flash point, freezing point and net heat of combustion are measured at laboratory standard conditions. In the scope of this work the values utilized for the training and validation data of density, kinematic viscosity and surface tension are used from extended temperature range of -40 to 25 °C. This was done for reasons of availability of the data, since most of the available pure compound data of the mentioned data sources is measured at 20 °C and 25 °C. Because most of the viscosity data is available as dynamic

viscosity the measurements were divided by density measurements at a temperature ± 0.1 °C of the viscosity measurement. For viscosity measurements at temperatures with no corresponding density measurement, the trained density models from this work were used to compute a density value at the corresponding temperature. Because the mean relative error of the predicted density is around 1 %, as will be later shown, the potential error is assumed to be neglectable.

A detailed summary about the number of data points for pure compounds, conventional and synthetic fuels for each property can be found in *Table 2*. Column #FL (left sub column) lists the number of unique fuels and column #DP (right sub column) the total number of datapoints.

Fuel type	Density		Kinematic viscosity		Surface tension		Flash point		Freezing point		Net heat of combustion	
	#FL	#DP	#FL	#DP	#FL	#DP	#FL	#DP	#FL	#DP	#FL	#DP
Pure compounds	1866	27481	1090	4003	279	1494	273	732	383	1120	1495	1802
Conventional jet fuels	82	485	80	204	67	181	78	78	76	76	80	83
Synthetic jet fuels	50	185	47	91	13	40	44	46	40	42	37	43
All	1998	28151	1217	4298	359	1715	395	856	499	1238	1612	1928

Table 2: Number of fluids and data points for each fuel type and the respective fuel property

To remove outliers the dataset for each property was processed with the modified Z-score by Iglewicz and Hoaglin, an outlier detection method for small sample sizes as recommended by NIST²⁸. Datapoints with a Z-score greater 3.5 were removed from the dataset. For the temperature dependent properties the outlier detection was carried out for measurements in steps 0.1 °C.

2.4 Partitioning of correlation data

To systematically investigate the effects of the presence of fuels in the training data on the predictive capability of the models to predict properties of jet fuels, we consider three different training datasets resulting in three different models for each property: (1) a dataset of pure compounds (Base dataset), (2) a dataset of pure compounds and conventional fuels (Fuel dataset 1) and (3) a dataset consisting of data from pure compounds, conventional and synthetic fuels (Fuel dataset 2). By comparing the predictions for these three datasets, we assess the effect of the presence of jet fuels in the training data to approximate the mixing behavior. Furthermore, we test the ability of the models to correctly estimate uncertainty due to different isomers in the GCxGC bin for both conventional and synthetic fuels. By assessing the predictions of the Base model, we additionally investigate if the M-QSPR method

[CX3H1R]	[cH0]	[CX2H1]#[CX2H0]
[CX3H0R]	[C]=[C](!CX1)!CX1	[CX2H0]#[CX2H0]
[cX3H1](:*):*	[C]=[C]([C])!C	[R]
[cX3H0](:*)(:*)*	[C]=[C]([C])[C]	

Table 3: Utilized SMARTS codes of molecular features for M-QSPR and QSPR modelling

To quantify chosen molecular features we used the RDKit Python package, one of this package's functionalities stands in counting occurrences of substructures based on the SMILES³². An example for the quantitative structures of a molecule is given in Figure 2 for 2,3-hydro-2-methyl-1H-idene. The number behind the SMARTS key shows the count of the molecular feature and describes the occurrences of a functional group in the molecule; e.g. for the feature [CX4H3] the count of 1 indicates, that the molecule contains 1 methyl-group, the count of 9 for the feature [R] indicates that the molecule furthermore contains 9 ring atoms etc.

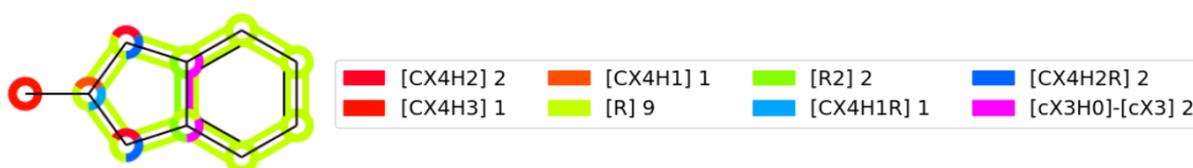


Figure 2: Quantified molecular features of 2,3-hydro-2-methyl-1h-idene, number behind SMART key shows count of molecular feature (For reasons of clarity only features that do appear in the properties are shown in the legend).

In the scope of this work only quantitative descriptions of substructures are considered for the correlation to confine the number of possible input features. In future work the effect of chemophysical descriptors as input features like the molar weight or the Van-der-Waals volume should be investigated.

3.1 Calculation of M-QSPR Representation

To calculate the M-QSPR representation of a jet fuel we compute the average occurrence of the considered quantitative structures in the fuel. Since the fuel composition is given as GCxGC measurement we calculate the mean occurrence of a structural features in a GCxGC bin (molecular family and number of C-atoms). We therefore average the occurrence of a quantitative structure for all molecules classified to the GCxGC bin. This process is repeated for all considered features in Table 3, which creates a matrix that lists the mean occurrence of the quantitative structures in each GCxGC bin. This mean occurrence matrix is multiplied with the molar fractions of the GCxGC measurements, summing up the values of the substructures for all GCxGC bins to compute the M-QSPR representation of the fuel as a vector of the considered quantitative substructures. Figure 3 schematizes this process for the methyl groups ([CX4H3]) of iso-alkanes with 8 carbon atoms.

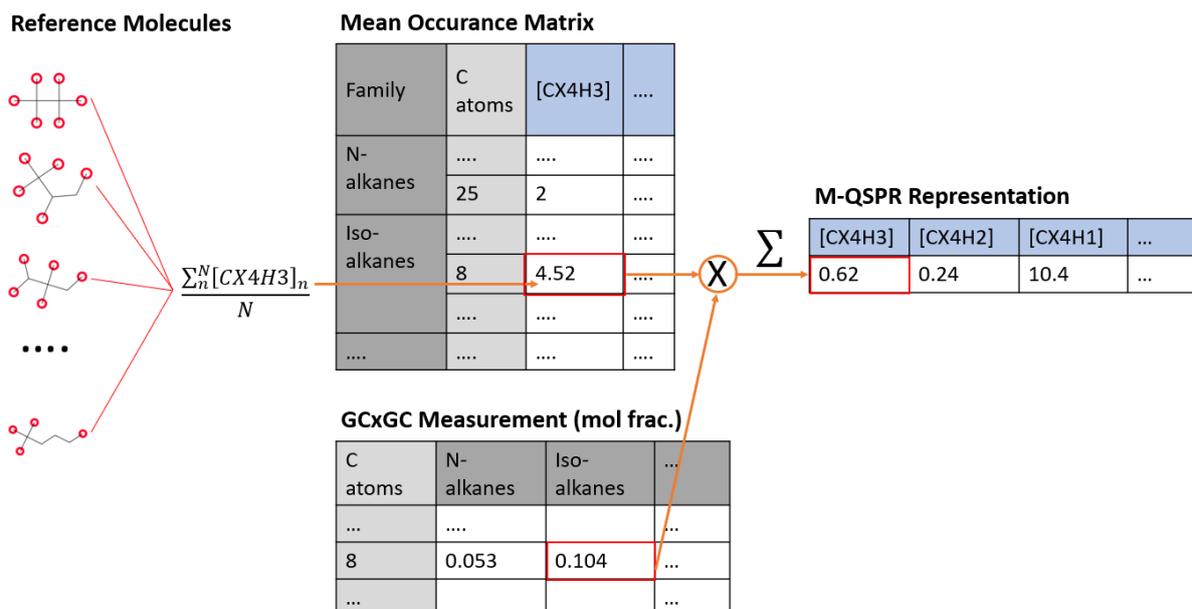


Figure 3: Schematic figure of mean occurrence matrix calculation from GCxGC with molar fractions

In the scope of this work we assume, that the mean occurrence matrix is representative for all considered conventional and synthetic fuels, regardless of their production pathway. If further knowledge of possible molecules exists for a production pathway, a unique mean occurrence matrix could be calculated from those molecules.

For GCxGC bins of fuels that do not appear in the mean occurrence matrix, because no representative molecules are available in the utilized pure compound database, see Table 1, the values of the mean occurrence of bins with a carbon number plus or minus one are used. E.g. for cyclo-aromatics C8, molecules from cyclo-aromatics C9 are used.

3.2 Chemical space and application domain comparison

The M-QSPR pseudo-structure allows the representation of fuels in the chemical space of the considered molecular features with dimensions for each feature. The unified representation raises the question for a comparison of the chemical spaces spanned by the considered dataset and thereby a comparison of the application domain of the three datasets outlined in 2.4 Partitioning of correlation data. To allow a visual comparison of the chemical spaces of the considered pure compounds and the M-QSPR representation of the fuels, we provide parallel line plots with filled areas between the minimum and the maximum value for a feature in Figure 4. The chemical space of the pure compounds is colored in orange, the space of conventional fuels and synthetic fuels is colored in blue and green respectively. For reasons of clarity features with zero values for pure compounds and fuels are not shown.

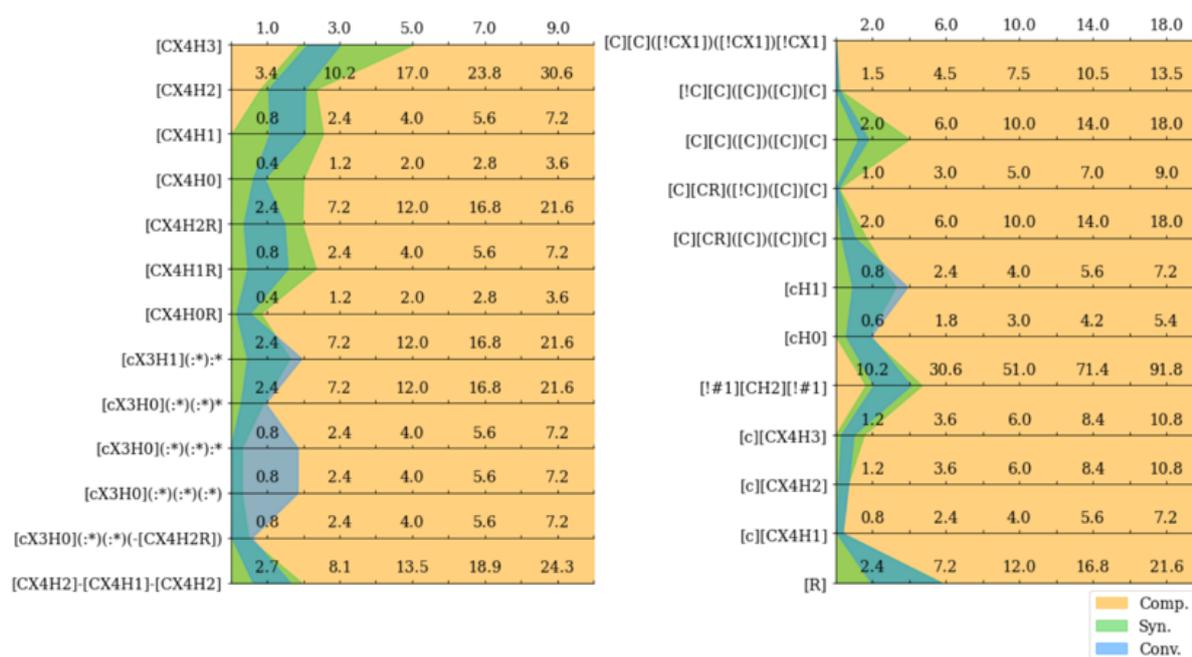


Figure 4: Parallel line plots with filled area for chemical space between minimum and maximum value in dataset

The plot illustrates that for all considered features the chemical space and thereby the application domain of the pure compounds contains the chemical spaces occupied by the pseudo-structure representation of conventional and synthetic fuels. With the exemption of the features that quantify cyclic molecule structures the chemical space of the synthetic fuels furthermore contains the space of the conventional fuels. As outlined in 2.1 Jet fuel data, this is due to the lower compositional variance of conventional fuel compositions. Figure 4 allows a visual comparison of the chemical spaces but does not allow an a priori estimation of the influence of the extend of the chemical spaces on the predictive capability of the models. We therefore investigate the influence of the training data composition on the predictive behavior for the different properties quantitatively in the following sections.

4. Probabilistic Machine Learning Model

4.1 Monte-Carlo Dropout Neural Network

For the correlation of the fuel properties with the pseudo-structure representations we utilized a Deep Neural Network algorithm with the Monte-Carlo Dropout technique (MCNN). The MCNN is a flexible Machine Learning algorithm that was developed by Gal and Ghahramani³³. It uses the popular regularization dropout technique of deactivating network neurons randomly not only during the training but also during the prediction. The prediction of a test dataset is thereby repeated multiple times, each time deactivating neurons of the network randomly, which results in varying outputs,

producing a distribution. Figure 5 illustrates the functionality of the MCNN during the prediction stage. Gal and Ghahramani describe this distribution as Bayesian approximation that captures noise of the training data and uncertainty of the prediction due to dissimilarity of training and test data. For our use case the distribution describes uncertainty due to multiple isomers for a GCxGC bin, the uncertainty associated to the approximative description of molecular structures, the measurement uncertainty and dissimilarity of training and test data. This is valuable additional information, since the predictive capability of Machine Learning models strongly depends on the data used for training and validation. The predicted distribution reflects the uncertainty of the prediction, with a narrow distribution for a certain prediction and a wider distribution for an uncertain one³⁴. The distribution itself is described by the median value y_{pred} calculated as median from predicted distribution and an upper and lower prediction interval (PI) y_{PI}^u and y_{PI}^l calculated as percentiles from the predicted distribution. The PI describe the range where a set fraction of the values are expected in^{34,35}. Figure 5 illustrates the calculation of y_{pred} as well as the PI y_{PI}^u and y_{PI}^l . We assess the validity of the distribution and the associated PI as part of the predictive capability assessment of the models. The MCNN was written using the python library pytorch³⁶.

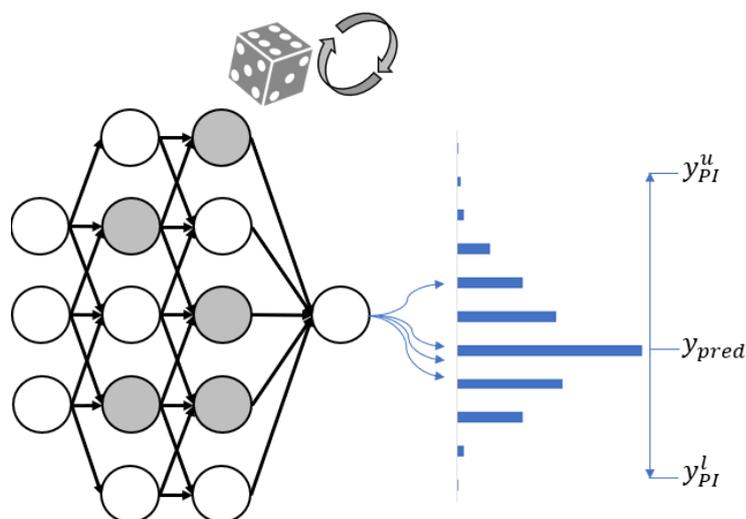


Figure 5: Schematic representation of Monte-Carlo Dropout Neural Network dropout functionality during prediction. Network neurons are deactivated randomly (grey) to generate a distribution of prediction values

4.2 Predictive capability assessment

For the assessment of the predictive capability of the trained models we utilize an assessment method, which we introduced in a recent publication⁶. The method was developed for probabilistic models and quantifies the predictive capability based on three metrics that measure the (1) accuracy and (2) precision of the prediction as well as (3) validity of the PI. For the accuracy we use the Mean Absolute Error (MAE), that calculates the mean of the absolute errors of the prediction value $y_{pred,i}$ and the actual test value $y_{test,i}$ for all measurements i , see Equation 1. The MAE was chosen to prevent

potential zero denominator calculations. To check the validity of the PI we calculate the Prediction Interval Coverage Probability (PICP), see *Equation 2*, which calculates the average probability that a measured test value lies inside between the lower PI y_{PI}^l and the upper PI y_{PI}^u of the prediction. The variable c_i therefore is a Boolean value; it is 1 if $y_{PI}^l < y_{test} < y_{PI}^u$ and 0 otherwise. If the PICP and the set confidence level of PI of the prediction are comparable, predictions do on average lie inside the PI and the PI can be considered valid. If this is true for training and testing, the PI is considered reliable. The Normalized Mean Prediction Interval Width (NMPIW) measures the precision of the prediction by calculating the mean width of the PI relative to a reference width $\Delta_{ref,i}$, see *Equation 3*.

$$MAE = \frac{1}{n_{Test}} \sum_{i=1}^{n_{Test}} y_{pred,i} - y_{test,i} \quad \text{Equation 1}$$

$$PICP = \frac{1}{n_{Test}} \sum_{i=1}^{n_{Test}} c_i * 100 \% \quad \text{Equation 2}$$

$$NMPIW = \frac{1}{n_{Test}} \sum_{i=1}^{n_{Test}} \frac{y_{PI,i}^u - y_{PI,i}^l}{\Delta_{ref,i}} * 100 \% \quad \text{Equation 3}$$

As reference width we chose the reproducibility limits of the property measurement methods. This allows a comparison of the predicted uncertainty with the uncertainty of the measurements.

For cases where the calculated PICP and the set confidence level are not comparable and therefore the PI are not valid, we introduce an additional fourth metric. The Mean Absolute Error of Outliers (MAOE), measures the mean deviation of outliers to the next upper or lower PI, see *Equation 4*. This error can be used to increase the PI in order to enclose the outliers on average and achieve a PICP of close to 100 %. For PICP values below 100 % the MAOE needs to be calculated iteratively. An illustration example will be given in the results of the density prediction in section 5.1 Density.

$$MAOE = \frac{1}{n_{test}} \sum_{i=1}^{n_{test}} c_i * \min(|y_{PI,i}^u - y_{test,i}|, |y_{test,i} - y_{PI,i}^l|) \quad \text{Equation 4}$$

All four metrics are percentage values and allow a comprehensive comparison of the predictive capability of probabilistic models.

4.3 Hyperparameter optimization

Since the training process of Machine Learning model is essentially an optimization problem ³⁴, the choice of the start parameters of this optimization, the hyperparameters of the models, are of great importance for the fit of the model. For neural networks the hyperparameters define the topology of the network e.g. the number of neurons per layer, as well as the training conditions e.g. the learning

rate. We furthermore consider the choice of an adequate scaler for input and output data a hyperparameter of the model. The complete list of hyperparameters as well as their considered options for the optimization can be found in the supplementary material S6. For the optimization of the weights and biases of neural network during the training we used the adaptive moment estimation (ADAM)³⁷ optimizer for all models.

These hyperparameters must be optimized before the actual cross-validation of the model. For this we used the principle of Bayesian optimization with a Gaussian Process Regressor with a Matern kernel function from the python library scikit-optimize³⁸, to optimize the hyperparameters in an optimization loop over 30 iterations. For each loop a 4-fold-cross-validation is computed. As data for the cross-validation optimization we utilized 800 stratified random samples with set fraction of 30% for fuels and 70% for pure compounds due to computational limitation. For the pure compounds the stratified samples were picked from groups based on the hydrocarbon family and the number of C atoms in the molecule in steps of 5, creating 35 classes, e.g. iso-alkanes 1-5 and iso-alkanes 6-10. For fuels the samples were picked from the conventional and synthetic fuel class. This stratified sampling approach guarantees the presence of both pure compounds and fuels in the data for the hyperparameter optimization. The number of samples and the range of the C atoms for the pure compounds was chosen arbitrary and proved to return sufficient results. No further investigations with other numbers of samples or the range of C atoms were undertaken.

As loss function for the hyperparameter optimization we utilized a custom function based on the Root Mean Squared Error (RMSE), see *Equation 5*. The loss function furthermore considers the two predictive capability metrics PICP and NMPIW, see *Equation 6*. For the PICP a target confidence level of 95 % was set for the confidence level of the model. The NMPIW relates the precision of the prediction to the reproducibility limits of the measurements. To constrain the maximum influence of the precision of the predictions on the hyperparameter optimization we set a maximum of 2 in *Equation 6*. Without this constraint the hyperparameter optimization could tend to dominantly focus on the precision and disregard the validity of the PI. The optimization aims at returning hyperparameters of an optimal model, that is as accurate and precise as possible with PI that comply as closely as possible to the set PICP target of 95 %. The most optimal model will probably not comply to all set thresholds of accuracy, validity and precision since the loss function describes a trade-off between the metrics. The introduction of additional weights for the different metrics in the loss function could force the compliance of one metric to a critical threshold, if required for the use case. Since the goal of this work is a comparison of the models and not their adequacy for a specific use case, the presented loss function is regarded as sufficient.

$$RMSE = \sqrt{\frac{1}{n_{Test}} \sum_{i=1}^{n_{Test}} (y_{pred,i} - y_{test,i})^2}$$

Equation 5

$$loss_{opt} = RMSE \left(1 + \frac{\max(0, 95 - PICP)}{100} + \min \left(2, \frac{\max(0, NMPIW - 100)}{100} \right) \right)$$

Equation 6

The hyperparameter configuration that yield the lowest average $loss_{opt}$ for a model in the cross-validation over all 30 iterations, is chosen as optimal configuration for the training with all data points of the three datasets.

5. Results and Discussion

To investigate the M-QSPR method we assess the predictive capability of the three models and the dependency of the method on the training data based on the prediction of 82 conventional and 50 synthetic fuels. For each property and dataset, the hyperparameters of the models were optimized using the process explained in 4.3 Hyperparameter optimization, before the training and cross-validation on the full datasets. The trained models thus differ only by the basis of the utilized training datasets.

We review the influence of the presence of fuel measurements in the training data by comparing the predictions of three models, trained on different datasets: 1) on a dataset of pure compounds (Base model), 2) on a dataset of pure compounds and conventional fuels (Fuel model 1), and 3) on a dataset consisting of data from pure compounds, conventional and synthetic fuels (Fuel model 2). We thereby investigate the ability of the M-QSPR Machine Learning model to approximate the mixing behavior for different fuel properties and fuel types with and without the presence of fuels of different fuel types in the training data. The prediction results are assessed with respect to state-of-the art mixing rules utilized for QSPR methods, that approximate jet fuels as mixtures of pure compounds, to explain possible errors of the Base models.

The assessment is presented for the investigated properties: density, freezing point, flash point, net heat of combustion, surface tension and kinematic viscosity according to the following structure. The prediction results of the models are presented as unity plots. The predicted value is thereby plotted against the corresponding measurement. The estimated prediction intervals (PI) for a 95 % confidence

level are indicated as error bars. Values of conventional fuels are displayed in blue, values of synthetic fuels are displayed in green. If prediction and measurement are in perfect accordance, the marker lies on the unity line, plotted in black. As additional frame of reference the measurement reproducibility, taken from CRC Report No. AV-23-15/17³⁹, see *Table 4*, are indicated as dashed gray lines in the figures. The predictive capability metrics are calculated separately for the conventional and the synthetic fuels and are provided in tables. For the Base model the metrics are calculated from the predictions of the fuels as hold-out data, meaning the fuel data was not used in cross-validation of the Base model. Since for the other two models (Fuel model 1 and Fuel model 2) fuels are part of the cross-validation data, the predictions are therefore taken from the testing stage of the 4-fold cross-validation of the optimized models. The validity metric PICP of the models is calculated with respect to the reproducibility limits, therefore if the measurement lies within the predicted PI +/- the reproducibility limit, the PI are considered valid. Results of the cross-validation of the pure compounds are not discussed as part of this work. The corresponding unity plots and the calculated metrics of the pure compounds can be found in the supplementary material S7. For the calculation of the PICP of the pure compounds it was assumed, that the reproducibility of the test methods of the jet fuels also apply for the measurement of the pure compounds.

Property	ASTM Method	Reproducibility
Density, kg/m ³	D4052	0.52
Flash point, °C	IP 170	3.2
Freezing point, °C	D5972	0.8
Net heat of combustion, MJ/kg	D4809	0.324
Surface Tension, mN/m	D971	0.1*X
Viscosity, mm ² /s	D445	0.019*X

Table 4: Minimum Standard Error of Prediction of selected properties as reported and computed relative error based on reference jet fuel, X variable stands for measured property value.

5.1 Density

The results of the density prediction for the three models are given in the unity plots in *Figure 6*. In all three plots the predictions of models lie close to the unity line. Therefore, all models regardless of the training data are able to compute physically correct results. This shows that all models calculate a valid mixing behavior regardless of their training data.

The predictive capability metrics in the bar plot in *Figure 7* reflect the observations from the unity plot. Accuracy, validity and precision are similar for all models. With a MAE of around 2 kg/m³ for conventional fuels and 5.5 to 6 kg/m³ for synthetic fuels, the accuracy is comparable with the results in the literature; 1.6-2.2 kg/m³^{35,40} for conventional and 5.06 kg/m³ for synthetic fuels¹⁵. However, the

direct comparison of the results of this work with the results reported in the literature is questionable, since the number and the composition of the validation datasets are not identical. The PICP of the synthetic fuels around 74-79 % compared to PICP of the conventional fuels of over 95 % is striking (black horizontal line in *Figure 7*). This means, that only the PI of the predictions of conventional fuels can be considered valid and reliable. For the synthetic fuels the PI are not valid and the prediction interval is not expressive. To make the PI valid, the calculated MAOE has to be considered. For the example of the predictions for the synthetic fuels with the Fuel model 2 a PICP of 74.45 % is calculated, therefore only 74.45 % of the measurements lie on average inside the prediction intervals. To statistically enclose all measurements the MAOE of 4.94 kg/m³ from the table in *Figure 7* has to be counted in. E.g. for a prediction with a median of 800 kg/m³, an upper PI of 805 kg/m³ and a lower PI of 795 kg/m³ the upper and the lower PI have to be increased to 809.94 kg/m³ and 790.06 kg/m³. The lower PICP is due to distinct outliers of the synthetic fuels, see *Figure 6*. Those are outliers originate from synthetic fuels with a large fraction of iso-alkanes at specific number of carbon atoms and probably consist of only a few molecules that are not further identified with the GCxGC measurement. With the exception of the distinct outliers at 726 and 790 kg/m³ the outliers are physically possible. This is asserted by calculating the possible value range for those fuels considering only molecules within lower and upper 2.5 % percentiles predicted densities. The densities for those compounds are estimated with the Base model and the range bulk density of the fuels is calculated using the recommended linear mixing rule, see *Equation 7*. These estimated lower and upper bounds of the physically possible densities are not shown in the scope of this work.

The NMPIW is comparable for all three considered models. The high values of up to over 3000 % illustrate, that the estimated uncertainty exceed the reproducibility of measurement method by a factor of 30. Since utilized composition and property measurements originate from different laboratories this high uncertainty is comprehensible, even for a relatively large dataset like this.

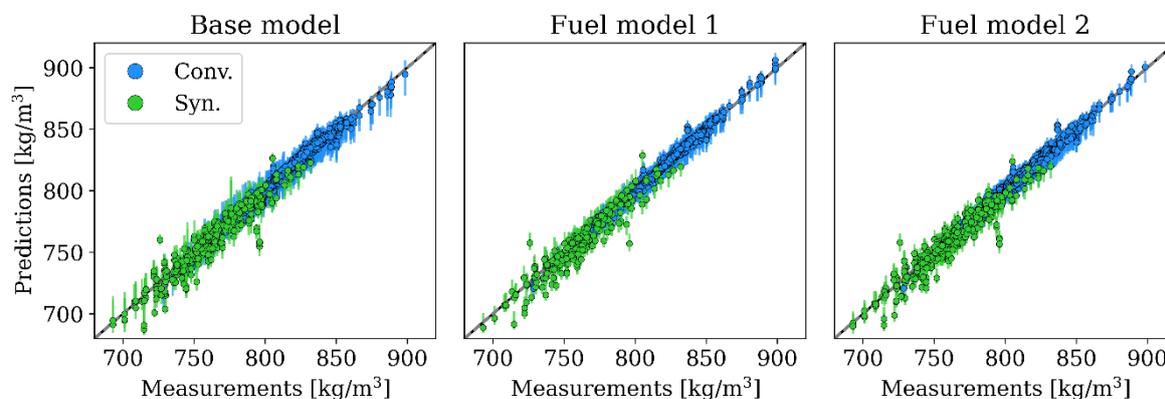


Figure 6: Unity plots of density prediction for three M-QSPR models, unity line (solid black line), reproducibility limits (black-dashed grey line)

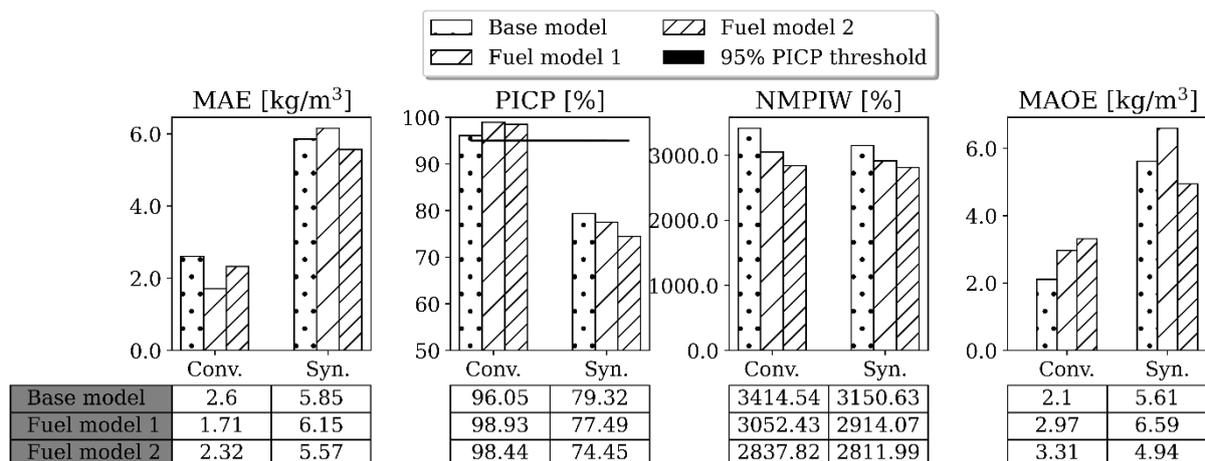


Figure 7: Bar plots and tables of predictive capability metrics of density predictions

The high accuracy of all models and therefore the good modelling of the mixing behavior, regardless of the training dataset, is comprehensible considering the utilized modelling principle of the M-QSPR method. The M-QSPR representations of the fuels are calculated by computing the linear average of the considered quantitative structures, weighted by the molar fractions of the GCxGC bins. This principle is similar to the mixing rule recommended by the literature. The literature recommends a linear mixing rule where the fuel density ρ_{mix} is computed as the sum of the density of the contained pure compound ρ_i weighted by the mass fraction w_i , see Equation 7⁴¹.

$$\rho_{mix} = \sum_i w_i * \rho_i \quad \text{Equation 7}$$

The results show, that the modelling of the density of conventional and synthetic jet fuels with the M-QSPR method is possible even without the presence of fuels in the dataset. The overall predictive capabilities are comparable. For the prediction of synthetic fuels, that are probably composed of only a few unidentified molecules, significant outliers can occur. To statistically enclose those outliers in the PI of the predictions, the MAOE has to be considered.

5.2 Freezing point

In contrast to the results of the density, the predictions of the Base model in Figure 8 show strong deviations from the unity line both for conventional and synthetic fuels. The Base model strongly underpredicts freezing point of the fuels. However, the results of the models trained both on pure components and fuels lie closer to the unity line. The presence of fuels in the training data is therefore essential for the accurate predictions of freezing point using M-QSPR models. The comparison of the MAE with results from other methods in the literature show comparable results for the conventional fuels with a MAE of 3.26 °C to 1-4 °C⁴⁰ and poorer results for the synthetic fuels with a MAE of 12.4 °C

to 5.06 °C¹⁵. Again, the exact comparability of the results of this work and the literature mentioned is questionable, since the number and the composition of the validation datasets are not identical.

The predicted PI in the unity plots are larger compared to the density prediction for all models, indicating a greater uncertainty in the property prediction. Furthermore, a strong horizontal scattering of the predictions with partly significant deviations is visible for all models, especially for the predictions of synthetic fuels. The comparison of the predictive capability metrics of the models in *Figure 9* shows a clear increase in accuracy and the validity of the PI for the Fuel model 1 and the Fuel model 2. Apart from the predictions of the conventional fuels of Fuel model 1 and Fuel model 2, the PICP of the synthetic fuels are largely inferior to 95 %, the PI are therefore not valid and reliable. Likewise to the density, the estimated MAOE has to be considered to increase the PI to enclose the outliers.

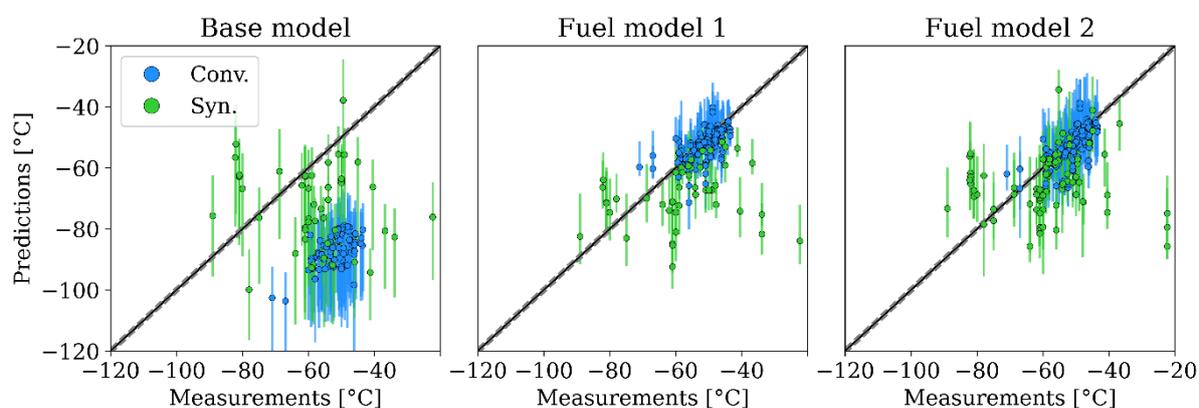


Figure 8: Unity plots of freezing point prediction for three M-QSPR models, unity line (solid black line), reproducibility limits (dashed grey line)

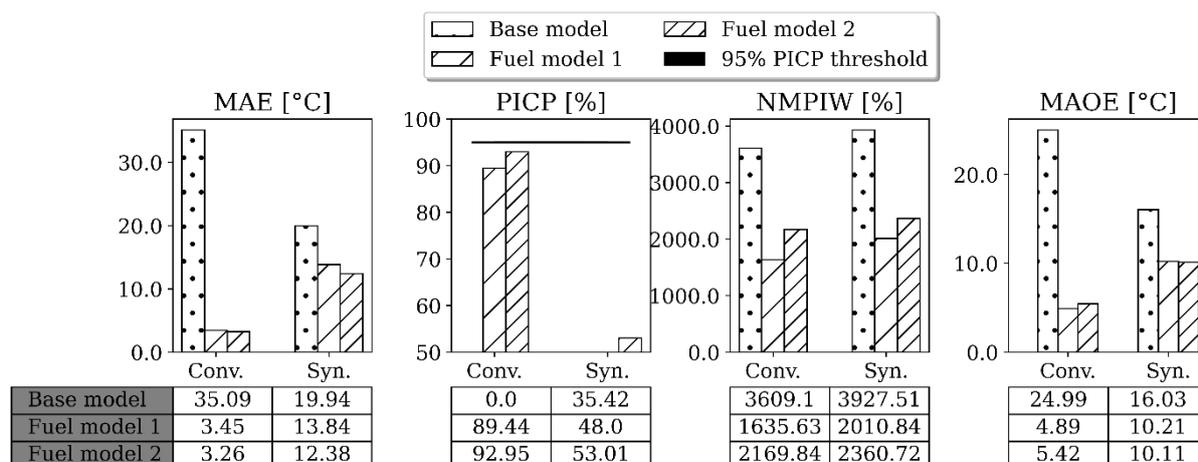


Figure 9: Bar plots and tables of predictive capability metrics of freezing point predictions

As the freezing point is a difficult fuel property to model with functional group count descriptors, the deviations of the Base model are comprehensible. The freezing of a multicomponent mixture is a

complicated process. Fuel components with higher freezing point start to crystallize and serve as crystallization seed for components with lower freezing point values⁴². As the principle of the M-QSPR method relies on a linear average of the substructures, that assumes no interactions, this could explain the strong underestimation of the Base model. For QSPR models that calculate the freezing point as mixture of pure compounds, the literature recommends a mixing rule that uses freezing indices $I_{fr,i}$ of pure compounds to calculate the freezing point of the mixture $T_{fl,mix}$ Equation 8 to Equation 10⁴². The freezing index of the mixture $I_{fr,mix}$ is thereby calculated under consideration of the volume fraction v_i .

$$I_{fr,i} = 3.23 * 10^{-6} * 1.067^{T_{fr,i}} \quad \text{Equation 8}$$

$$I_{fr,mix} = \sum_i v_i * I_{fr,i} \quad \text{Equation 9}$$

$$T_{fl,mix} = 193.798 + 15.379 * \ln(I_{fr,mix}) \quad \text{Equation 10}$$

The mixing rule illustrates the complexity of the freezing point modelling problem. The ability of the M-QSPR approach to directly incorporate fuel data in the training data resolves the need of a complicated mixing rule by learning the mixing rule implicitly. This verifiably increases the predictive capability of the models.

However, the scattering observed for the outliers of the synthetic fuels in unit plots indicate a systematic error. This scattering could be due to the smaller dataset of the freezing point with 499 unique pure compounds and fuels compared to the density 1998 and a disproportion in variance between the input and the output data. The variance in the M-QSPR representation of the fuels might not be sufficient to predict the values in the variance in the freezing points. The addition of further chemo-physical features could increase the variance and reduce the observed error. If the horizontal scattering is not reduced further by the presence of additional chemo-physical features, a more detailed representation of the fuel composition, up to the molecular level, might be necessary. Due to the small amount of data however a definite statement is difficult.

5.3 Flash point

The unity plots for the flash point predictions are displayed in *Figure 10*. Predictions of all models follow the unity line. Similar to the freezing point however, significant horizontal scattering is observed, especially for synthetic fuel predictions. The outliers correspond predominantly to fuels with a large fraction of iso-alkanes. The PI are significantly wider compared to the density prediction, indicating greater prediction uncertainty. The comparison of the predictive capability metrics in *Figure 11* shows similar results for all models with similar accuracies for all models. The calculated MAE for conventional fuels of 3.5°C is comparable to results reported in the literature of 2.5-4°C⁴⁰ and slightly poorer for the synthetic fuels, 6-7 °C to 3°C¹⁵. The different number and compositions of the validation data make a

direct comparison questionable. The PICP of the synthetic fuels for all models do not comply to the set limit of 95 %, the PI are therefore not valid and reliable and the use of the MAOE is necessary to enclose the observed outliers. Higher uncertainty as well as the scattering and deviation can probably be attributed to noise in the data due to higher reproducibility uncertainties, see *Table 4*, the smaller dataset compared to other properties and similarly to the freezing point to a disproportion in the variance of the input and the output of the models. This could, similarly to the freezing point, be the reason for the lower accuracies. Furthermore, the test methods of the flash point measurements for the pure compounds could differ from the one of the fuels. The reproducibility for the pure compounds could therefore surpass the reproducibility of the measurement of the jet fuels, inducing an even greater noise in the training data. Since information about the reproducibility of pure compound measurements are unknown, the impact cannot be estimated.

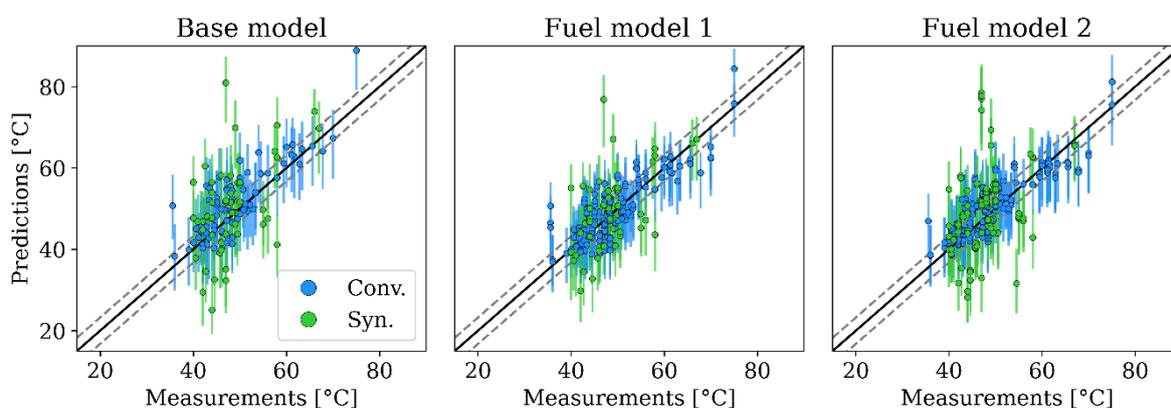


Figure 10: Unity plots of flash point prediction for three M-QSPR models, unity line (solid black line), reproducibility limits (dashed grey line)

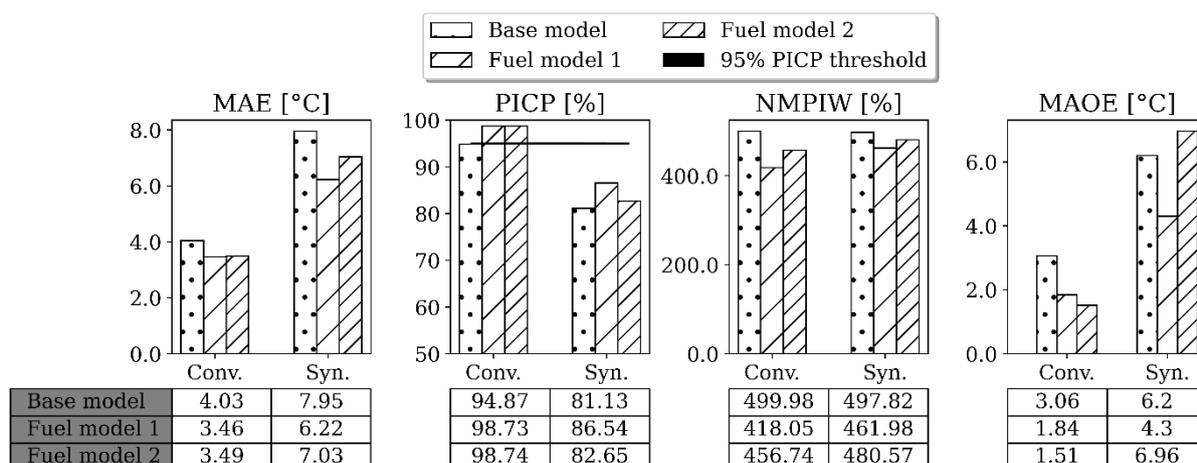


Figure 11: Bar plots and tables of predictive capability metrics of flash point predictions

For QSPR based methods, that calculate the flash point for a fuel as mixture of pure compounds, a similar method to the one of the freezing point is recommended in the literature^{43,44}. The flash point of the fuel $T_{fl,mix}$ is calculated on the basis of the ignition indices of the pure compounds $I_{fl,i}$ which are calculated from the $T_{fl,i}$ in K, see *Equation 11 to Equation 13*. The consideration of the mixing rule can therefore not explain similar prediction accuracies of the model

$$\lg(I_{fl,i}) = -6.1188 + \frac{2414}{T_{fl,i} + 503.71} \quad \text{Equation 11}$$

$$I_{fl,mix} = \sum_i v_i * I_{fl,i} \quad \text{Equation 12}$$

$$T_{fl,mix} = \frac{2414}{6.1188 + \lg(I_{fl,mix})} + 42.59 \quad \text{Equation 13}$$

Similarly to the freezing point, the observed horizontal scattering could potentially be reduced by increasing the dataset and by the addition of further chemo-physical QSPR features beyond the structural features included in this work. If the scattering does not decrease a more detailed representation of the fuel composition that differentiates between isomers might be necessary, especially for synthetic fuels with a large fraction of iso-alkanes.

5.4 Net heat of combustion

The predictions of the net heat of combustion are displayed in the unity plots in *Figure 12*. Most of the predictions for conventional and synthetic fuels lie inside the reproducibility limits of ± 0.324 MJ/kg for all three models. The M-QSPR method can therefore correctly approximate the mixing behavior of fuels even for models solely trained on pure compound data. However, the comparison of the unity plots of the models shows a systematic deviation for the conventional fuels predictions of the Base model, this could be due to a systematic difference in the measurement methods, since the measurement method of the pure compounds is often not listed in the utilized databases. For the models that were trained on fuels the systematic deviation decreases. For a group of synthetic fuels in the value range of 43.8 to 44.3 MJ/kg a systematic horizontal deviation is visible. This systematic error was also observed for fuels with similar net heat of combustion by Yang et. al¹⁵. These datapoints belong to synthetic fuels with a large fraction of iso-alkanes. Likewise to the freezing point, this systematic error might be due to insufficient variance in the M-QSPR representation or the influence of distinct isomers that are not identified by GCxGC measurement.

The comparison of the predictive capability metrics in *Figure 13* shows a clear increase of the predictive capability with the addition of fuel data to the dataset, especially with respect to the accuracy of the models for the prediction of conventional fuels. The calculated MAE of 0.09-0.12 MJ/kg is comparable to results reported in the literature 0.02-0.3 MJ/kg^{15,40}. The PICP of the predictions comply to the set

95% for all models. The precision expressed by the NMPIW is almost identical for both fuel groups and all models and similar to the reported reproducibility of the measurement method.

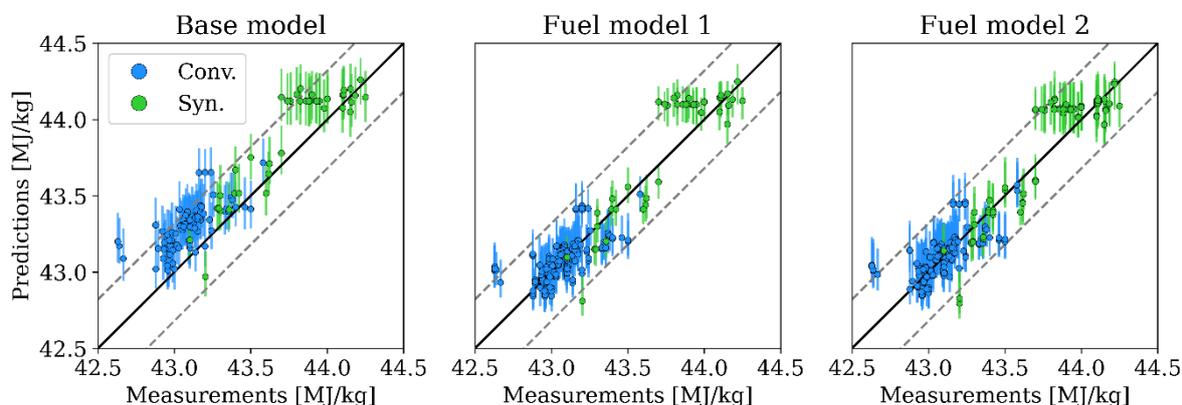


Figure 12: Unity plots of net heat of combustion prediction for three M-QSPR models, unity line (solid black line), reproducibility limits (dashed grey line)

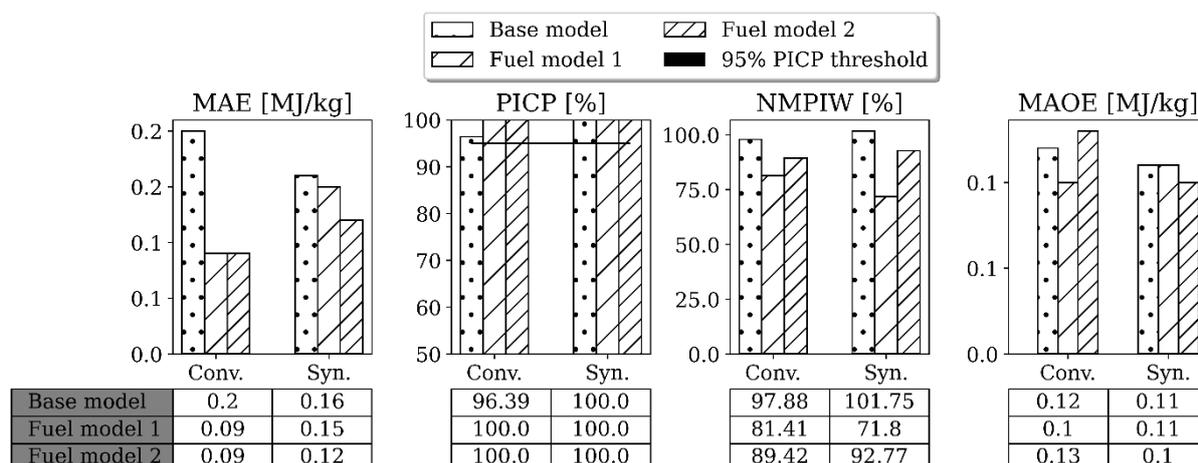


Figure 13: Bar plots and tables of predictive capability metrics of net heat of combustions predictions

The literature recommends a linear mixing rule weighted by mass fraction w_i to calculate the net heat of combustion of a fuel HOC_{mix} as mixture of pure components with individual net heat of combustion HOC_i , see Equation 14⁴³.

$$HOC_{mix} = \sum_i w_i * HOC_i \quad \text{Equation 14}$$

The increase in accuracy for the Fuel model 1 and the Fuel model 2 can probably be attributed to the correction of a systematic offset of the measurement methods of pure compound and fuels or an adjustment of the mean molecule selection to better fit the correlation property on the fuels data in the training dataset. The observed systematic error for synthetic fuels could potentially be reduced by the addition of further chemo-physical descriptors. If the error is not reducible, the GCxGC representation of the fuel composition might not provide enough variance to correctly model the net

heat of combustion for synthetic fuels. A more detailed description of the fuel up to the molecular level might be necessary, likewise to the freezing point and the flash point.

5.4 Surface tension

The results of the surface tension show a similar behavior as the ones of the density. The predictions shown in *Figure 14* lie close to the unity line and most of the predictions lie inside the uncertainty region of the reproducibility of 10 % of the measured value. All three models calculate the correct mixing behavior regardless of the training data composition. The predictive capability metrics reflect this observation. The accuracy, validity and precision are similar for all three models. The MAE of 0.36-0.62 mN/m is similar to values reported in the literature 0.49 mN/m²¹. However, it has to be noted that the direct comparison of the accuracies with results in the literature is questionable, since the number and composition of the validation data differs. The set PICP 95 % is met by every model. The NMPIW corresponds to 80-90 % of the reproducibility of the measurement method.

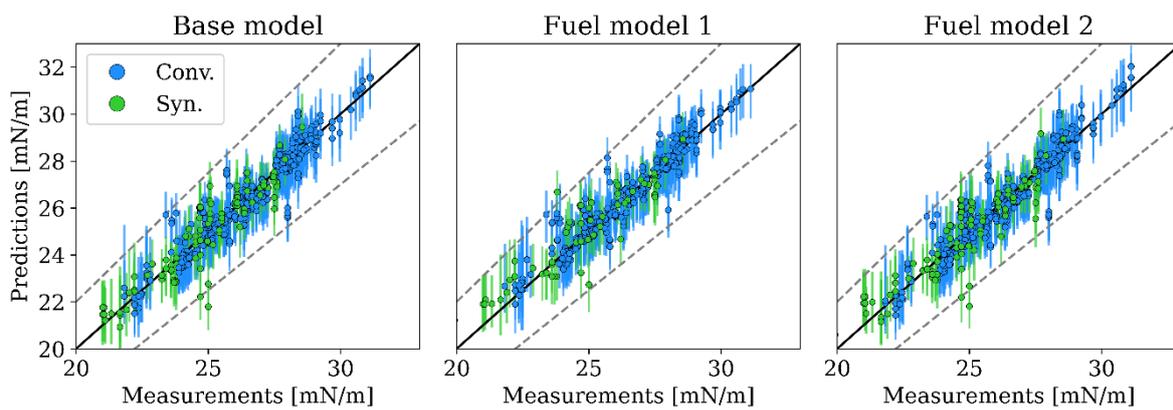


Figure 14: Prediction results of the surface tension as unity plots and comparison of predictive capability metrics as bar plots

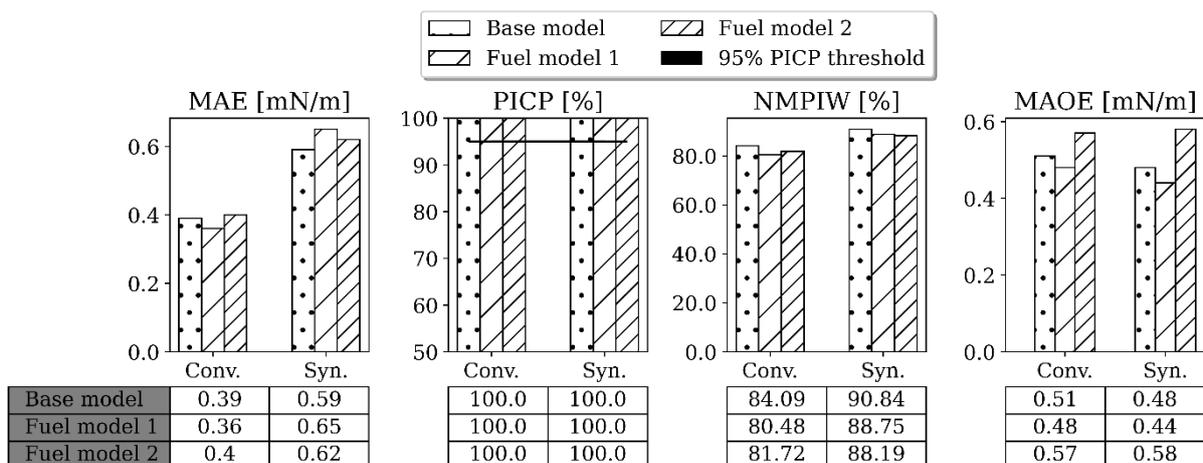


Figure 15: Bar plots and tables of predictive capability metrics of surface tension predictions

Likewise to density and net heat of combustion, the literature recommends a linear mixing rule for the calculation of the surface tension of fuel σ_{mix} as mixture of pure compounds σ_i weighted by the mass fraction w_i , see *Equation 15*⁴⁵. The results illustrate, that the modelling of the surface tension with the M-QSPR method is possible for conventional and synthetic jet fuels even without their presence in the training data.

$$\sigma_{mix} = \sum_i w_i * \sigma_i \quad \text{Equation 15}$$

5.5 Kinematic viscosity

The results for the prediction of the kinematic viscosity are displayed in *Figure 16*. For the Base model a systematic non-linear deviation is visible in the unity plot for values above 8 mm²/s, a trend that was also observed by Yang et al. for the kinematic viscosity predictions of synthetic fuels¹⁵. These values correspond to the low temperature range below -20 °C and can probably be attributed to low temperature behavior of jet fuels. Similar to the freezing point the presence of components with a higher freezing point strongly influence the viscosity due to wax formation at higher temperatures⁴². This results in higher viscosities at low temperatures, which the Base model, trained solely on pure compounds, can not describe based on the underlying linear average of the M-QSPR method. The addition of both conventional and synthetic fuels to the training data significantly improves the predictions at low temperatures. The predictive capability metrics in *Figure 17* illustrate this by a significant increase in the accuracy with a lower MAE and in the validity of the PI with a higher PICP. However, with exception of the results for synthetic fuels of the Fuel model 2 the PICP are not reliable. Therefore, the PI have to be increased with the computed MAOE, see 5.1 Density, to make the intervals valid and applicable quantitatively. The calculated MAE of 0.27-0.68 mm²/s is better compared to values reported in the literature 0.68 (-20 °C) to 2.38 (-40°C)¹⁵. However, since the number and composition of considered fuels differs from the ones reported in the literature, a direct comparison is questionable.

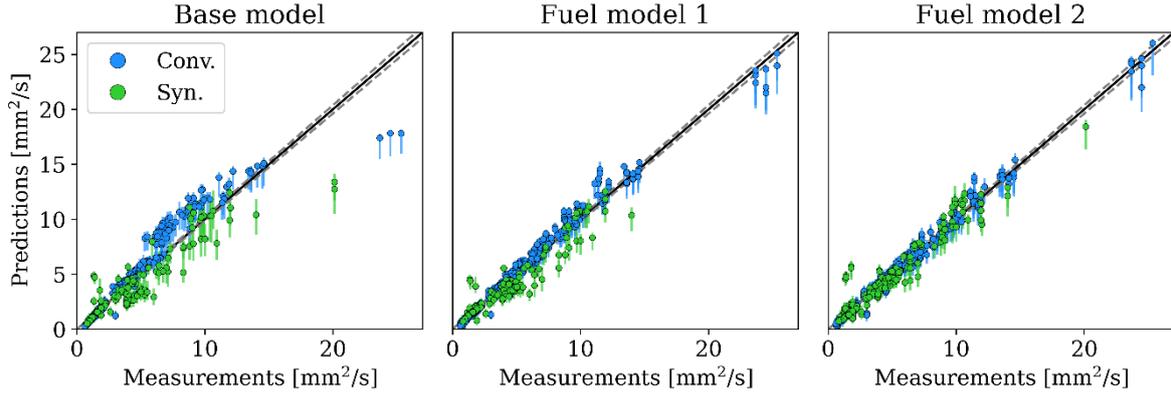


Figure 16: Prediction results of the kinematic viscosity as unity plots and comparison of predictive capability metrics as bar plots

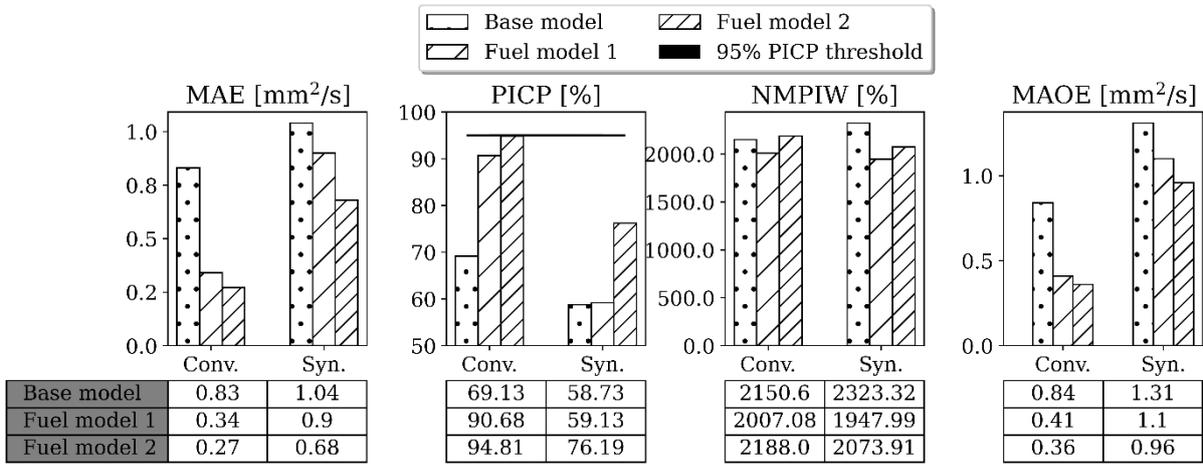


Figure 17: Bar plots and tables of predictive capability metrics of kinematic viscosity predictions

The poor results of the Base model are comprehensible due to the mentioned low temperature behavior of jet fuels. To model the kinematic viscosity with the QSPR method for fuels as mixtures of pure components the Grunberg-Nissan mixing rule is recommended in the literature, see Equation 16⁴⁶.

$$\ln(v_{mix}) = \sum_i w_i * \ln(v_i) + \frac{1}{2} \sum_{i=1}^n \sum_{j=1, j \neq i}^n w_i w_j G_{ij} \quad \text{Equation 16}$$

This is a logarithmic mixing rule that calculates the viscosity of the fuel v_{mix} as the sum of the logarithmic viscosities of the pure compounds v_i weighted by the mass fraction w_i . To consider the interactions of the fuel components the mixing rule contains binary interaction coefficients G_{ij} , to account for the mentioned mutual influence of fuel components at low temperatures. In practice, these binary interactions are often unknown and are therefore not considered in calculation of the mixing rule⁴¹. The results of the kinematic viscosity clearly show the benefit of the M-QSPR method of directly

incorporating fuel measurements in the training data. The mixing rule and interaction are then learned directly and hence, evidently improve the predictive capability of the models.

6. Summary and Outlook

Model-based tools which are able to predict critical fuel properties based on measurements from small fuel volumes are regarded a key enabler technology in the field of jet fuel screening and development. They bear the potential to significantly reduce time and cost for the approval of new synthetic jet fuel candidates. The high accuracy requirements for fuel screening require accurate predictions over the complete application domain of the models. The vast range of possible fuel compositions thereby complicates the model development. Models based on the Quantitative Structure-Property Relationship (QSPR) method have the potential to model fuels over the whole range of potential fuel compositions due to the amount of data available for pure compounds. However, QSPR models are not able to directly train on data of fuels and consequently need to rely on empirical mixing rules to calculate the bulk property of fuels. Furthermore, the established QSPR models utilize deterministic models that do not compute potential uncertainties in the predictions that stem from uncertainties in model inputs and training data like unidentified isomers in the fuel compositions or measurement uncertainties. The consideration of those uncertainties however is needed to correctly quantify the predictive capability of a model.

We presented a Mean Quantitative Structure-Property Relationship (M-QSPR) method with Monte-Carlo dropout neural network (MCNN), a probabilistic Machine Learning algorithm. The M-QSPR allows the training on both pure compounds and fuels. This drastically increases the amount of training data and makes empirical mixing rules for M-QSPR models obsolete. The probabilistic correlation models allow the prediction of uncertainties due to unidentified isomers, measurement noise and dissimilarity of training and test data without the need of time-consuming sampling of the predictions. In the context of this study, the method was applied for the modelling of jet fuels based on low volume measurements of two-dimensional gas chromatography (GCxGC). The presented M-QSPR approach represents fuels as pseudo-structures. These representations are computed from averaged quantitative structural features of potential molecules in the fuels, weighted by molar fractions from the GCxGC bins (molecular family and number of C-atoms) that the molecules are classified to. To determine the pseudo-structures, 1866 possible representative molecules from 7 hydrocarbon families are considered to compute a representation with 47 structural features. This fuel representation was correlated with the properties density, viscosity, surface tension, freezing point, flash point and net heat of combustion, which are critical properties for the screening and

development of jet fuel candidates. We investigated the ability of the M-QSPR method to predict the properties and approximate the mixing behavior both with and without the presence of fuels in the training data with three different datasets: 1) only pure component data, 2) data from pure components and conventional fuels, and 3) data from pure components, conventional fuels and synthetic fuels. All three models are tested/cross-validated on the selection of 82 conventional fuels and 50 synthetic fuels. The prediction results of the three models were compared on four metrics that quantify accuracy as well as the precision and reliability of prediction intervals of the probabilistic models.

For prediction of the density, surface tension and net heat of combustion the M-QSPR method yield highly accurate results even without the presence of fuels in the training data. This was explained by the linear mixing behavior of these properties, which corresponds to the way quantitative substructures of molecules are averaged by M-QSPR method. For the freezing point and kinematic viscosity at temperatures below -20°C , the fuel presence was found to be essential to correctly estimate the mixing behavior. The presence of fuels allows the models to directly learn low the temperature behavior of the fuel, caused by complex interactions by the fuel components. It was thereby observed that the presence of conventional fuels was often sufficient to achieve accurate results. The addition of synthetic fuels shows only small additional improvements in the predictive capability of the models. For freezing point, flash point and the net heat of combustion significant prediction outliers and systematic errors were observed, especially for synthetic fuels with a high fraction of iso-alkanes. This could be due to the smaller dataset compared to the other properties or a significant disproportion between the variance of the input and output data. The variance in the M-QSPR representation could not be sufficient to correctly estimate within the variance of the property values. This could be improved by including chemo-physical descriptors beyond the considered functional group descriptors as part of the input features. If the systematic errors are not reduced the deviations are most likely caused by the influence of isomers that are not further identified by the GCxGC measurements. A more detailed description of the fuel composition down to the molecular level might be necessary to accurately predict the listed properties for corresponding synthetic fuels.

The results show the clear benefit of the M-QSPR method to train models on data from both pure compounds and fuels to robustly and accurately predict properties over the vast range of possible jet fuel compositions. The ability of the utilized probabilistic models to predict an uncertainty distribution was found to deliver valuable additional information. Even though the predicted prediction intervals did not contain the required 95 % of the measurements for several cases, they can be used for the uncertainty quantification of the predictions under consideration of the calculated mean error of outliers. The accuracies of the property predictions are comparable with results of other methods

reported in the literature. A direct comparison of the results is however questionable, since the number and composition of the validation fuels differ. A systematic comparison of presented M-QSPR method with other methods like the QSPR method or the method of direct correlation should be the next step. Models from different methods utilizing the same GCxGC measurements (hydrocarbon-families and carbon atom range) should be validated on a set of representative conventional and synthetic reference jet fuels based on unified predictive capability metrics. This would allow an unbiased comparison of the different methods and reveal their advantages and disadvantages. Furthermore, the considered structural features should be extended by chemo physical molecular descriptors e.g. the Van-der-Waals volume to investigate potential improvements for the modelling with the M-QSPR method.

8. Supporting Information

Supporting Information	Description
S1	Comparison of modelling methods for jet fuel property prediction
S2	Considered molecular families for GCxGC measurements and classification criteria
S3	Description of utilized structural molecular features
S4, S5	Comparison of the chemical spaces and application domain of molecules generated by MOLGEN and the utilized pure compound database
S6	Tables of utilized parameters for hyperparameter optimization
S7	Testing results of cross-validation for pure compounds prediction

9. Acknowledgement

The research presented in this paper has been performed in the framework of the JETSCREEN project (JET fuel SCREENING and optimization) and has received funding from the European Union Horizon 2020 Programme under grant agreement n° 723525.

10. Abbreviations

ATJ: Alcohol to Jet; CHCJ: Catalytic Hydrothermal Conversion Jet fuel; Comp: Pure compounds; Conv.: Conventional fuels; CPK: Cycloparaffinic Kerosene; FT: Fischer-Tropsch; GCxGC: Two-dimensional gas chromatography; HEFA: Hydroprocessed esters and fatty acids; MCNN: Monte-Carlo Dropout neural network; MAE: Mean Absolute Error; MAOE: Mean Absolute Outlier Error; M-QSPR: Mean

Quantitative Structure-Property Relation; NMPIW: Normalized Mean Prediction Interval Width; PI: Prediction interval; PICP: Prediction Interval Coverage Probability; QSPR: Quantitative structure property relation; SMILES: Simplified Molecular-Input Line-Entry System; SMARTS: SMILES arbitrary target specifications; Syn.: Synthetic fuels

References

- (1) E-Well'Com. *JETSCREEN : JET Fuel SCREENING and Optimization - JETSCREEN will develop a screening and optimization platform for alternative fuels*. <https://www.jetscreen-h2020.eu/> (accessed 2021-02-24).
- (2) Colket, M.; Heyne, J.; Rumizen, M.; Gupta, M.; Edwards, T.; Roquemoire, W. M.; Andac, G.; Boehm, R.; Lovett, J.; Williams, R.; Condevaux, J.; Turner, D.; Rizk, N.; Tishkoff, J.; Li, C.; Moder, J.; Friend, D.; Sankaran, V. Overview of the National Jet Fuels Combustion Program. *AIAA Journal* **2017**, *55* (4), 1087–1104. DOI: 10.2514/1.J055361.
- (3) Heyne, J.; Rauch, B.; Le Clercq, P.; Colket, M. Sustainable aviation fuel prescreening tools and procedures. *Fuel* **2021**, *290*, 120004. DOI: 10.1016/j.fuel.2020.120004.
- (4) D02 Committee. *Practice for Evaluation of New Aviation Turbine Fuels and Fuel Additives*; ASTM International, West Conshohocken, PA.
- (5) Vozka, P.; Modereger, B. A.; Park, A. C.; Zhang, W. T. J.; Trice, R. W.; Kenttämä, H. I.; Kilaz, G. Jet fuel density via GC × GC-FID. *Fuel* **2019**, *235*, 1052–1060. DOI: 10.1016/j.fuel.2018.08.110.
- (6) Hall, C.; Rauch, B.; Bauder, U.; Le Clercq, P.; Aigner, M. Predictive Capability Assessment of Probabilistic Machine Learning Models for Density Prediction of Conventional and Synthetic Jet Fuels. *Energy Fuels* **2021**, *35* (3), 2520–2530. DOI: 10.1021/acs.energyfuels.0c03779.
- (7) Villanueva, N.; Flaconnèche, B.; Creton, B. Prediction of Alternative Gasoline Sorption in a Semicrystalline Poly(ethylene). *ACS combinatorial science* **2015**, *17* (10), 631–640. DOI: 10.1021/acscombsci.5b00094. Published Online: Sep. 18, 2015.
- (8) Saldana, D. A.; Starck, L.; Mougín, P.; Rousseau, B.; Ferrando, N.; Creton, B. Prediction of Density and Viscosity of Biofuel Compounds Using Machine Learning Methods. *Energy Fuels* **2012**, *26* (4), 2416–2426. DOI: 10.1021/ef3001339.
- (9) Saldana, D. A.; Starck, L.; Mougín, P.; Rousseau, B.; Creton, B. On the rational formulation of alternative fuels: melting point and net heat of combustion predictions for fuel compounds using machine learning methods. *SAR and QSAR in environmental research* **2013**, *24* (4), 259–277. DOI: 10.1080/1062936X.2013.766634. Published Online: Apr. 10, 2013.

- (10) Saldana, D. A.; Starck, L.; Mougin, P.; Rousseau, B.; Pidol, L.; Jeuland, N.; Creton, B. Flash Point and Cetane Number Predictions for Fuel Compounds Using Quantitative Structure Property Relationship (QSPR) Methods. *Energy Fuels* **2011**, *25* (9), 3900–3908. DOI: 10.1021/ef200795j.
- (11) Fortunato, M. Jetscreen program quantitative assessment of the jet fuel physical, chemical and thermophysical properties and development of low and high fidelity models. In *Proceedings of the 16th International Conference on Stability Handling and Use of Liquid Fuels 2019: 16th International Conference on Stability Handling and Use of Liquid Fuels; Hilton Long Beach Hotel Long Beach; United States; 8 September 2019 through 12 September 2019; Code 156514*.
- (12) NIST. *NIST Standard Reference Database 103a*. <https://www.nist.gov/mml/acmd/trc/thermodata-engine/srd-nist-tde-103a> (accessed 2021-02-26).
- (13) *DIPPR 801 Database*. <https://www.aiche.org/dippr/events-products/801-database> (accessed 2021-02-26).
- (14) Tim Nelson; Global Industry Manager - Fuels and Chemicals Phenomenex. Application of GC×GC-VUV and GC×GC-FID for the analysis of common gasoline samples, middle distillates and crude oil distillation cuts using Zebron® ZB-35HT and ZB-1PLUS GC Columns.
- (15) Yang, Z.; Kosir, S.; Stachler, R.; Shafer, L.; Anderson, C.; Heyne, J. S. A GC × GC Tier α combustor operability prescreening method for sustainable aviation fuel candidates. *Fuel* **2021**, *292*, 120345. DOI: 10.1016/j.fuel.2021.120345.
- (16) Li, R.; Herreros, J. M.; Tsolakis, A.; Yang, W. Novel Functional Group Contribution Method for Surrogate Formulation with Accurate Fuel Compositions. *Energy Fuels* **2020**, *34* (3), 2989–3012. DOI: 10.1021/acs.energyfuels.9b04270.
- (17) Li, R.; Herreros, J. M.; Tsolakis, A.; Yang, W. Machine learning regression based group contribution method for cetane and octane numbers prediction of pure fuel compounds and mixtures. *Fuel* **2020**, *280*, 118589. DOI: 10.1016/j.fuel.2020.118589.
- (18) Steinmetz, D.; Arriola González, K. R.; Lugo, R.; Verstraete, J.; Lachet, V.; Mouret, A.; Creton, B.; Nieto-Draghi, C. Experimental and Mesoscopic Modeling Study of Water/Crude Oil Interfacial Tension. *Energy Fuels* **2021**, *35* (15), 11858–11868. DOI: 10.1021/acs.energyfuels.1c00834.
- (19) Ajmani, S.; Rogers, S. C.; Barley, M. H.; Livingstone, D. J. Application of QSPR to mixtures. *Journal of chemical information and modeling* **2006**, *46* (5), 2043–2055. DOI: 10.1021/ci050559o.
- (20) Gaudin, T.; Rotureau, P.; Fayet, G. Mixture Descriptors toward the Development of Quantitative Structure–Property Relationship Models for the Flash Points of Organic Mixtures. *Ind. Eng. Chem. Res.* **2015**, *54* (25), 6596–6604. DOI: 10.1021/acs.iecr.5b01457.
- (21) Wang, Y.; Yan, F.; Jia, Q.; Wang, Q. Distributive structure-properties relationship for flash point of multiple components mixture. *Fluid Phase Equilibria* **2018**, *474*, 1–5. DOI: 10.1016/j.fluid.2018.07.005.

- (22) Coordinating Research Council. CRC Report No. 647 **2006**.
- (23) Voigt, C.; Kleine, J.; Sauer, D.; Moore, R. H.; Bräuer, T.; Le Clercq, P.; Kaufmann, S.; Scheibe, M.; Jurkat-Witschas, T.; Aigner, M.; Bauder, U.; Boose, Y.; Borrmann, S.; Crosbie, E.; Diskin, G. S.; DiGangi, J.; Hahn, V.; Heckl, C.; Huber, F.; Nowak, J. B.; Rapp, M.; Rauch, B.; Robinson, C.; Schripp, T.; Shook, M.; Winstead, E.; Ziemba, L.; Schlager, H.; Anderson, B. E. Cleaner burning aviation fuels can reduce contrail cloudiness. *Commun Earth Environ* **2021**, *2* (1). DOI: 10.1038/s43247-021-00174-y.
- (24) ChemSpider | Search and share chemistry. <https://www.chemspider.com/> (accessed 2021-02-26).
- (25) PubChem. *PubChem*. <https://pubchem.ncbi.nlm.nih.gov/> (accessed 2021-02-26).
- (26) molgen group. *MOLGEN*. <https://www.molgen.de/> (accessed 2021-02-26).
- (27) Henze, H. R.; Blair, C. M. THE NUMBER OF ISOMERIC HYDROCARBONS OF THE METHANE SERIES. *J. Am. Chem. Soc.* **1931**, *53* (8), 3077–3085. DOI: 10.1021/ja01359a034.
- (28) 1.3.5.17. *Detection of Outliers*. <https://www.itl.nist.gov/div898/handbook/eda/section3/eda35h.htm> (accessed 2021-02-26).
- (29) Nieto-Draghi, C.; Fayet, G.; Creton, B.; Rozanska, X.; Rotureau, P.; Hemptinne, J.-C. de; Ungerer, P.; Rousseau, B.; Adamo, C. A General Guidebook for the Theoretical Prediction of Physicochemical Properties of Chemicals for Regulatory Purposes. *Chemical reviews* **2015**, *115* (24), 13093–13164. DOI: 10.1021/acs.chemrev.5b00215. Published Online: Dec. 1, 2015.
- (30) Creton, B. Chemoinformatics at IFP Energies Nouvelles: Applications in the Fields of Energy, Transport, and Environment. *Molecular informatics* **2017**, *36* (10). DOI: 10.1002/minf.201700028. Published Online: Apr. 18, 2017.
- (31) *Daylight Theory: SMARTS - A Language for Describing Molecular Patterns*. <https://www.daylight.com/dayhtml/doc/theory/theory.smarts.html> (accessed 2021-09-13).
- (32) *RDKit*. <https://www.rdkit.org/> (accessed 2021-02-26).
- (33) Gal, Y.; Ghahramani, Z. *Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning*. <http://arxiv.org/pdf/1506.02142v6>.
- (34) Murphy, K. P. *Machine learning: A probabilistic perspective*; Adaptive computation and machine learning series; The MIT Press, 2012.
- (35) Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian Processes for Machine Learning*.
- (36) *PyTorch*. <https://pytorch.org/> (accessed 2021-10-26).
- (37) Kingma, D. P.; Ba, J. *Adam: A Method for Stochastic Optimization*. <http://arxiv.org/pdf/1412.6980v9>.
- (38) *scikit-optimize: sequential model-based optimization in Python — scikit-optimize 0.8.1 documentation*. <https://scikit-optimize.github.io/stable/> (accessed 2021-10-26).

(39) Thom, M. A. Review of Existing Test Methods Used for Aviation Jet Fuel and Additive Property Evaluations with Respect to Alternative Fuel Compositions.

(40) Shi, X.; Li, H.; Song, Z.; Zhang, X.; Liu, G. Quantitative composition-property relationship of aviation hydrocarbon fuel based on comprehensive two-dimensional gas chromatography with mass spectrometry and flame ionization detector. *Fuel* **2017**, *200*, 395–406. DOI: 10.1016/j.fuel.2017.03.073.

(41) Bell, D.; Heyne, J. S.; Won, S. H.; Dryer, F.; Haas, F. M.; Dooley, S. On the Development of General Surrogate Composition Calculations for Chemical and Physical Properties. In *55th AIAA Aerospace Sciences Meeting*; American Institute of Aeronautics and Astronautics: Reston, Virginia, 01092017. DOI: 10.2514/6.2017-0609.

(42) Coutinho, J. A. P. A Thermodynamic Model for Predicting Wax Formation in Jet and Diesel Fuels. *Energy Fuels* **2000**, *14* (3), 625–631. DOI: 10.1021/ef990203c.

(43) Flora, G.; Kosir, S. T.; Behnke, L.; Stachler, R. D.; Heyne, J. S.; Zabarnick, S.; Gupta, M. Properties Calculator and Optimization for Drop-in Alternative Jet Fuel Blends. In *AIAA Scitech 2019 Forum*; American Institute of Aeronautics and Astronautics: Reston, Virginia, 01072019. DOI: 10.2514/6.2019-2368.

(44) Wang, X.; Jia, T.; Pan, L.; Liu, Q.; Fang, Y.; Zou, J.-J.; Zhang, X. Review on the Relationship Between Liquid Aerospace Fuel Composition and Their Physicochemical Properties. *Trans. Tianjin Univ.* **2021**, *27* (2), 87–109. DOI: 10.1007/s12209-020-00273-5#Sec21.

(45) Wang, P.; Anderko, A.; Young, R. D. Modeling Surface Tension of Concentrated and Mixed-Solvent Electrolyte Systems. *Ind. Eng. Chem. Res.* **2011**, *50* (7), 4086–4098. DOI: 10.1021/ie101915n.

(46) Hauck, F. R.; Yang, Z.; Kosir, S. T.; Heyne, J. S.; Landera, A.; George, A. Experimental Validation of Viscosity Blending Rules and Extrapolation for Sustainable Aviation Fuel. In *AIAA Propulsion and Energy 2020 Forum*; American Institute of Aeronautics and Astronautics: Reston, Virginia, 08242020. DOI: 10.2514/6.2020-3671.

TOC Graphic

