

# SUPPORTING INFORMATION

## Data-Driven, Explainable Machine Learning Model for Predicting Volatile Organic Compounds' Standard Vaporization Enthalpy

*José Ferraz-Caetano<sup>1</sup>, Filipe Teixeira<sup>2</sup>, M. Natália D. S. Cordeiro<sup>1</sup>*

1) LAQV-REQUIMTE – Department of Chemistry and Biochemistry – Faculty of Sciences, University of Porto - Rua do Campo Alegre, S/N, 4169-007 Porto, Portugal

2) CQUM – Centre of Chemistry, University of Minho, Campus de Gualtar, 4710-057 Braga, Portugal

**KEYWORDS:** Machine Learning; Supervised Learning; Standard Vaporization Enthalpy; Thermochemical prediction

## INDEX

Annex S1 – Repository and Code.....	<b>Error! Bookmark not defined.</b>
Annex S2 – Model Database Description.....	3
Annex S3 – RDKit Full Descriptor List .....	4
Annex S4 – ML Model Development.....	7
Annex S5 – Model Validation .....	8
Annex S6 – ML VOC Model Optimization.....	10
Bibliography .....	11

## Annex S1 – Repository and Code

**Table S1**

List of links for the model's database and developed code.

Model	Website
Vaporization Enthalpy Database	<a href="https://gitfront.io/r/ESA-Vanadium/HdCi2vFs5msj/VOC-EnthVapML/">https://gitfront.io/r/ESA-Vanadium/HdCi2vFs5msj/VOC-EnthVapML/</a>
Model Development, Statistical Results and External Database <sup>1</sup>	<a href="https://gitfront.io/r/ESA-Vanadium/HdCi2vFs5msj/VOC-EnthVapML/">https://gitfront.io/r/ESA-Vanadium/HdCi2vFs5msj/VOC-EnthVapML/</a>

### Statistical definitions

- Mathematical definition of the average absolute relative deviation (AARD), root mean square error (RMSE) and mean absolute error (MAE):

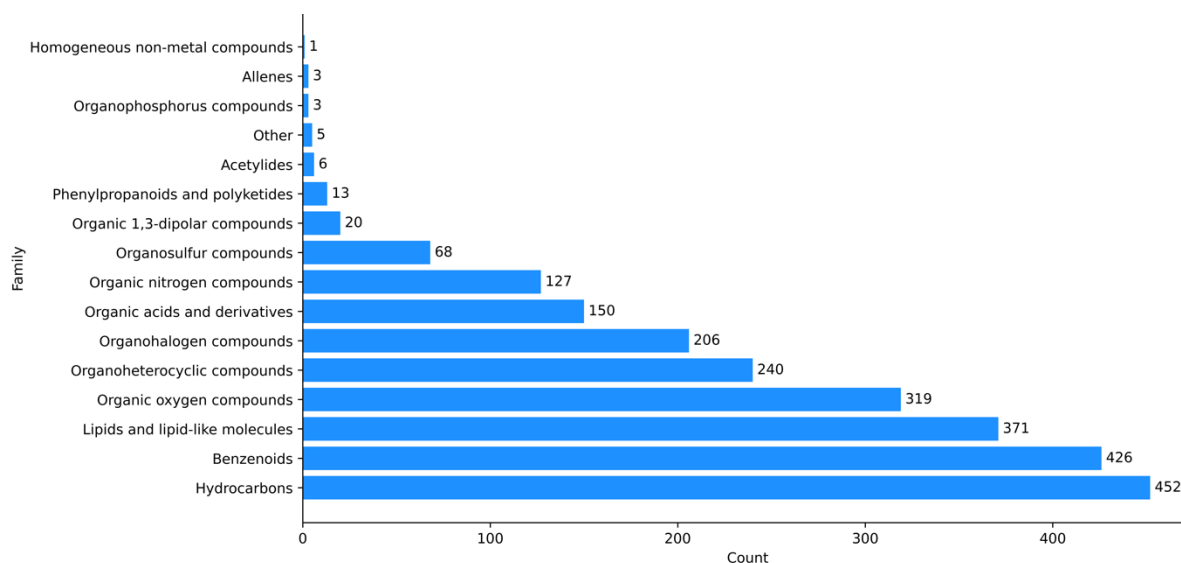
$$AARD = \frac{100}{N} \sum_i^N \frac{|prediction(i) - literature(i)|}{literature(i)}$$

$$RMSE = \sqrt{\frac{\sum_{i=1}^N (prediction(i) - literature(i))^2}{N}}$$

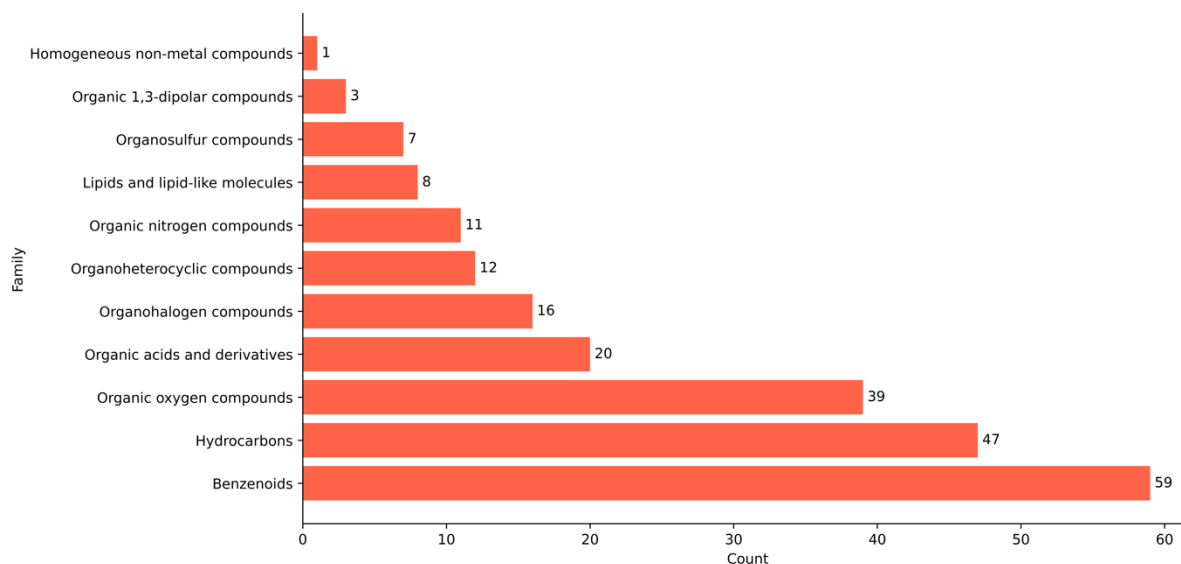
$$MAE = \frac{\sum_{i=0}^N |prediction(i) - literature(i)|}{N}$$

## Annex S2 – Model Database Description

The list of compounds used in the model's database and the VOC database<sup>1</sup> validation are presented in the Online Repository in Annex 1. In Figure S1 and S2, we present a graphical representation of different chemical families presented in the model database sorted by the overall number of entries.



**Figure S1.** Graphical representation of entries in the Model Database, sorted by molecular families.



**Figure S2.** Graphical representation of entries in the VOC Database, sorted by molecular families.

## Annex S3 – RDKit Full Descriptor List

Descriptors used in Machine Learning model development (Table S2). The respective bibliographical foundation of each descriptor is described in the RDKit WebBook Documentation<sup>2</sup>. The descriptor calculations were made by converting each SMILES string (representing each molecular entry), running on top of Python version 3.9<sup>3</sup> and using the RDKit package, version 2022.09.4<sup>4</sup>.

**Table S2**

List of Descriptors Used in Machine Learning model development

Descriptor	Type	Meaning
<b>BalabanJ</b>	Structural	Distance sum of the two end-vertex for each edge. BalabanJ index has been proven to be relevant to network branching.
<b>BertzCT</b>	Structural	Complexity index, considering both the variety of kinds of bond connectivities and atom types; information contents related to bond connectivity and atom type diversity.
<b><math>0\chi</math>, <math>1\chi</math></b>	Structural	This descriptor signifies a retention index (zero order) derived directly from gradient retention times.
<b><math>0\chi_n - 4\chi_n</math></b>	Structural	This descriptor signifies a retention index (zero order) derived directly from gradient retention times.
<b><math>0\chi_v - 4\chi_v</math></b>	Structural	This descriptor signifies atomic valence connectivity index (order 0).
<b>VSA_EState1 – VSA_EState11</b>	Surface Area	MOE-type (QSAR model) descriptors using EState indices and surface area contributions.

<b>Hall Kier <math>\alpha</math></b>	Structural	Descriptor fingerprint that displays the difference between active and inactive molecules.
<b>HeavyAtomCount</b>	Structural	The number of heavy atoms in the molecule.
<b>HeavyAtomMolWt</b>	Structural	The average molecular weight of the molecule ignoring hydrogens.
<b><math>\kappa_1, \kappa_2, \kappa_3</math></b>	Structural	This descriptor signifies # $\kappa$ shape index: $(n-1) \times 2 / m^2$
<b><i>Max and Min AbsEStateIndex</i></b>	Electronic	Maximum and minimum Absolute E-State
<b><i>Max and Min AbsPartialCharge</i></b>	Electronic	Maximum and minimum Absolute Partial Charge
<b><i>Max and Min EStateIndex</i></b>	Electronic	Maximum and minimum E-State
<b><i>Max and Min PartialCharge</i></b>	Electronic	Maximum and minimum Partial Charge
<b>Mol logP</b>	Structural	Wildman-Crippen logP value.
<b>MolMR</b>	Electronic	Wildman-Crippen molar refractivity.
<b>MolWt</b>	Structural	The average molecular weight of the molecule.
<b>NHOH Count</b>	Structural	The number of NHs or OHs.
<b>NO Count</b>	Structural	The number of Nitrogens and Oxygens.
<b><i>n<sub>alicarb</sub></i></b>	Structural	The number of aliphatic carbocycles.
<b><i>n<sub>alihet</sub></i></b>	Structural	The number of aliphatic heterocycles.
<b><i>n<sub>alirig</sub></i></b>	Structural	The number of aliphatic rings.

<b><i>n<sub>arocarb</sub></i></b>	Structural	The number of aromatic carbocycles.
<b><i>n<sub>arohet</sub></i></b>	Structural	The number of aromatic heterocycles.
<b><i>n<sub>arorig</sub></i></b>	Structural	The number of aromatic rings.
<b><i>n<sub>Ha</sub></i></b>	Structural	The number of Hydrogen Bond Acceptors.
<b><i>n<sub>Hd</sub></i></b>	Structural	The number of Hydrogen Bond Donors.
<b><i>n<sub>het</sub></i></b>	Structural	The number of Heteroatoms.
<b><i>n<sub>radele</sub></i></b>	Structural	The number of radical electrons.
<b><i>n<sub>rot</sub></i></b>	Structural	The number of Rotatable Bonds.
<b><i>n<sub>satrig</sub></i></b>	Structural	The number of saturated rings.
<b><i>n<sub>ele</sub></i></b>	Structural	The number of valence electrons.
<b>PEOE_VSA1 – PEOE_VSA14</b>	Surface Area	MOE-type (QSAR model) descriptors using partial charges and surface area contributions.
<b>SMR_VSA1 – SMR_VSA10</b>	Surface Area	MOE-type (QSAR model) descriptors using MR contributions and surface area contributions.
<b>SlogP_VSA1 – SlogP_VSA12</b>	Surface Area	MOE-type (QSAR model) descriptors using logP contributions and surface area contributions.
<b>TPSA</b>	Surface Area	The total polar surface area of a molecule based upon fragment calculations.
<b>VSA_EState1 – VSA_EState10</b>	Surface Area	MOE-type (QSAR model) descriptors using EState indices and surface area contributions.

76

77

78

79

## Annex S4 – ML Model Development

Initial model setup for  $\Delta_{\text{vap}}H_{\text{m}}^{\circ}$  prediction was assessed with increasing training loads to determine a train/test split. Considering the trade-off between result accuracy, model overfitting and data representability, we have selected a 60:40 train/test split for model optimization. We thus calculated the statistical performance of the Random Forest Regressor, Gradient Boosting Regressor, Support Vector Machines and MLP Neural Network algorithms, comparing with Linear Regressor, as presented in Table S3.

**Table S3**

Statistical results for model performance for each algorithm for a 60:40 train/test split: prediction scores, mean absolute error, average absolute relative deviation, random mean square error and coefficient of determination.

Regressor Algorithm	R <sup>2</sup> Test	R <sup>2</sup> Train	Score Test	MAE / kJ mol <sup>-1</sup>	RMSE / kJ mol <sup>-1</sup>	AARD / %
Random Forest	0.969 ± 0.004	0.9951 ± 0.0007	0.968 ± 0.005	3.5 ± 0.2	5.8 ± 0.4	5.8 ± 0.2
Gradient Boosting	0.969 ± 0.005	0.988 ± 0.001	0.968 ± 0.005	3.8 ± 0.1	5.8 ± 0.4	6.3 ± 0.2
MLP Neural Network	0.96 ± 0.02	0.98 ± 0.02	0.95 ± 0.03	3.7 ± 0.5	6.8 ± 2.0	6.7 ± 1.9
Linear Regression	0.93 ± 0.01	0.93 ± 0.01	0.93 ± 0.01	6.1 ± 0.3	8.9 ± 0.5	10.3 ± 0.4
SVM	0.55 ± 0.05	0.54 ± 0.03	0.51 ± 0.05	10.0 ± 0.5	22.1 ± 2.5	16.4 ± 0.7

The RF, GB, MLP algorithms compare favorably in overall accuracy with reported models for thermodynamic predictions, which average MAE values<sup>5-7</sup> under 4 kJ mol<sup>-1</sup> (even models using a single family dataset<sup>8</sup>). Our results strongly suggest that all three models benefit from a wide and diverse set of chemical species, while RF, GB and MLP present more accurate predictions than the SVM algorithm and linear regression.

Considering our best supervised algorithms (GB, RG and MLP), we argue a positive outlook on further model development using these three algorithms. Key difference between a Neural Network and our best supervised algorithms stems from a wide variation of RMSE, calculated with a high prediction error, not compatible with sustained recurrent calculations. Hence a sustained perspective on model optimization using the ensemble algorithms.

## Annex S5 – Model Validation

To evaluate the wide scope of this model, we have predicted the standard vaporization enthalpy of a known database used for non-ML prediction methods<sup>1</sup>. This external database is composed of several hydrocarbon-based molecules, not drawn from our original model database. Table S4 and Figure S10 display an evaluation of our model against previously reported studies.

**Table S4**

Statistical results for model performance of each optimized algorithm for  $\Delta_{\text{vap}}H_{\text{m}}^{\circ}$  prediction of the external dataset.

Regressor Algorithm	R <sup>2</sup> Test	MAE / kJ mol <sup>-1</sup>	RMSE / kJ mol <sup>-1</sup>	AARD / %
Random Forest	0.964	2.815	3.860	4.748
Gradient Boosting	0.950	3.531	4.643	6.030
Linear Regressor	0.965	2.144	2.700	4.240
Support Vector Machine	0.953	2.669	3.567	4.889
<b>MLP Neural Network</b>	<b>0.982</b>	<b>1.379</b>	<b>2.124</b>	<b>2.480</b>

**Table S5**

Ranking of method performance for  $\Delta_{\text{vap}}H_{\text{m}}^{\circ}$  predictions across the same external dataset.

Model	AARD / %
Guthrie and Taylor <sup>11</sup>	10.6
Ducros <sup>12</sup>	7.0
Domalski and Hearing <sup>13</sup>	6.4
Constantinou and Gani <sup>14</sup>	5.5
Chickos <sup>15</sup>	4.6
Kolská <sup>16</sup>	4.1
Gharagheizi <sup>1</sup>	3.1
Santos and Leal <sup>17</sup>	2.6
<b>Our ML Model</b>	<b>2.5</b>

Our ML model has predicted the standard vaporization enthalpies of the selected molecules in the database with an AARD in line with the best previous attempts. The low AARD displayed by our model is even more significant if we consider the background of the compared studies. For example, more accurate models can be only applied for predictions of hydrocarbons (as Santos and Leal's model), which relinquishes any attempt to apply such model to any generic



molecule. Also, some models only work if specific types of chemical descriptors are known (such as Abraham's solvation features), or depend on functional/fragment group values, while our proposed model has a faster, wider scope of chemical diversity in the original dataset. For example, while Gharagheizi presents one of the most diverse initial datasets, it does not include representative counts of key chemical families like Lipids, Salts, Organometallic and Organic acid compounds, as we include in our database.

## Annex S6 – ML VOC Model Optimization

Prior to model development we evaluated each tested algorithms performance for predicting VOCs  $\Delta_{\text{vap}}H_{\text{m}}^{\circ}$ . We present the results in Table S6, where we see that the RF yielded the best results. We then used a pipeline strategy to automatically select and tune each algorithms' best features for  $\Delta_{\text{vap}}H_{\text{m}}^{\circ}$  determinations, scoring with the lowest prediction MAE possible. Sample code is presented in the Annex S1. Using the selected hyper-parameters, the model is optimized by removing offsetting descriptors, which do not contribute positively towards model prediction ( $\text{PI} > 0\%$ ). Our models were recalculated accordingly and statistical results with value evaluations from the firsts iterations are presented in the manuscript's Table 1.

**Table S6**

Statistical results for model performance of each optimized algorithm for VOC  $\Delta_{\text{vap}}H_{\text{m}}^{\circ}$  prediction.

Regressor Algorithm	R <sup>2</sup> Test	MAE / kJ mol <sup>-1</sup>	RMSE / kJ mol <sup>-1</sup>	AARD / %
<b>Random Forest</b>	<b>0.946</b>	<b>3.072</b>	<b>4.785</b>	<b>6.362</b>
Gradient Boosting	0.943	3.754	5.116	7.885
Linear Regressor	0.933	3.636	5.315	8.049
Support Vector Machine	0.651	8.005	12.085	17.561
MLP Neural Network	0.937	3.167	5.286	6.322

### List of optimal hyperparameters for RF algorithm for model optimization:

Random Forest Regressor : n\_estimators= 300, min\_samples\_split= 2, min\_samples\_leaf= 1, max\_features='sqrt', max\_depth= 20, bootstrap= False, random\_state=47

## Bibliography

1. Gharagheizi, F.; Ilani-Kashkouli, P.; Acree, W. E.; Mohammadi, A. H.; Ramjugernath, D., A group contribution model for determining the vaporization enthalpy of organic compounds at the standard reference temperature of 298K. *Fluid Phase Equilibria* **2013**, *360*, 279-292.
2. RDKit: Open-source cheminformatics - Descriptptor Webbook. <https://www.rdkit.org/docs/GettingStartedInPython.html#list-of-available-descriptors> (accessed March 1, 2022).
3. Python Software Foundation - Python Language Reference, version 3.9.8. <http://www.python.org>.
4. Landrum, G. RDKit: Open-source cheminformatics 2022\_09\_4 (Q3 2022) Release - January 16, 2023. <http://www.rdkit.org/> (accessed January 18, 2023).
5. Dobbelaere, M. R.; Plehiers, P. P.; Van de Vijver, R.; Stevens, C. V.; Van Geem, K. M., Learning Molecular Representations for Thermochemistry Prediction of Cyclic Hydrocarbons and Oxygenates. *The Journal of Physical Chemistry A* **2021**, *125* (23), 5166-5179.
6. Yalamanchi, K. K.; van Oudenhoven, V. C. O.; Tutino, F.; Monge-Palacios, M.; Alshehri, A.; Gao, X.; Sarathy, S. M., Machine Learning To Predict Standard Enthalpy of Formation of Hydrocarbons. *The Journal of Physical Chemistry A* **2019**, *123* (38), 8305-8313.
7. Aldosari, M. N.; Yalamanchi, K. K.; Gao, X.; Sarathy, S. M., Predicting entropy and heat capacity of hydrocarbons using machine learning. *Energy and AI* **2021**, *4*, 100054.
8. Benson, S. W.; Buss, J. H., Additivity Rules for the Estimation of Molecular Properties. Thermodynamic Properties. *The Journal of Chemical Physics* **1958**, *29* (3), 546-572.
9. Churchill, B.; Acree, W. E.; Abraham, M. H., Development of Abraham model expressions for predicting the standard molar enthalpies of vaporization of organic compounds at 298.15 K. *Thermochimica Acta* **2019**, *681*, 178372.

- 217 10. Low, K.; Coote, M. L.; Izgorodina, E. I., Explainable Solvation Free Energy Prediction  
218 Combining Graph Neural Networks with Chemical Intuition. *Journal of Chemical Information*  
219 *and Modeling* **2022**, 62 (22), 5457-5470.
- 220 11. Guthrie, J. P.; Taylor, K. F., Additivity methods for estimating heats of vaporization of  
221 organic compounds. *Canadian Journal of Chemistry* **1983**, 61 (3), 602-607.
- 222 12. Ducros, M.; Gruson, J. F.; Sannier, H., Estimation des enthalpies de vaporisation des  
223 composés organiques liquides. Partie 1. Applications aux alcanes, cycloalcanes, alcènes,  
224 hydrocarbures benzeniques, alcools, alcanes thiols, chloro et bromoalcanes, nitriles, esters,  
225 acides et aldehydes. *Thermochimica Acta* **1980**, 36 (1), 39-65.
- 226 13. Domalski, E. S.; Hearing, E. D., Estimation of the Thermodynamic Properties of C-H-  
227 N-O-S-Halogen Compounds at 298.15 K. *Journal of Physical and Chemical Reference Data*  
228 **1993**, 22 (4), 805-1159.
- 229 14. Constantinou, L.; Gani, R., New group contribution method for estimating properties  
230 of pure compounds. *AIChE Journal* **1994**, 40 (10), 1697-1710.
- 231 15. Chickos, J. S.; Hesse, D. G.; Liebman, J. F., Estimating vaporization enthalpies of  
232 organic compounds with single and multiple substitution. *The Journal of Organic Chemistry*  
233 **1989**, 54 (22), 5250-5256.
- 234 16. Kolská, Z.; Růžička, V.; Gani, R., Estimation of the Enthalpy of Vaporization and the  
235 Entropy of Vaporization for Pure Organic Compounds at 298.15 K and at Normal Boiling  
236 Temperature by a Group Contribution Method. *Industrial & Engineering Chemistry Research*  
237 **2005**, 44 (22), 8436-8454.
- 238 17. Santos, R. C.; Leal, J. P., A Review on Prediction Methods for Molar Enthalpies of  
239 Vaporization of Hydrocarbons: The ELBA Method as the Best Answer. *Journal of Physical*  
240 *and Chemical Reference Data* **2012**, 41 (4), 043101.

241