

An Artificial Intelligence Approach to Predict Physical Properties of Liquid Hydrocarbons

Márton Virt^{1*}, Victor Zaghini Francesconi², Marius Drexler², Ulrich Arnold², Jörg Sauer², Máté Zöldy¹

¹ Department of Automotive Technologies, Faculty of Transportation Engineering and Vehicle Engineering, Budapest University of Technology and Economics, Műgyetem rkp. 3., H-1111 Budapest, Hungary

² Institute of Catalysis Research and Technology (IKFT), Karlsruhe Institute of Technology (KIT), Hermann-von-Helmholtz-Platz 1, 76344 Eggenstein-Leopoldshafen, Germany

* Corresponding author, e-mail: virt.marton@edu.bme.hu

Received: 23 August 2024, Accepted: 22 October 2024, Published online: 06 November 2024

Abstract

Accurate physical property prediction of newly developed compounds is vital across various industrial sectors, particularly for the customization of fuels and additives. Artificial intelligence (AI) has recently emerged as a best practice in numerous industrial fields because of its capacity for swift and precise calculations. While conventional methods such as group contribution models have been used to estimate physical properties from molecular structure, AI offers significant potential for improving the predictive accuracy. Thus, this work focuses on developing an AI model to predict key properties – boiling points, melting points, and flashpoints – of various hydrocarbons, to demonstrate the AI's superior predictive capabilities. A dataset consisting of 202 organic compounds was created and multilayer perceptron (MLP) neural networks were employed to estimate these properties using atomic numbers, functional groups, and molecular complexity as inputs. The model's performance was evaluated and compared against conventional group contribution methods on the same dataset. The AI model was further tested on new acetal compounds, revealing its broader applicability in both fuel and chemical sectors. Results show that the AI outperformed conventional methods, excelling in 5 out of 8 hydrocarbon types for boiling points, 7 for melting points, and all 8 for flashpoints.

Keywords

artificial intelligence, property prediction, organic compounds, e-fuels

1 Introduction

The development of new compounds plays a pivotal role in many industrial sectors. An example is the field of fuel development, which has diverse requirements for the introduction of new fuels and additives. Modern mobility focuses on carbon neutrality, thus alternative technologies and fuels have to be investigated to replace fossil fuels [1]. A promising alternative is the e-fuel technology, since it provides fuels which are compatible to the existing transport and mobility infrastructure, offering a drop-in solution [2]. E-fuels are synthetic fuels produced with renewable electricity and resources that enable carbon neutrality through the carbon cycle [3]. Since fuel properties can be engineered precisely, they can also provide clean and efficient engine operation [4].

Generally, the e-fuel development for hard-to-electrify applications, such as ship transportation, heavy-duty, and

other off-road machinery is of high interest and importance. Since these sectors are utilizing diesel engines with critical pollutant emissions, oxygenates, especially oxymethylene ethers (OME), have attracted interest as renewably synthesized diesel substitute fuels enabling low soot formation and nitrogen oxide (NO_x) emissions [5–7]. OMEs are oligomeric acetals with the structure CH₃O(CH₂O)_nCH₃. Especially the chain length of $n = 3–5$ demonstrates properties similar to conventional diesel and showcases beneficial combustion behavior, leading to soot emission reduction and resolving the soot-NO_x target conflict [8].

OMEs are oxygenates with cleaner combustion properties than fossil diesel due to the lack of C–C bonds in the molecular structure. However, the use of neat OME has its challenges. Slight engine modifications are necessary due to a lower heating value and incompatibility of the OME-fuel

to sealing materials. Also, customization of fuel properties is of importance, such as safety (flash point, FP), applicability during winter (cold filter plugging point, CFPP), and lubricity (high frequency reciprocating rig value, HFRR) [8].

Regarding OME production, these compounds are originally synthesized from methanol and its derivatives, such as formaldehyde or dimethyl ether. Methanol is a platform chemical that can be produced from renewable resources, thus considerably reducing the carbon footprint of OME in comparison to fossil diesel [9, 10]. In recent studies, modification strategies were developed to customize the fuel properties of OMEs in order to achieve better compatibility with the diesel standard and the existing infrastructure [11–15]. Here, alcohols other than methanol are used in the production process.

The design of new possible e-fuel compounds is a challenging task that requires different predictive methods to estimate physical properties. Group contribution methods are well-established methods capable of performing such estimations. Here, the molecule to be analyzed is divided into smaller groups or fragments. These groups are assigned tabulated values, so that substance properties can be estimated using calculation rules. However, the simplicity of these methods can lead to inaccuracies, especially when considering complex molecules. Hence, better tools are required. Artificial intelligence (AI) is a recently emerged method that provides excellent predictive abilities. Therefore, it would be beneficial to apply this technology to molecular property prediction problems.

Multilayer perceptron (MLP) type artificial neural networks (ANN) are commonly applied to perform similar predictive tasks. This type of ANN has a simple structure. It has an input layer with the defined input features and an output layer with neurons that calculate the final value of the output features. Between the input and output layers, there is at least one hidden layer with a certain number of neurons. Every neuron of a layer is connected to all of the neurons of the next layer. An MLP network with at least one hidden layer is a universal approximator and therefore, it can be applied to nonlinear problems as well [16].

Utilization of cognitive tools such as neural networks could have a great potential to increase mobility efficiency and sustainability [17, 18]. There are previous researches that demonstrated the possibilities of AI applications in several fields [19], such as Santak and Conduit [20], who predicted the properties exclusively of alkanes with ANNs. They concluded that the boiling point, heat capacity, vapor pressure, melting point, flash point, and viscosity

can be accurately predicted for these simple hydrocarbons. Using image processing techniques, Xu et al. [21] predicted structure-dependent properties from 3D molecular images with convolutional neural networks. The boiling and melting point could accurately be predicted, and the models could also be used to predict some molecules' critical points. Recently, Pérez-Correa et al. [22] developed MLP networks to predict eight different physical properties, including the melting and boiling points of organic compounds. The predictions were accurate, however, they concluded that the melting point prediction requires additional molecular indicators, such as the special distribution of functional groups.

In this work, a dataset with 202 organic compounds was created to train and validate MLP-type ANNs that can accurately predict fundamental and safety-relevant properties, such as boiling point, melting point, and flash point, based on the molecular structure. The most important accuracy measures are reported and compared to the accuracy of well-known group contribution methods, such as the ones from Joback and Reid [23], Constantinou and Gani [24], Pérez Ponce et al. [25] and Stefanis et al. [26]. In addition, properties of OMEs, which belong to the acetal substance class, are predicted, since this substance class is generally not well represented in such studies. Furthermore, the properties of novel OME compounds, as described in [11], are predicted to evaluate the applicability of the presented method.

2 Materials and methods

2.1 The created dataset

The primary focus of the current research is e-fuel development, thus combustible liquid compounds were selected for the dataset. The data were collected from different online available chemical databases and publications [27–31]. Table 1 presents the distribution of different substance classes in the dataset.

Table 1 Number of different classes of compounds in the dataset

Compound class	Number of compounds
<i>n</i> -Alkanes	16
Isoalkanes	22
Cycloalkanes	14
Alkenes	24
Alcohols	38
Ketones	18
Ethers	32
Esters	38

The aim of the algorithm is to predict physical properties from the molecular structure. This is described with the C, H, and O atom numbers and the number of functional groups of the molecule. These are: $-\text{CH}_3$, $-\text{CH}_2-$, $>\text{CH}-$, $>\text{C}<$, $=\text{CH}_2$, $=\text{CH}-$, $=\text{C}<$, $-\text{OH}$ (alcohol), $-\text{O}-$ (ether), $=\text{O}$, $>\text{C}=\text{O}$, $-\text{CH}=\text{O}$ (aldehyde) and $-\text{COO}-$ (ester). The geometry also influences the properties, thus the Cactvs complexity value is used to include geometrical information [32]. The molar masses, H/C, and O/C ratios were also calculated from the structures and included in the dataset. These are the possible input features of the ANNs. Melting points, boiling points and flash points were selected as target parameters, thus 3 ANN models are created to predict these output features.

2.2 Neural network creation

To predict the melting points, boiling points and flash points, 3 MISO MLP ANNs are created based on our well-established methods described in [16]. The ANN models are created in a Python environment using the Keras and TensorFlow libraries. The procedure starts with the dataset processing. From the 202 sample compounds of the entire dataset, 151 were used to train, 30 to validate the models during hyperparameter optimization, and 21 to test the performance of the final networks. In order to reduce the number of possible input features, we do not distinguish between ring and non-ring functional groups; only the overall number of them is used. Thereby, the possible number of input features is reduced to 20 and their quality also improves as the number of nonzero data is reduced. Before starting the training process, both the input and output features are scaled into a range of 0 to 1.

First, the relevant input features have to be selected. This starts by removing redundant input features. Then, a recursive feature elimination (RFE) process is done to select the main parameters of interest. Adding or removing features from the automatically selected list can also be beneficial.

The Adam algorithm is employed on the *training dataset* to train the MLP models. This adaptive training method optimizes the network's weights and biases quickly to minimize the loss function, which was selected to be the mean squared error (MSE). Before each training iteration, the order of the training data is randomized to mitigate the impact of local minima. The randomness of the training is treated with repeated training and evaluation, thus the hyperparameter optimization is based on the average accuracies of 8 repetitions. To avoid the vanishing gradient problem, the Rectified Linear Unit (ReLU) activation function is used at the hidden

layers with He initialization. Since a regression problem is solved, linear activation functions are used in the output layers. An early stopping method is applied to avoid overfitting. The maximum number of epochs was chosen to be 500, with a patience of 50 epochs.

Architecture optimization is conducted using a constructive architecture selection method. During optimization, the network performances are evaluated with the *validation dataset*. A topology is accepted if the coefficient of determination (R^2) is at least 0.98 or the topological boundaries are reached. After identifying a good architecture, the final MLP model is created. The previously unseen *test dataset* is used to evaluate its performance; thus, direct and indirect data leakage is avoided. The model performances are evaluated based on three common measures applied in the field of AI. These are the root mean squared error (RMSE), the mean average error (MAE), and the coefficient of determination (R^2). Prediction error plots and learning curves are also applied to evaluate network accuracies. The described model creation process is summarized on Fig. 1.

2.3 The reference predictive methods

To evaluate the performance of the AI method, it is compared to conventional property estimation methods, such as group contribution methods. Here, a comparison with the methods of Joback and Reid [23], Constantinou and Gani [24], Pérez Ponce et al. [25] and Stefanis et al. [26] is made. The molecule to be analyzed is divided into smaller groups or fragments. These groups are assigned tabulated values so that the substance property can be estimated using a calculation rule. Some group contribution methods, such as the one from Constantinou and Gani [24], make a distinction between first-order and second-order estimates. In the second-order estimation, the groups into which the molecule is divided are larger than in the first-order estimation. Here, the interaction of different atoms

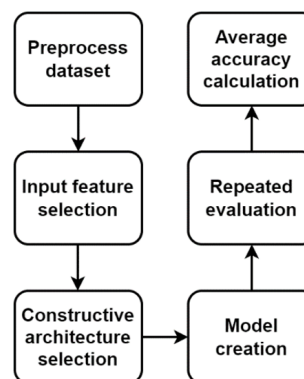


Fig. 1 The summary of the AI model creation process

or groups of atoms in the molecule is taken into account so that the estimation of the properties is more precise. Therefore, an adapted group contribution is used for the calculation rule if, for example, a CH_2 group is bound to an ether group, ester group, ketone group, etc. Regardless of whether a molecular fragment can be assigned to a second-order group, it is assigned to a first-order group. The tabulated values assigned to each of these groups are inserted into the corresponding calculation rule so that the property can be estimated.

3 Results and discussion

3.1 Performance evaluation of the created AI models

Section 3.1 demonstrates the performance of the created MISO MLP networks. First, the boiling point model is discussed. Through the manually assisted RFE method, the following features were selected as the model's input: C, O, $-\text{CH}_3$, $-\text{CH}_2-$, $>\text{CH}-$, $-\text{OH}$, $=\text{O}$, and the Cactvs complexity. The best network structure identified by the architecture selection method on the investigated range was 8–55–55–65–1. The three hidden layers with these relatively high numbers of neurons mean the relation between the input and output parameters is complex. Thus, high network capacity was required for accurate predictions.

Table 2 presents the investigated accuracy measures of the boiling point model. As expected, the model performs better on the training dataset than on the validation and test sets. However, the performance is similarly good for all sub-datasets, which suggests that the network achieved a good fit. Considering the large temperature range of 400 °C, the MAE of 13.34 °C and the test R^2 of 0.9528, good predictive abilities are achieved in the case of boiling points. The prediction error of the three datasets is presented in Fig. 2. The number of outliers is small; the predictions mostly remain close to the unity line. Overall, the boiling point model can be considered reliable.

In the case of the melting point model, the following input features were selected: C, O, $-\text{CH}_3$, $-\text{CH}_2-$, $-\text{OH}$, $=\text{O}$, and the Cactvs complexity. The network architecture is 7–20–20–1. Table 3 contains the accuracy measures of the model. All datasets have similar accuracies and the model has a lower accuracy when compared to the prediction of the boiling point. Fig. 3 presents the learning curves of the

Table 2 Accuracy measures of the boiling point model

Accuracy measure	Training	Validation	Test
RMSE (°C)	14.13	22.59	19.63
MAE (°C)	7.18	14.36	13.34
R^2 (-)	0.9711	0.9474	0.9528

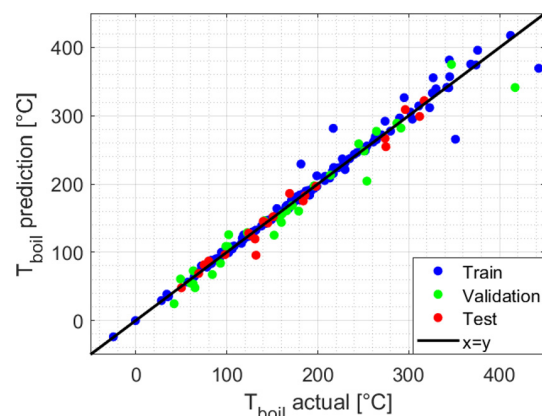


Fig. 2 Prediction error plot of the boiling point model for the train, validation and test datasets

Table 3 Accuracy measures of the melting point model

Accuracy measure	Training	Validation	Test
RMSE (°C)	24.51	23.96	24.77
MAE (°C)	18.27	16.68	19.58
R^2 (-)	0.7622	0.7755	0.8045

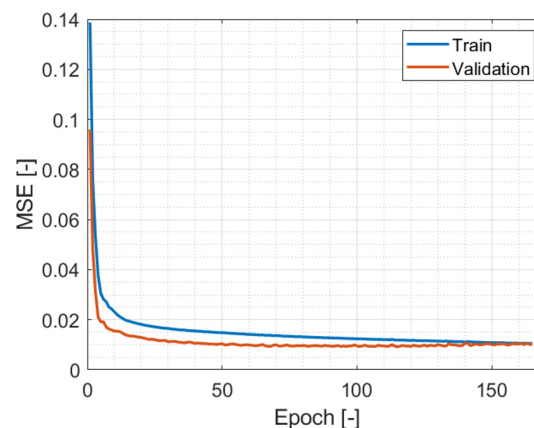


Fig. 3 Learning curves of the melting point model

model. The early stopping method terminated the training after 165 epochs. At the end of the training, the train and validation curves reached a similar level; thus, a good fit was achieved despite the moderate accuracies. The MAE for the test dataset is 19.58 °C, which is higher than that for the boiling point model, especially if the smaller temperature range of the melting points is considered. Note that this error is still lower than the error in the similar investigation of Pérez-Correa et al. [22], who achieved 26.23 °C MAE. The test R^2 is 0.8045, which is also low, especially when compared to the boiling point model.

Fig. 4 demonstrates the prediction accuracy of the network. It is discernable that the predictions are farther away from the unity line; however, the points mostly remain inside a ± 50 °C range around that line. This means that the main characteristics of the input-output relation were

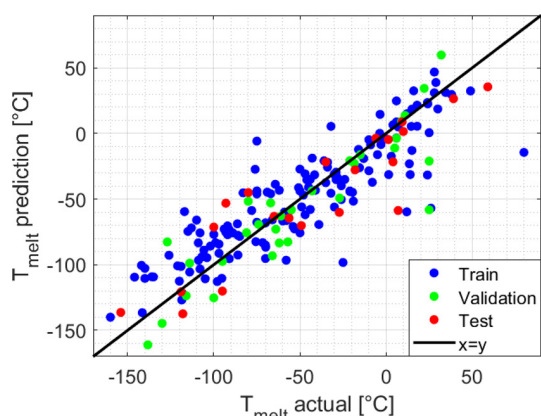


Fig. 4 Prediction error plot of the melting point model for the train, validation and test datasets

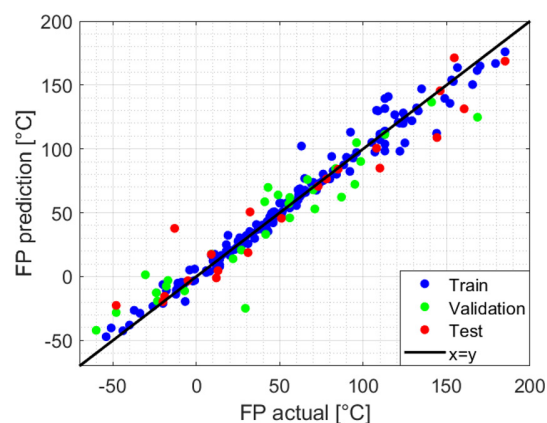


Fig. 5 Prediction error plot of the flash point model for the train, validation and test datasets

mapped by the network, but additional behaviors should also be mapped to improve accuracy. Since the model has only two hidden layers with a smaller number of neurons, one could assume that the accuracy could be improved by increasing the model capacity. However, we experienced overfitting in the case of higher-capacity networks: more neurons and more inputs led to better training accuracy, however the validation and test accuracies start to drop rapidly, which means that the generalization is no longer appropriate. Joback and Reid [23] concluded that the prediction of the melting point is problematic because it is highly dependent on the molecule geometry. Our dataset only consists of indirect information on geometry, such as the complexity value *Cactvs* and the number of functional groups, but exact connections and geometrical data are not present. For further improvements, more data is needed on the molecule geometry. Overall, the melting point model can still be considered applicable since a good fit was achieved, and the accuracy measures are still in an acceptable range.

The final predicted variable was the flash point. The model accuracies presented in Table 4 are generally good, however the difference between the training and validation accuracies suggests that the generalization is not optimal. The network architecture was 10–60–70–65–1, while the selected inputs were C, O, $-\text{CH}_2$, $=\text{CH}$, O, OH, $=\text{O}$, *Cactvs* complexity, H/C ratio, and O/C ratio. Fig. 5 demonstrates the prediction error of the flash point model. The predictions are relatively close to the unity line, and the number of outliers is small. Overall, the model has

Table 4 Accuracy measures of the flash point model

Accuracy measure	Training	Validation	Test
RMSE (°C)	8.79	18.45	18.21
MAE (°C)	5.80	14.06	13.51
R^2 (-)	0.9738	0.8864	0.9281

good accuracy. However, the system behavior could not be mapped as well as in the case of the boiling point model.

3.2 AI method's comparison to conventional methods

Section 3.1 demonstrates that the created MLP networks are accurate. To evaluate model performances, Section 3.2 compares the AI models with conventional group contribution methods.

Table 5 compares the prediction accuracies of the AI model and the methods from Joback and Reid [23] and Constantinou and Gani [24] for boiling point predictions. From the two reference methods, the method from Constantinou and Gani [24] provides more accurate results; however, the AI model outperforms both conventional methods. The MAE is only 8.89 °C in the case of the AI method, while for the method of Constantinou and Gani [24] a MAE of 19.82 °C is achieved.

The model's performance regarding the prediction of the melting point can be seen in Table 6. In addition to the previous methods, the model from Pérez Ponce et al. [25] is also

Table 5 Comparison of boiling point models' performances on the full dataset

Accuracy measure	AI	Joback and Reid [23]	Constantinou and Gani [24]
RMSE (°C)	16.30	42.69	35.14
MAE (°C)	8.89	21.98	19.82
R^2 (-)	0.9647	0.7583	0.8362

Table 6 Comparison of melting point models' performances on the full dataset

Accuracy measure	AI	Joback and Reid [23]	Constantinou and Gani [24]	Pérez Ponce et al. [25]
RMSE (°C)	24.46	40.43	36.96	44.4173
MAE (°C)	18.17	29.19	25.54	32.6455
R^2 (-)	0.7700	0.3711	0.4745	0.2412

taken as a reference since the prediction of the melting point is associated with greater difficulty. It is discernible that the model from Constantinou and Gani [24] is again the most accurate among the reference methods, while the model from Pérez Ponce et al. [25] has the worst performance. The AI outperforms the reference methods here as well. The MAE is 18.17 °C for the AI method, while the Constantinou and Gani [24] method provided a MAE of 25.54 °C.

Table 7 contains the accuracies of the flash point models. As a reference method, Stefanis et al. [26] provide a formula to calculate the flashpoints. In comparison with the AI method, the difference is high. The AI model has a 7.83 °C MAE, while the reference method has 14.40 °C as MAE. Thus, the AI outperforms the conventional method.

A comparison of the various methods' performance for the different classes of compounds (cf. Table 1) is shown in Fig. 6. Evaluation of the predicting ability of the models by compound class illustrates again that the AI model generally outperforms the conventional methods with little exception. It becomes evident that values for oxygenates, in particular, can be predicted with higher accuracy when compared to the conventional methods. This aspect seems especially promising in the context of the development of synthetic fuels to substitute fossil fuels, as generally novel compounds like e-fuels are considered and the use of oxygenates (e.g., OME) to reduce the amount of emissions produced during the combustion process is evaluated.

3.3 Predicting the properties of OMEs

The performance of the AI's predictive methods was tested in a real research application by forecasting the properties of acetals as a new compound class not included in the data set so far. The training dataset of the AI contained many types of different organic compounds in order to create a general mapping between their molecular features and their respective physical properties. Thus, other types of compounds, such as the currently investigated acetals, may be predicted sufficiently accurate as well.

The investigated methods were applied to a set of 18 acetals, including the aforementioned modified OME. Novel acetals [11] were also included in the comparison.

Table 7 Comparison of flash point models' performances on the full dataset

Accuracy measure	AI	Stefanis et al. [26]
RMSE (°C)	11.97	20.82
MAE (°C)	7.83	14.40
R^2 (-)	0.9548	0.8630

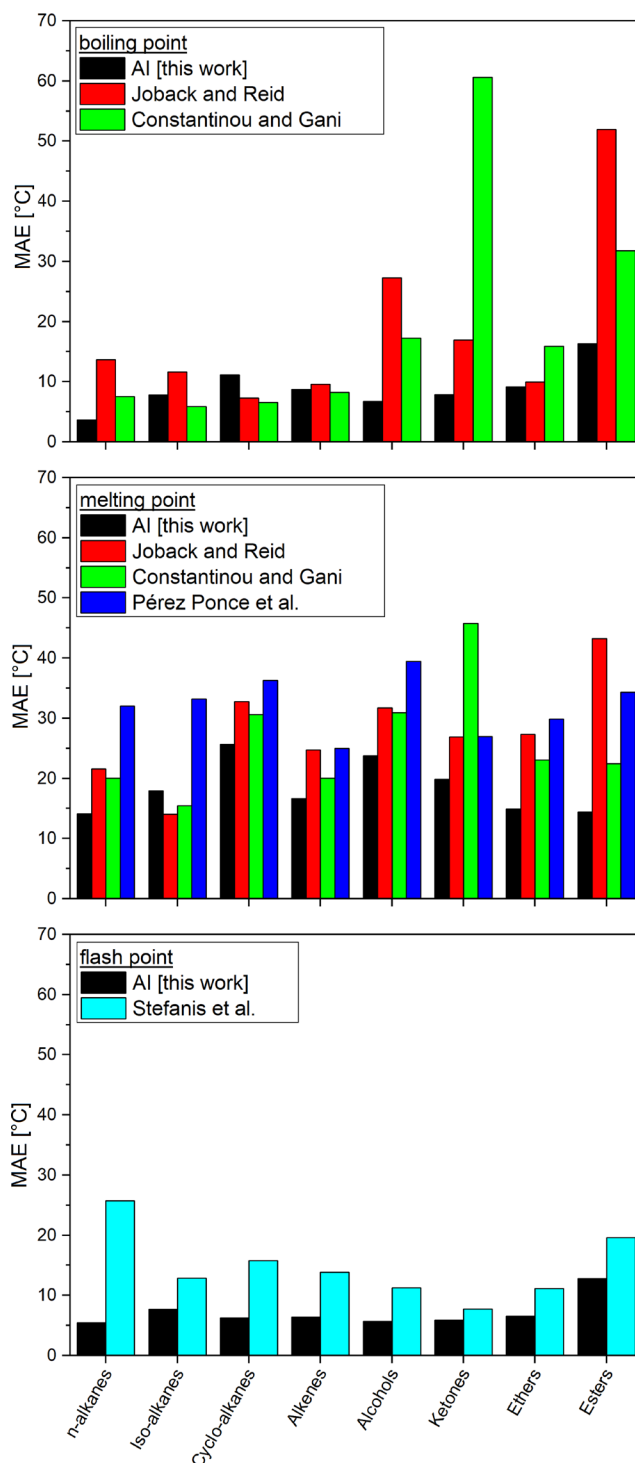


Fig. 6 MAEs of the AI model as well as group contribution methods by compound class for all predicted properties

These novel substances are acetals with branches in their molecule, a structure that has not been extensively studied so far. The literature values for the boiling, melting and flash points can be found in [8, 11, 33–39]. Fig. 7 compares the prediction accuracy of the methods for all properties, while the calculated accuracy measures can be found in Table 8.

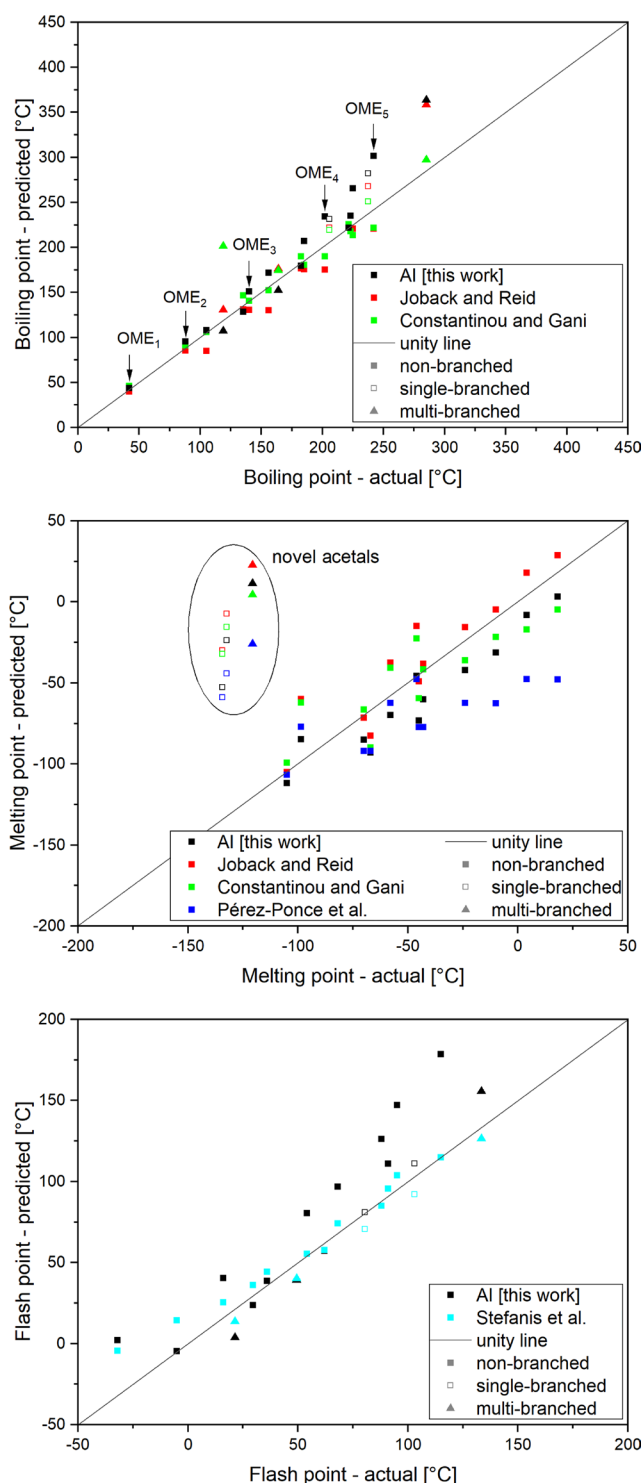


Fig. 7 Prediction error plot for boiling point, melting point and flash point for the acetal set for all applied methods

It is discernable that the AI's good performance decreased with this unseen type of substance class. In the case of boiling point, the group contribution methods outperformed the AI, although the highest R^2 value could be achieved. Detailed evaluation revealed an increase in prediction error within one group of acetals with increasing chain length,

Table 8 Comparison of the prediction methods performances on the acetal dataset

Parameter		Number of compounds	RMSE (°C)	MAE (°C)	R^2 (-)
Boiling point	AI	18	30.23	21.55	0.9546
	Joback and Reid [23]		22.71	15.32	0.9187
	Constantinou and Gani [24]		21.53	12.26	0.8839
Melting point ¹	AI	12	17.04	15.38	0.9087
	Joback and Reid [22]		17.26	12.92	0.8667
	Constantinou and Gani [24]		18.68	16.04	0.7683
	Pérez Ponce et al. [25]		35.39	29.24	0.6727
Flash point	AI	17	27.58	21.13	0.8624
	Stefanis et al. [26]		10.64	8.51	0.9554

¹ Novel, branched acetals from [8] have been excluded from the data set for the melting point prediction for all methods as they were deemed unsuitable for taking the molecular structure into account properly due to insufficient input features.

as illustrated in Fig. 7 for the group of OME₁₋₅. This pattern could indicate that the -O- group is insufficient to represent higher acetals with multiple connected functional groups of this type properly, therefore increasing the prediction error with each added chain segment. Regarding the melting points, the AI, as well as the method of Joback and Reid [23], perform similarly well, although it has to be noted that none of the applied methods seems suitable for the prediction of novel, branched acetals. Due to the bad prediction for all models and the small size of the data set for acetals as a new substance class, these data points have been excluded from the calculation of the accuracy measures in Table 8. As the AI method is not able to predict the melting points of these compounds, the implemented Cactvs complexity value does not seem to contain sufficient information about the geometry, thus other alternatives should be investigated. In the case of the flash points, the method of Stefanis et al. [26] is more accurate than the AI, which might change if acetals are added to the training data set. From the results of Table 8, it can be concluded that the AI models could achieve a generalization and the predictions remained competitive, even for an unseen type of molecule to some extent. However, modifications in the algorithm could help to improve this in the future. Further steps should include expanding the database to increase the predicting capabilities of the method as well as employing additional geometry features to reflect the influence of

different functional groups and molecule structures for the prediction of different substance properties.

4 Conclusion

This paper presents an AI approach to predict the physical properties of different compounds based on their molecular structures. The method of AI creation, as well as the performance analysis of the trained neural networks, were considered. The performance of the AI models was compared to conventional group contribution methods, and the model was applied to acetals as a new compound class to test a possible application scenario.

It can be concluded that AI can be effectively used to predict physical properties based on the molecular structure of substances. The created models were accurate, and the group contribution methods were outperformed regarding the range of investigated organic compounds used as a database. The presented AI method had the best performance in case of 5 out of 8 hydrocarbon types for boiling points, 7 for melting points, and all 8 for flash-points. The good performance of the melting point prediction is one of the most important findings, since conventional methods can hardly predict it accurately.

References

- [1] Singh, G., Esmailpour, M., Ratner, A. "Effect of carbon-based nanoparticles on the ignition, combustion and flame characteristics of crude oil droplets", *Energy*, 197, 117227, 2020.
<https://doi.org/10.1016/j.energy.2020.117227>
- [2] Cipriano, E., da Silva Major, T. C. F., Pessela, B., Chivanga Barros, A. A. "Production of Anhydrous Ethyl Alcohol from the Hydrolysis and Alcoholic Fermentation of Corn Starch", *Cognitive Sustainability*, 1(4), 2022.
<https://doi.org/10.55343/cogsust.36>
- [3] Prentice, I. C., Farquhar, G. D., Fasham, M. J. R., Goulden, M. L., Heimann, M., Jaramillo, V. J., Kheshgi, H. S., ... Yool, A. "The Carbon Cycle and Atmospheric Carbon Dioxide", In: Houghton, J. T., Ding, Y., Griggs, D. J., Noguer, M., van der Linden, P. J., Dai, X., Maskell, K., Johnson, C. A. (eds.) *Climate Change 2001: The Scientific Basis*, Cambridge University Press, 2001, pp. 183–238. ISBN 0521 80767 0 [online] Available online: <https://www.ipcc.ch/site/assets/uploads/2018/02/TAR-03.pdf> [Accessed: 12 December 2023]
- [4] Lindstad, E., Lagemann, B., Rialland, A., Gamlem, G. M., Valland, A. "Reduction of maritime GHG emissions and the potential role of E-fuels", *Transportation Research Part D: Transport and Environment*, 101, 103075, 2021.
<https://doi.org/10.1016/j.trd.2021.103075>
- [5] Omari, A., Heuser, B., Pischinger, S., Rüdinger, C. "Potential of long-chain oxymethylene ether and oxymethylene ether-diesel blends for ultra-low emission engines", *Applied Energy*, 239, pp. 1242–1249, 2019.
<https://doi.org/10.1016/j.apenergy.2019.02.035>
- [6] Park, W., Park, S., Reitz, R. D., Kurtz, E. "The effect of oxygenated fuel properties on diesel spray combustion and soot formation", *Combustion and Flame*, 180, pp. 276–283, 2017.
<https://doi.org/10.1016/j.combustflame.2016.02.026>
- [7] Tan, Y. R., Botero, M. L., Sheng, Y., Dreyer, J. A. H., Xu, R., Yang, W., Kraft, M. "Sooting characteristics of polyoxymethylene dimethyl ether blends with diesel in a diffusion flame", *Fuel*, 224, pp. 499–506, 2018.
<https://doi.org/10.1016/j.fuel.2018.03.051>
- [8] Lautenschütz, L., Oestreich, D., Seidenspinner, P., Arnold, U., Dinjus, E., Sauer, J. "Physico-chemical properties and fuel characteristics of oxymethylene dialkyl ethers", *Fuel*, 173, pp. 129–137, 2016.
<https://doi.org/10.1016/j.fuel.2016.01.060>
- [9] Mantei, F., Kopp, S., Holfelder, A., Flad, E., Kloeters, D., Kraume, M., Salem, O. "Suitable commercial catalysts for the synthesis of oxymethylene dimethyl ethers", *Reaction Chemistry & Engineering*, 8(4), pp. 917–932, 2023.
<https://doi.org/10.1039/D2RE00508E>
- [10] Voelker, S., Deutz, S., Burre, J., Bongartz, D., Omari, A., Lehrheuer, B., Mitsos, A., Pischinger, S., Bardow, A., von der Assen, N. "Blend for all or pure for few? Well-to-wheel life cycle assessment of blending electricity-based OME₃₋₅ with fossil diesel", *Sustain Energy & Fuels*, 6(8), pp. 1959–1973, 2022.
<https://doi.org/10.1039/D1SE01758F>

The method was also applied to a real research scenario, where the properties of compounds of a new class were predicted. The new substance class of acetals was not previously part of the training dataset. However, the AI could still make accurate predictions when compared to group contribution methods. Therefore, a good generalization was achieved. Regarding future investigations, further improvement of the method is possible. The main focus should be on finding better features that contain more comprehensive information on the molecular geometry of the substances. Then, the improved models could be applied to enhance the prediction of existing properties and predict more complex properties, such as cetane numbers or cold-properties.

Acknowledgement

The research leading to this result was funded by the KTI_KVIG_8-1_2021.

This work was supported by AVL Hungary Kft.

Financial support from the Helmholtz Association within the research program "Materials and Technologies for the Energy Transition" (MTET) is gratefully acknowledged.

- [11] Drexler, M., Haltenort, P., Zevaco, T. A., Arnold, U., Sauer, J. "Synthesis of tailored oxymethylene ether (OME) fuels *via* transacetalization reactions", *Sustainable Energy & Fuels*, 5(17), pp. 4311–4326, 2021.
<https://doi.org/10.1039/D1SE00631B>
- [12] Jacobs, S., Döntgen, M., Alqaity, A. B. S., Hesse, R., Kruse, S., Beeckmann, J., Kröger, L. C., Morsch, P., Leonhard, K., Pitsch, H., Heufer, K. A. "A Comprehensive Experimental and Kinetic Modeling Study of the Combustion Chemistry of Diethoxymethane", *Energy & Fuels*, 35(19), pp. 16086–16100, 2021.
<https://doi.org/10.1021/acs.energyfuels.1c01988>
- [13] Bartholet, D. L., ArellanoTreviño, M. A., Chan, F. L., Lucas, S., Zhu, J., St. John, P. C., Alleman, T. L., McEnally, C. S., Pfefferle, L. D., Ruddy, D. A., Windom, B., Foust, T. D., Reardon, K. F. "Property predictions demonstrate that structural diversity can improve the performance of polyoxymethylene ethers as potential bio-based diesel fuels", *Fuel*, 295, 120509, 2021.
<https://doi.org/10.1016/j.fuel.2021.120509>
- [14] Arellano-Treviño, M. A., Baddour, F. G., To, A. T., Alleman, T. L., Hays, C., Luecke, J., Zhu, J., McEnally, C. S., Pfefferle, L. D., Foust, T. D., Ruddy, D. A. "Diesel fuel properties of renewable polyoxymethylene ethers with structural diversity", *Fuel*, 358, 130353, 2024.
<https://doi.org/10.1016/j.fuel.2023.130353>
- [15] Arellano-Treviño, M. A., Alleman, T. L., Brim, R., To, A. T., Zhu, J., McEnally, C. S., Hays, C., Luecke, J., Pfefferle, L. D., Foust, T. D., Ruddy, D. A. "Blended fuel property analysis of butyl-exchanged polyoxymethylene ethers as renewable diesel blendstocks", *Fuel*, 322, 124220, 2022.
<https://doi.org/10.1016/j.fuel.2022.124220>
- [16] Virt, M., Zöldy, M. "Cost Efficient Training Method for Artificial Neural Networks based on Engine Measurements", *Acta Polytechnica Hungarica*, 21(7), pp 123–145, 2024.
<https://doi.org/10.12700/APH.21.7.2024.7.8>
- [17] Zöldy, M., Baranyi, P., Török, Á. "Trends in Cognitive Mobility in 2022", *Acta Polytechnica Hungarica*, 21(7), pp. 189–202, 2024.
<https://doi.org/10.12700/APH.21.7.2024.7.11>
- [18] Filina-Dawidowicz, L., Stankiewicz, S., Čižiūnienė, K., Matijošius, J. "Factors influencing intermodal transport efficiency and sustainability", *Cognitive Sustainability*, 1(1), 2022.
<https://doi.org/10.55343/cogsust.9>
- [19] Barrachina, D. G.-L., Boldizar, A., Zoldy, M., Torok, A. "Can Neural Network Solve Everything? Case Study of Contradiction In Logistic Processes With Neural Network Optimisation", In: 2019 Modern Safety Technologies in Transportation (MOSATT), Kosice, Slovakia, 2019, pp. 21–24. ISBN 978-1-7281-5084-0
<https://doi.org/10.1109/MOSATT48908.2019.8944120>
- [20] Santak, P., Conduit, G. "Predicting physical properties of alkanes with neural networks", *Fluid Phase Equilibria*, 501, 112259, 2019.
<https://doi.org/10.1016/j.fluid.2019.112259>
- [21] Xu, Y., Huang, X., Li, C., Wei, Z., Wang, M. "Predicting structure-dependent properties directly from the three dimensional molecular images via convolutional neural networks", *AIChE Journal*, 68(8), e17721, 2022.
<https://doi.org/10.1002/aic.17721>
- [22] Pérez-Correa, I., Giunta, P. D., Francesconi, J. A., Mariño, F. J. "Artificial neural network for the prediction of physical properties of organic compounds based on the group contribution method", *The Canadian Journal of Chemical Engineering*, 101(8), pp. 4771–4783, 2023.
<https://doi.org/10.1002/cjce.24788>
- [23] Joback, K. G., Reid, R. C. "Estimation Of Pure-Component Properties From Group-Contributions", *Chemical Engineering Communications*, 57(1–6), pp. 233–243, 1987.
<https://doi.org/10.1080/00986448708960487>
- [24] Constantinou, L., Gani, R. "New group contribution method for estimating properties of pure compounds", *AIChE Journal*, 40(10), pp. 1697–1710, 1994.
<https://doi.org/10.1002/aic.690401011>
- [25] Pérez Ponce, A.A., Salfate, I., Pulgar-Villaruel, G., Palma-Chilla, L., Lazzús, J. A. "New group contribution method for the prediction of normal melting points", *Journal of Engineering Thermophysics*, 22(3), pp. 226–235, 2013.
<https://doi.org/10.1134/S1810232813030065>
- [26] Stefanis, E., Constantinou, L., Panayiotou, C. "A Group-Contribution Method for Predicting Pure Component Properties of Biochemical and Safety Interest", *Industrial & Engineering Chemistry Research*, 43(19), pp. 6253–6261, 2004.
<https://doi.org/10.1021/ie0497184>
- [27] Yanowitz, J., Ratcliff, M. A., McCormick, R. L., Taylor, J. D., Murphy, M. J. "Compendium of Experimental Cetane Numbers", National Renewable Energy Laboratory, Golden, CO, USA, Rep. NREL/TP-5400-67585, 2014. [online] Available at: <https://www.nrel.gov/docs/fy17osti/67585.pdf> [Accessed: 12 July 2024]
- [28] IFA "GESTIS database of German Social Accident Insurance", [online] Available at: <https://gestis-database.dguv.de/> [Accessed: 12 July 2024]
- [29] Chemical Book "Chemical Book database", [online] Available at: <https://www.chemicalbook.com/> [Accessed: 12 July 2024]
- [30] Royal Society of Chemistry "ChemSpider database", [online] Available at: <https://www.chemspider.com/> [Accessed: 12 July 2024]
- [31] National Library of Medicine "PubChem database", [online] Available at: <https://pubchem.ncbi.nlm.nih.gov/> [Accessed: 12 July 2024]
- [32] Bertz, S. H. "The first general index of molecular complexity", *Journal of the American Chemical Society*, 103(12), pp. 3599–3601, 1981.
<https://doi.org/10.1021/ja00402a071>
- [33] Arellano-Treviño, M. A., Bartholet, D., To, A. T., Bartling, A. W., Baddour, F. G., Alleman, T. L., Christensen, E. D., ..., Ruddy, D. A. "Synthesis of Butyl-Exchanged Polyoxymethylene Ethers as Renewable Diesel Blendstocks with Improved Fuel Properties", *ACS Sustainable Chemistry & Engineering*, 9(18), pp. 6266–6273, 2021.
<https://doi.org/10.1021/acssuschemeng.0c09216>
- [34] An, G., Xia, Y., Xue, Z., Shang, H., Cui, S., Lu, C. "Combination of Theoretical and Experimental Insights into the Oxygenated Fuel Poly(oxymethylene) Dibutyl Ether from n-Butanol and Paraformaldehyde", *ACS Omega*, 7(3), pp. 3064–3072, 2022.
<https://doi.org/10.1021/acsomega.1c06452>

- [35] Lide, D. R. (ed.) "CRC Handbook of Chemistry and Physics: A Ready-Reference Book of Chemical and Physical Data", CRC Press, 2004. ISBN 9780849304859
- [36] Boyd, R. H. "Some physical properties of polyoxymethylene dimethyl ethers", *Journal of Polymer Science*, 50(153), pp. 133–141, 1961.
<https://doi.org/10.1002/pol.1961.1205015316>
- [37] Deutsch, D., Oestreich, D., Lautenschütz, L., Haltenort, P., Arnold, U., Sauer, J. "High Purity Oligomeric Oxymethylene Ethers as Diesel Fuels", *Chemie Ingenieur Technik*, 89(4), pp. 486–489, 2017.
<https://doi.org/10.1002/cite.201600158>
- [38] Patrini, R., Marchionna, M. "A process for the selective production of dialkyl-polyformals", Milan, Italy, EP 1 505 049 A1 (Ceased), 2005.
- [39] Lucas, S. P., Chan, F. L., Fioroni, G. M., Foust, T. D., Gilbert, A., Luecke, J., McEnally, C. S., Serdoncillo, J. J. A., Zdanowicz, A. J., Zhu, J., Windom, B. "Fuel Properties of Oxymethylene Ethers with Terminating Groups from Methyl to Butyl", *Energy & Fuels*, 36(17), pp. 10213–10225, 2022.
<https://doi.org/10.1021/acs.energyfuels.2c01414>