

Supporting Information: Fast Evaluation of the Adsorption Energy of Organic Molecules on Metals via Graph Neural Networks

Sergio Pablo-García^{1,2,†}, Santiago Morandi^{3,4,†}, Rodrigo A.
Vargas-Hernández^{1,5}, Kjell Jorner^{1,2,6}, Núria López^{3,*}, and Alán
Aspuru-Guzik^{1,2,*}

¹Department of Chemistry, University of Toronto, Lash Miller Chemical
Laboratories 80 St. George Street, ON M5S 3H6, Toronto, Canada

²Department of Computer Science, University of Toronto, Sandford Fleming
Building, 40 St. George Street, ON M5S 2E4, Toronto, Canada

³Institute of Chemical Research of Catalonia, The Barcelona Institute of Science
and Technology, Av. Països Catalans 16, 43007, Tarragona, Spain

⁴Department of Physical and Inorganic Chemistry, Universitat Rovira i Virgili,
Campus Sescelades, N4 Block, C. Marcel·lí Domingo 1, 43007, Tarragona, Spain

⁵Vector Institute for Artificial Intelligence, 661 University Ave. Suite 710, ON
M5G 1M1, Toronto, Canada

⁶Department of Chemistry and Chemical Engineering, Chalmers University of
Technology, Kemigården 4, SE-412 96, Gothenburg, Sweden

[†]Authors Equally Contributed to this Work

*Corresponding Authors: email: nlopez@icq.es; alan@aspuru.com

October 5, 2022

Notes

S1	FG-dataset	4
S2	Adsorption Conformational Search	5
S3	Automation of DFT Data Generation	5
S4	Graph Representation Algorithm	7
S5	GNN Model Architecture	9
S6	GNN Training	10
S7	BM-dataset	12
S8	Hyperparameter Optimization	13

Figures

S1	Example of graph data structure representation.	15
S2	Data cleaning workflow applied to the raw graph FG-dataset.	16
S3	Graph representations of the FG-dataset without metal-adsorbate connections. a) Distribution by chemical family and b) by metal.	17
S4	GNN model architecture before hyperparameter optimization.	18
S5	GNN model architecture after hyperparameter optimization.	19
S6	Stratified data splitting procedure.	20
S7	Cross validation approach for estimating the GNN generalization performance.	21
S8	Learning rate and MAE of the train/validation/test sets during the training processes of the cross validation.	22
S9	a) Box-plot of the test error distribution sorted by metal and b) mean error and standard error of the mean of the predictions obtained by the cross validation models.	23
S10	BM-dataset: Biomass molecules and metals.	24
S11	BM-dataset: Polyurethane molecules and metals.	25
S12	BM-dataset: Plastic molecules and metals.	26

Tables

S1	FG-dataset: Hydrocarbons, alcohols, aldehydes, ketones and ethers.	27
S2	FG-dataset: Carbonates, carboxylic acids and esters.	28
S3	FG-dataset: Amines and imines.	29
S4	FG-dataset: Amidines.	30

S5	FG-dataset: Thiols, thials, thioketones and thioethers.	31
S6	FG-dataset: Amides.	32
S7	FG-dataset: Oximes.	32
S8	FG-dataset: Carbamate esters.	33
S9	FG-dataset: Aromatic molecules.	34
S10	Standard deviation and standard error of the mean of the prediction error of the models obtained by cross validation sorted by chemical family. Values are reported in eV.	35
S11	Standard deviation and standard error of the mean of the prediction error of the models obtained by cross validation sorted by metal (the last row considers the gas-phase subset). Values are reported in eV.	35
S12	Summary of the performed hyperparameter optimization. All the pooling-wise options refer to the GMT implementation in Pytorch Geometric (PyG).	36
S13	Adsorption energy of BM-dataset samples obtained with the optimized GNN and compared to DFT. Values are reported in eV.	37

Notes

Note S1: FG-dataset

The functional groups dataset (FG-dataset) created to develop the GNN model presented in this work includes 204 closed-shell molecules. For each of them, the FG-dataset contains the gas-phase molecule and the molecule adsorbed on 12 metal surfaces (Ag, Au, Cd, Cu, Ir, Ni, Os, Pd, Pt, Rh, Ru, Zn). One adsorption configuration is included for each molecule/metal combination except for the aromatics with one ring, which comprise two different configurations for each metal. This leads to a total of 2849 DFT samples in the FG-dataset. The set of molecules represents the most common functional groups: for each chemical family, all the existing configurations containing up to C_4 are included, except aromatics where $C_{>4}$ are considered. The chemical families included in the FG-dataset are the following:

- Alkanes, alkenes and alkynes. **Table S1**
- Alcohols, aldehydes, ketones and ethers (1 Oxygen). **Table S1**
- Carbonates, carboxylic acids and esters (2–3 Oxygens). **Table S2**
- Amines and imines (1 Nitrogen). **Table S3**
- Amidines (2 Nitrogens). **Table S4**
- Thiols, thials, thioketones and thioethers (1 Sulfur). **Table S5**
- Amides and oximes (1 Oxygen + 1 Nitrogen). **Table S6–S7**

- Carbamate esters (2 Oxygens + 1 Nitrogen). **Table S8**
- Aromatic molecules with up to two rings containing O, N and S. **Table S9**

Note S2: Adsorption Conformational Search

Even simple C₂₋₃ adsorbates could have approximately 100 conformations. Thus, for the initial DFT adsorption geometries we followed a simplified conformational analysis based on the heuristic rules devised in Refs. [1, 2]. These rules can be summarised as follows:

- (i). The unsaturated bonds were placed close to the surface.
- (ii). Heteroatoms (O, N and S) were placed close to the surface.
- (iii). Carbon tails face the surface.
- (iv). If the intermediate did not converge to a reasonable structure, the molecule was readjusted manually, trying up to 6 conformations that preserve the rules (i-iii).

Note S3: Automation of DFT Data Generation

To build and obtain the E_{DFT} of the adsorbate/metal combinations included in the FG-dataset, we executed the following procedure:

- (i). First, the metallic surfaces were built starting from their respective bulks.

- (ii). A list of SMILES [3] of the molecules in the FG-dataset has been fed to Open Babel [4] in order to generate the .xyz geometry files of the molecules. The applied force field is the MMFF94 [5].
- (iii). The .xyz geometry files have been converted with pymatgen [6] to POSCAR files representing the molecules in cubic cells with dimension of 20 Å. The gas-phase molecules have been calculated with this initial geometry.
- (iv). The relaxed molecules were copied and placed at a distance of 2.0–2.3 Å from the Rh slab, our reference metal. The conformational rules explained in **Note S2** have been followed.
- (v). The resulting structures were sent to our computational resources and relaxed with VASP. Once full convergence has been reached, the geometry of the relaxed molecules on Rh has been extracted and automatically placed on the slabs of the all the other metals preserving the adsorbate distance from the Rh slab.
- (vi). The obtained structures were checked for conformation errors, and if needed, manually built and relaxed again. Typical problems encountered were related to the fragmentation of the adsorbate during the relaxation due to the unstable initial geometry.

For the BM-Dataset, the molecules were manually built, adsorbed on the metal surfaces and relaxed using VASP, including the gas phase molecules.

Structures obtained from both datasets were uploaded to ioChem-BD.[7, 8]

Note S4: Graph Representation Algorithm

To convert the three-dimensional structures obtained from DFT to their respective graph representation, we applied a modified version of the algorithm presented in Ref. [9], implemented in pyRDTP. To define the neighbors of each atom, the algorithm reads the relaxed three-dimensional atomic positions from the geometry file (CONTCAR) generated by VASP, and applies the Voronoi tessellation method. This method creates a partition of the three-dimensional space, defining a region for each atom that consists of all points of the space closer to that atom than to any other. Two atoms are considered connected if they share a Voronoi facet and their distance is less than the sum of the atomic covalent radii plus a tolerance distance. In this work we used as covalent radii those provided by Cordero et al. [10] multiplied by a scaling factor of 1.5 for metals, and a tolerance parameter of 0.5 Å to help in the detection of metal-adsorbate connections.

Once the connectivity is defined, the graphs are generated representing the atoms as nodes and the detected connections as edges. The metal atoms not directly connected to the adsorbate are not considered during the graph generation process. The atomic elements are embedded to the nodes using the one-hot encoding approach as implemented in Scikit-Learn (**Figure S1**).[11] This step is needed to convert categorical variables (as atomic elements) into ML-suitable data structures. This algorithm is applied to all the samples in both FG-dataset and BM-dataset.

For a fraction of the adsorption systems, the graph conversion results into inaccurate representations, as for specific geometries the M—A distance is so

high that the algorithm is unable to properly define their connectivity. This is due to the fact that the algorithm is based on a purely geometrical criterion that defines the edges relying on a set of empirical covalent radii, which are fixed for each element and do not account for all the possible phenomena occurring in catalytic systems. Thus, the strategy for minimizing the amount of bad representations has been to fine tune the tolerance parameter and the covalent radii scaling factor. In order to discard inaccurate graph representations and properly curate the graph dataset to be suitable for the GNN model training, we implemented the following sieves (**Figure S2**):

- (i). A first filter that discards the graphs representing adsorption configurations without the presence of metal atoms.
- (ii). A second filter that verifies the correct connectivity of C and H atoms within the molecules: Connectivity of carbon atoms is properly defined if the number of edges connecting it to other atoms in the molecule is equal or less than 4, while hydrogen atoms are correctly connected if its number of edges is exactly one.
- (iii). A filter to prevent the inclusion in the dataset of DFT samples that contain more than one adsorbate on the slab or with a final geometry in which the adsorbate has dissociated in multiple fragments.
- (iv). A last filter for removing duplicate graphs deriving from the presence of stereo-isomers adsorbed in the same configuration on the metal surface.

The amount of graphs pruned out after the first two filters intrinsically depends on the graph conversion algorithm: a higher applied tolerance pa-

parameter would reduce the number of discarded graphs after the first filter, but at the same time it increases the number of removed graphs in the second filter, due to the creation of nonphysical connections within the adsorbates.

Note S5: GNN Model Architecture

The architecture of the GNN model developed in this work is shown in **Figure S4**. The input graphs are represented by a set of node feature vectors, each of them being a 17-dimensional array (12 metals + C, H, O, N and S) needed for representing the chemical element via the one-hot encoder, and by the graph connectivity in coordinate format (**Figure S1**). In order to transform this mathematical graph representation into the prediction of the DFT energy of the associated chemical system, the following transformations are applied in the listed order:

- (i). First, each node feature vector is transformed via one fully connected layer which increases the dimensionality of the vector from 17 to 128, $\mathbb{B}^{17} \rightarrow \mathbb{R}^{128}$, where \mathbb{B} denotes the Boolean space used to define the nodes via one hot encoding.
- (ii). Then, two additional fully connected layers are applied, keeping the dimensionality of the nodes ($\mathbb{R}^{128} \rightarrow \mathbb{R}^{128}$). Up to now, no information about the graph connectivity is exploited.
- (iii). Three GraphSAGE [12] convolutional layers are applied to all the nodes to capture the information from the neighbors by exploiting the graph connectivity. Between each convolution, a fully-connected layer is ap-

plied.

- (iv). Finally, the information embedded in the nodes is compressed into a global graph representation with the Graph Multiset Transformer (GMT),[13] a global pooling layer which returns back a scalar value, namely the prediction of the DFT energy of the chemical system represented by the initial input graph. $\mathbb{R}^{N \times 128} \rightarrow \mathbb{R}$, where N is the number of nodes in the graph.

The activation function used in the node-level layers (all except the pooling layer) is the rectified linear unit (ReLU). [14] In total, the developed GNN model contains 398,081 trainable parameters, 224,513 of them belonging to the pooling layer due to its high internal complexity and the remaining parameters equally distributed among the other layers.

Note S6: GNN Training

The model training has been performed minimizing the mean absolute error (MAE) as loss function with the ADAM algorithm as optimizer.[15] The learning rate (LR) value is steered with the reduce-on-plateau scheduler. The initial learning rate has been set to 10^{-3} and is reduced exponentially by the scheduler every time in which there is no improvement after 5 epochs (patience) in the MAE of the validation set. The minimum LR possible has been set to 10^{-6} , while the decrease factor has been set to 0.7. In each training 250 epochs are performed. During each epoch, the training set is fed to the model in batches of 32 samples, performing a backward propagation and updating the model parameters after each batch.

The performance of the GNN model has been assessed by applying a stratified splitting of the FG-dataset by chemical families followed by a 5-fold cross validation. The first (**Figure S6**) allows a proper distribution of all the chemical families among the splits, while the latter provides a robust estimation of the GNN generalization performance. The cross validation approach follows the process depicted in **Figure S7**: After the partition in 5 stratified splits, each split is employed as test set: for each split used as test set, 4 splits are left and each of them is employed as validation set. This leads to a final validation approach consisting of 20 independent learning processes, each one performed with a unique combination of the splits among train, validation and test set. The generalization performance of the model is finally assessed by averaging the MAE of the absolute errors of the learning processes performed employing the same test set. **Figure S8** shows the values of the learning rate and MAE among the trained models. It is important to mention that after the creation of the train, validation and test sets in each learning process, a standardized target scaling is applied using the mean and standard deviation of the energy values of the samples from the train and validation sets, discarding the samples from the test set as its inclusion would lead to a data leakage. The target scaling is essential to ensure a stable learning process: if the target variable among the graph dataset has a large spread of values, it may result in large error gradients causing model parameters to change dramatically, making the learning process unstable.

Note S7: BM-dataset

The big molecules dataset (BM-dataset) used for testing the GNN model includes three industrially relevant materials:

- Biomass. Ref. [16], **Figure S10**
- Polyurethane precursors. Ref. [17], **Figure S11**
- Plastics. Ref. [18] **Figure S12**

These three groups consist of complex chemical structures that can be seen as combinations of the functional groups present in the FG-dataset. For each group, 5 representative molecules of larger size compared to those in the FG-dataset have been selected and relaxed through DFT to simulate the gas-phase and adsorption configuration on 2 metals chosen according to the existing applications and studies. The BM-dataset (45 samples, 30 adsorptions + 15 gas-phase) is used as an additional test for assessing the GNN performance on samples coming from a distribution distinct from that used to build the model (FG-dataset). For example, for the plastic group we represent polyethylene (PE), polypropylene (PP, both syndio- and isotactic), polystyrene (PS) and polyethylene terephthalate (PET) as molecules composed by a reasonable number of monomers. We generated the DFT adsorption systems of these molecules on Pt and Ru metal surfaces as these represent potential candidates for applications related to chemical recycling technology.

Note S8: Hyperparameter Optimization

Before testing the proposed model on the BM-dataset, we performed a hyperparameter optimization study to explore the vast space defined by all the variables that can affect the final performance of the GNN model. The hyperparameters can be defined as all the variables that are not model parameters, but that affect model performance at the same extent. These can be divided in two groups:

- Training-related: they define the training process. Examples are learning rate, optimization algorithm, batch size, etc.
- Model-related: they define the architecture of the model and include kind of activation function (ReLU, Tanh, etc.), number and depth of layers, bias inclusion, etc.

We adopted the hyperband asynchronous algorithm (ASHA) [19] implemented in RayTune. [20] ASHA combines random search and aggressive early stopping in order to optimize the hyperparameters and is based on the proved fact that in order to find the best hyperparameter settings, just a small amount of iterations (epochs) is sufficient to discriminate between bad models and promising candidates.

The hyperparameter space, shown in **Table S12** has been investigated picking randomly 2000 different settings and for each of those, a model training has been run by ASHA with a grace period of 15 epochs (e.g., the poorly-performing models are discarded after having trained them for a minimum of 15 iterations) and a maximum of 200 for the best ones. The final hyper-

parameter setting is the one that minimizes the MAE of the energy of the samples belonging to the BM-dataset.

Figures

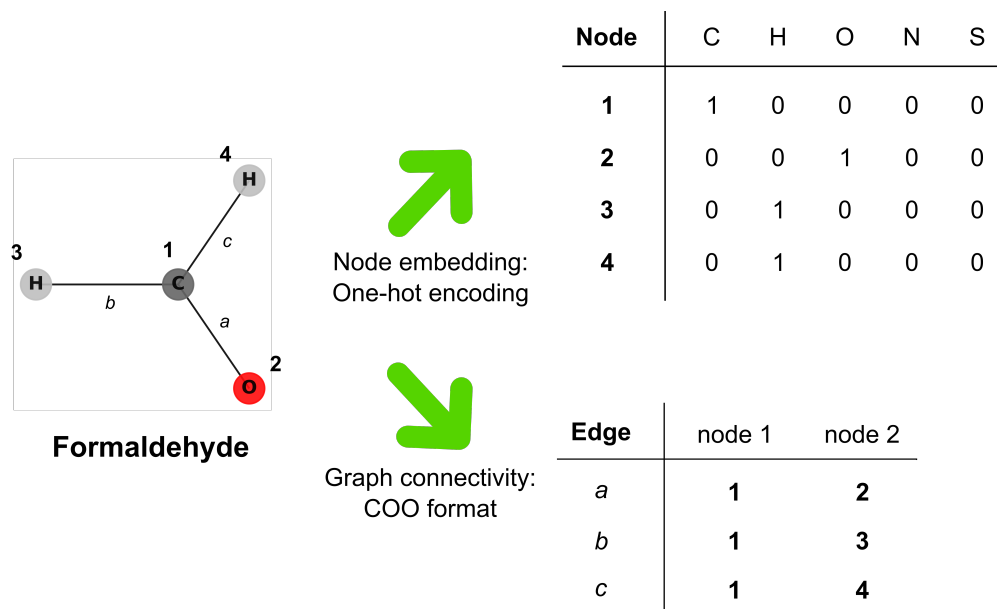


Figure S1: Example of graph data structure representation.

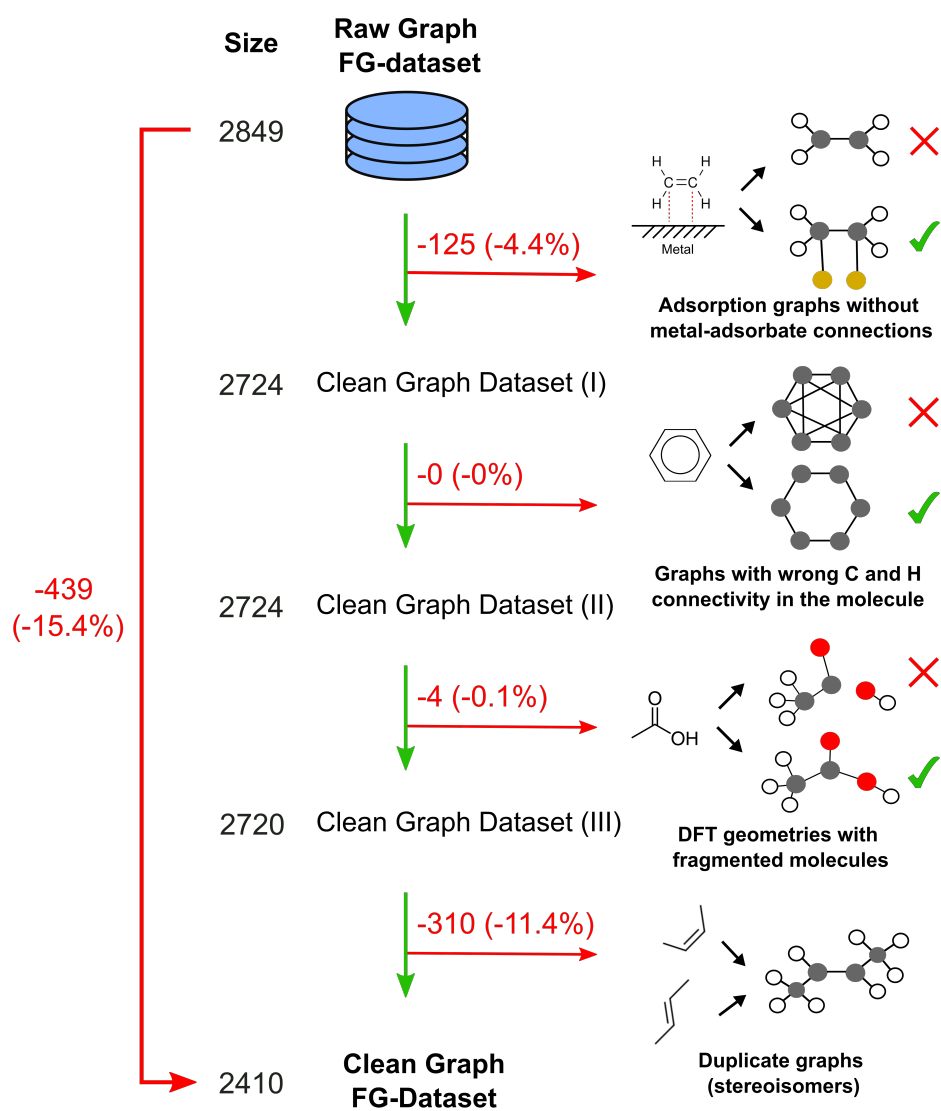


Figure S2: Data cleaning workflow applied to the raw graph FG-dataset.

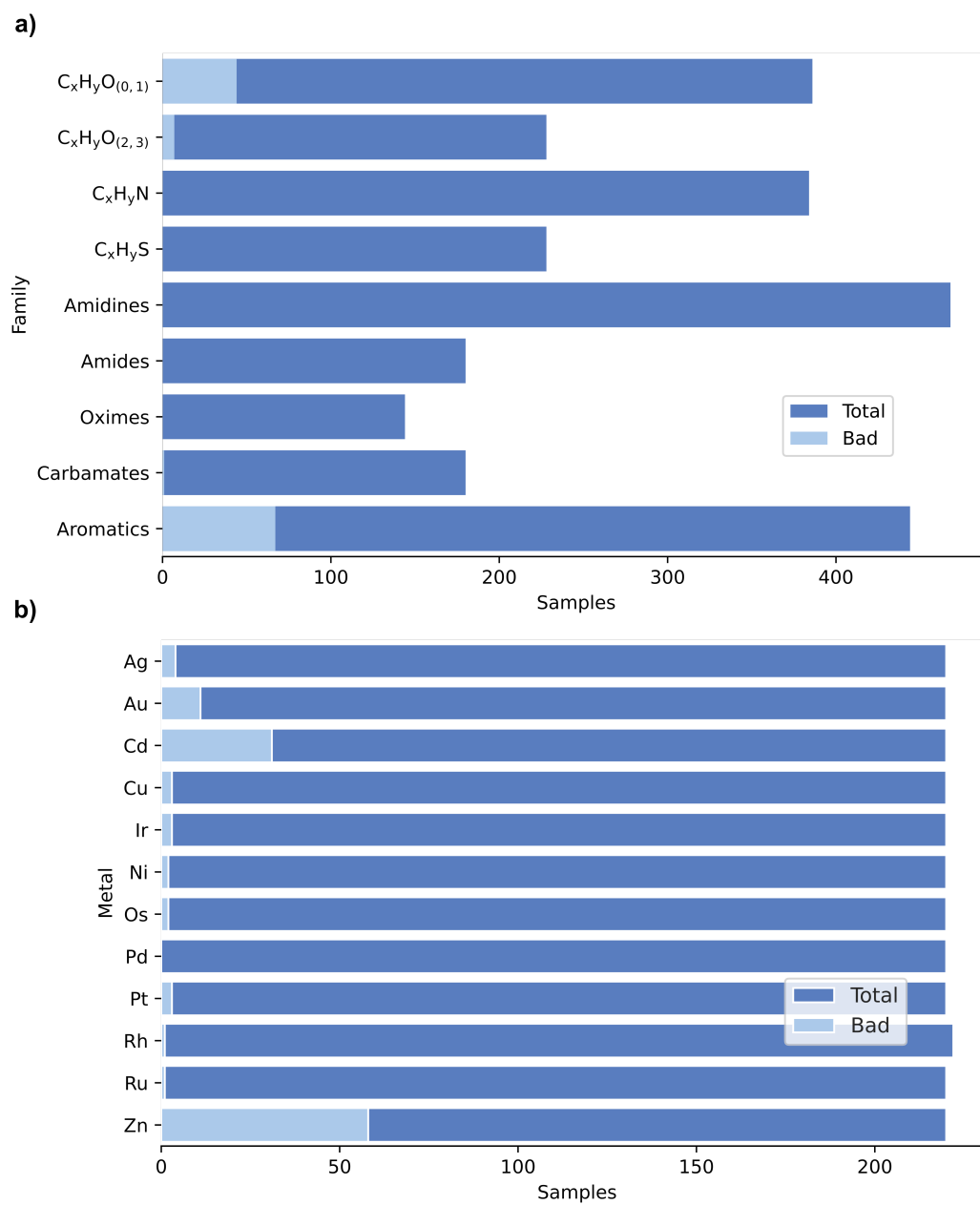


Figure S3: Graph representations of the FG-dataset without metal-adsorbate connections. a) Distribution by chemical family and b) by metal.

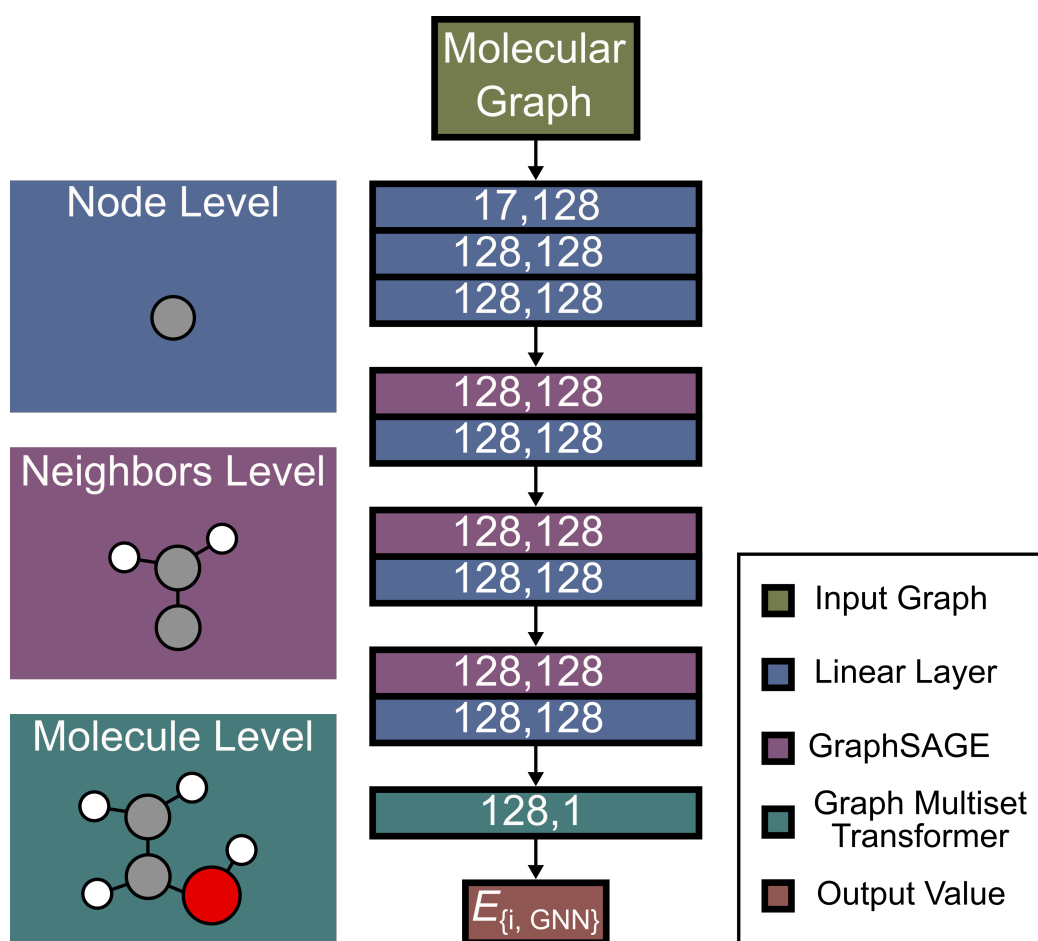


Figure S4: GNN model architecture before hyperparameter optimization.

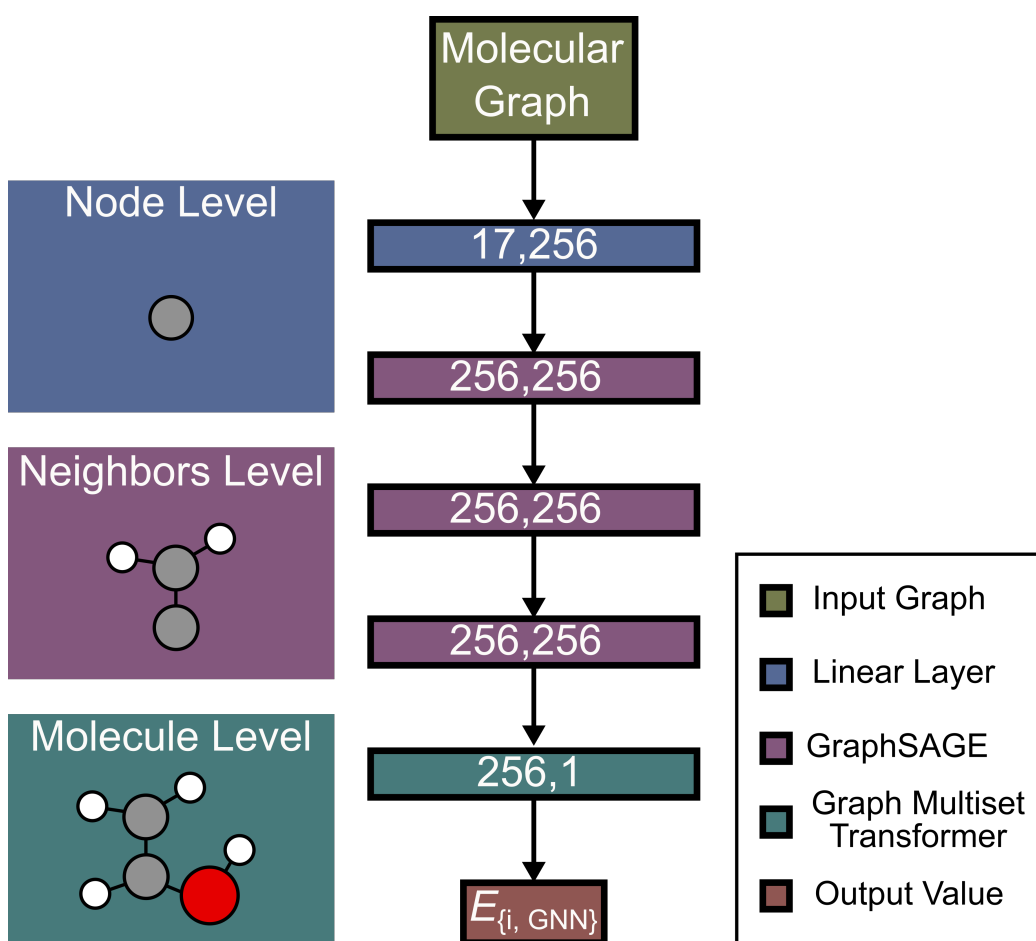


Figure S5: GNN model architecture after hyperparameter optimization.

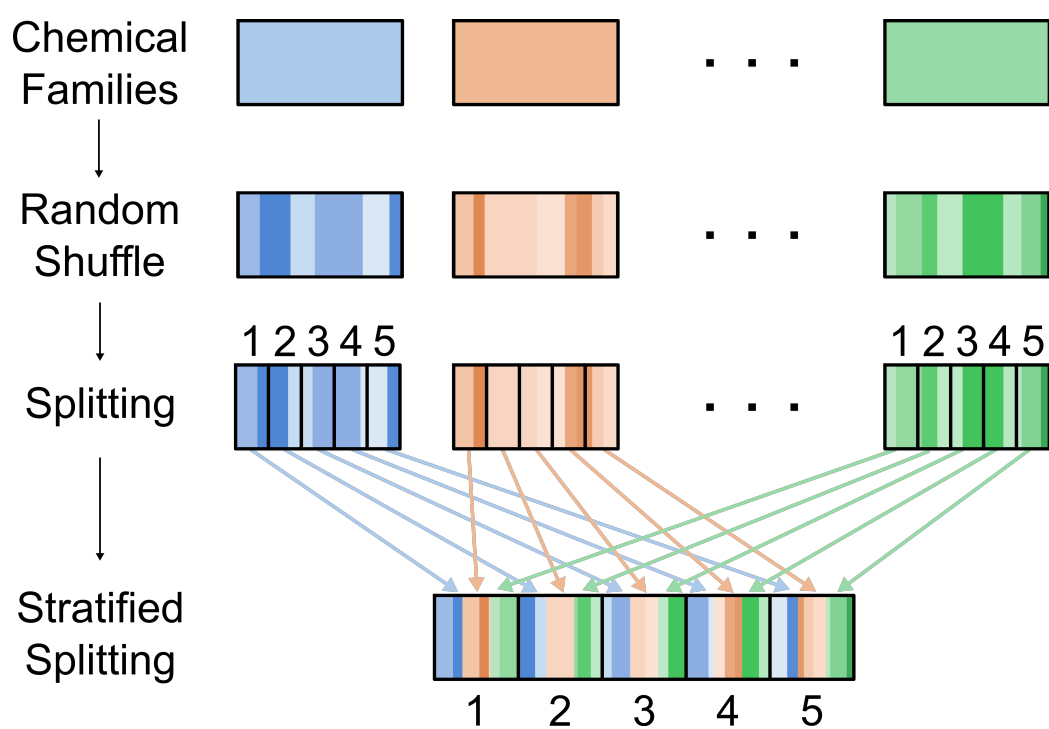


Figure S6: Stratified data splitting procedure.

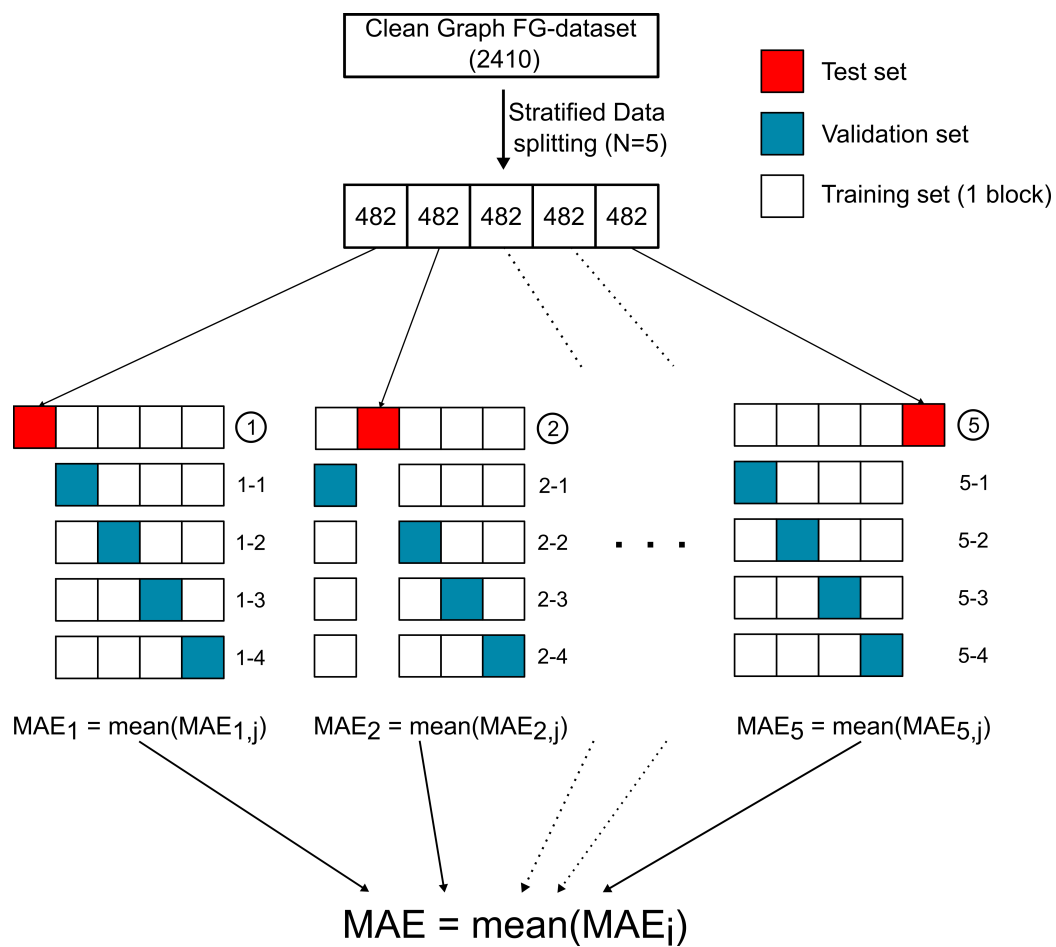


Figure S7: Cross validation approach for estimating the GNN generalization performance.

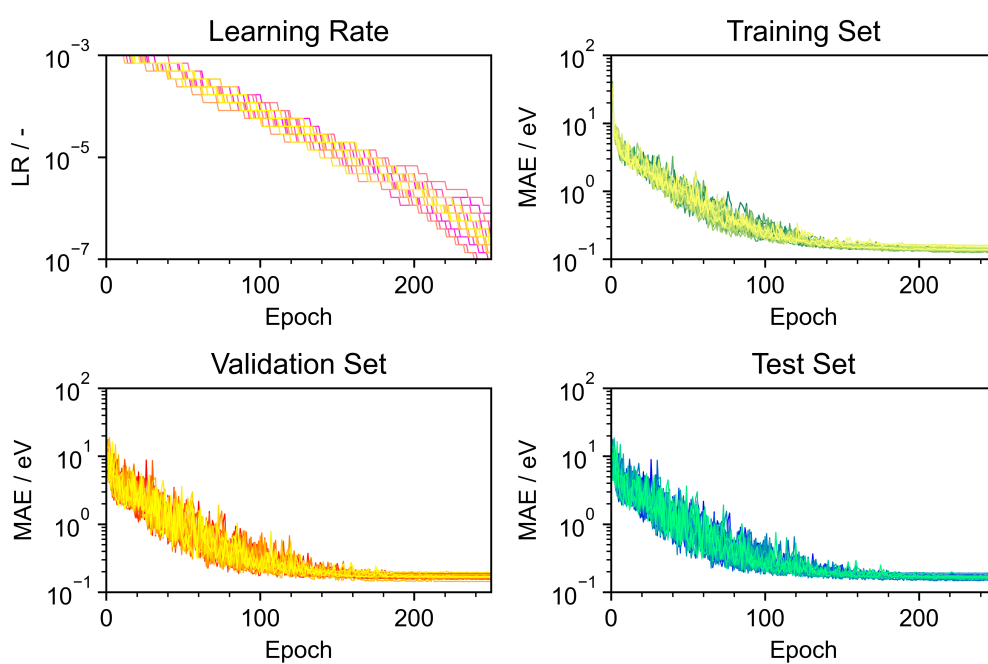


Figure S8: Learning rate and MAE of the train/validation/test sets during the training processes of the cross validation.

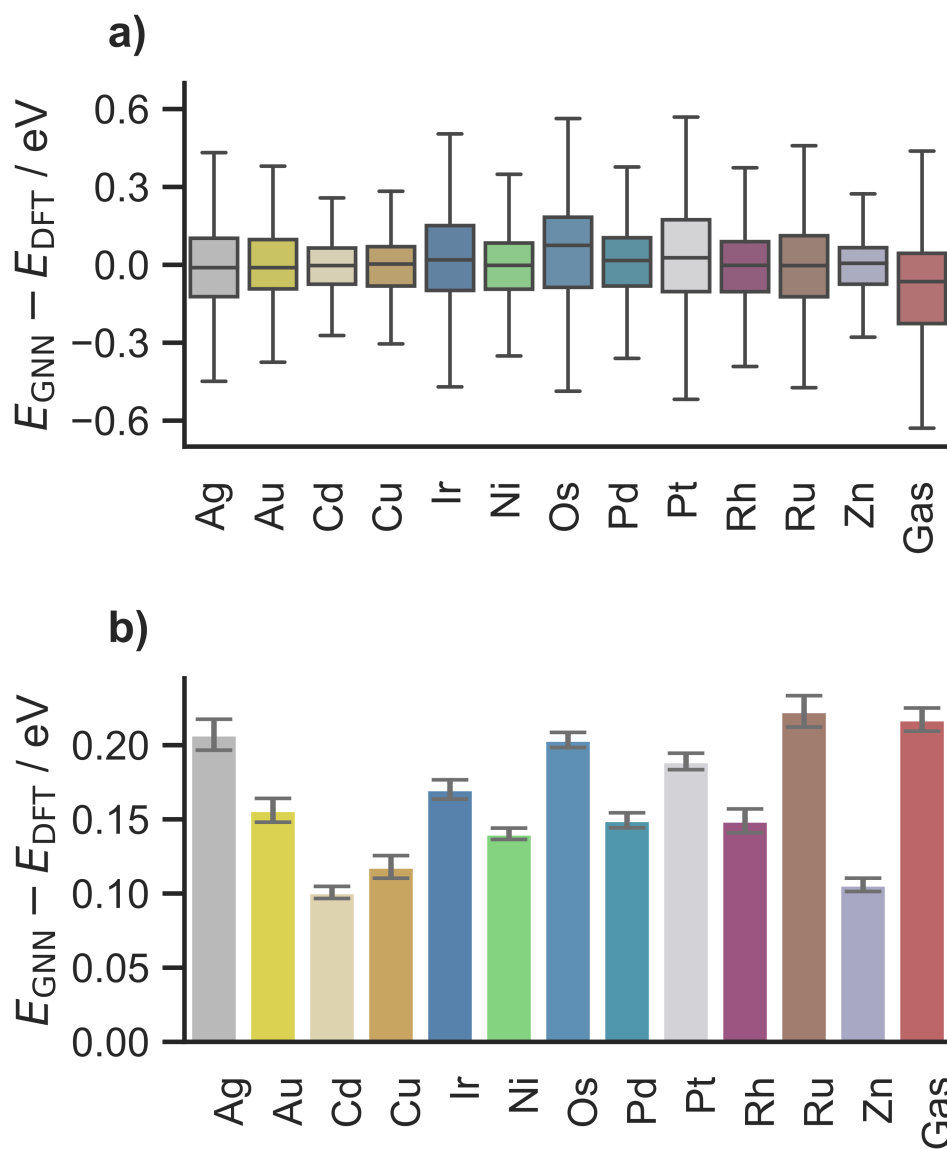
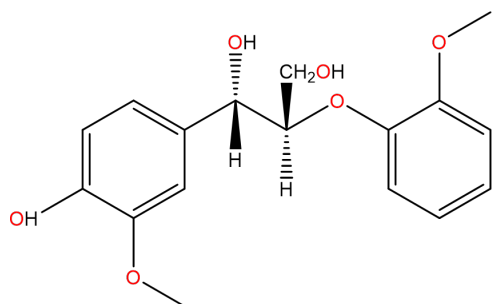
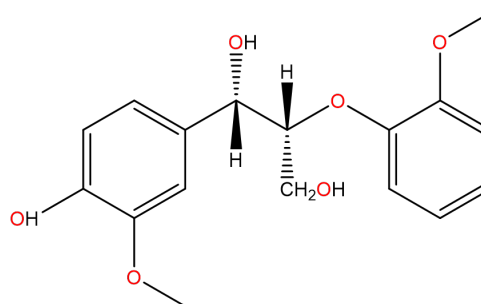


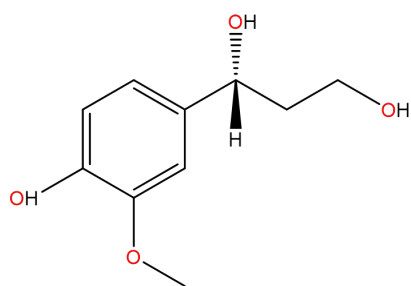
Figure S9: a) Box-plot of the test error distribution sorted by metal and b) mean error and standard error of the mean of the predictions obtained by the cross validation models.



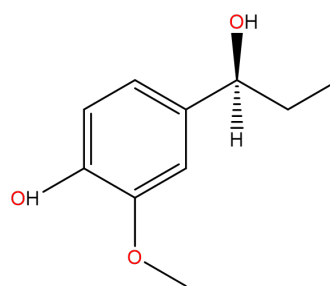
(1*S*,2*S*)-1-(4-hydroxy-3-methoxyphenyl)-2-(2-methoxyphenoxy)propane-1,3-diol



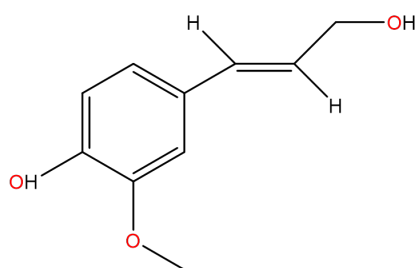
(1*S*,2*R*)-1-(4-hydroxy-3-methoxyphenyl)-2-(2-methoxyphenoxy)propane-1,3-diol



(*R*)-1-(4-hydroxy-3-methoxyphenyl)propane-1,3-diol



(*S*)-4-(1-hydroxypropyl)-2-methoxyphenol

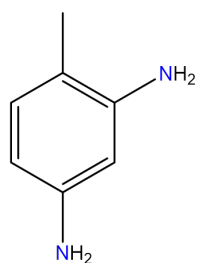


(*E*)-4-(3-hydroxyprop-1-en-1-yl)-2-methoxyphenol

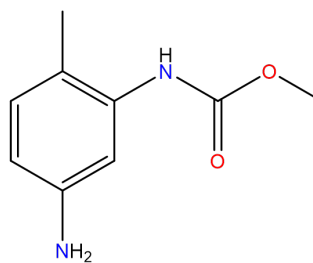
Surfaces



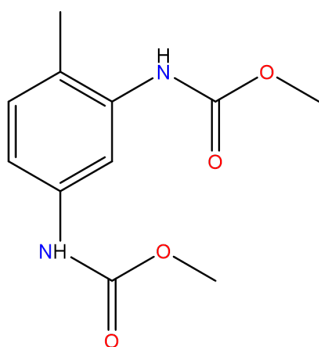
Figure S10: BM-dataset: Biomass molecules and metals.



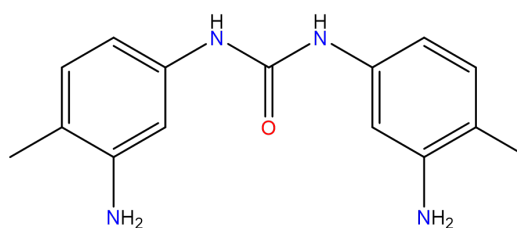
4-methylbenzene-1,3-diamine



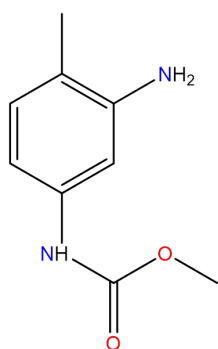
methyl (5-amino-2-methylphenyl)carbamate



dimethyl (4-methyl-1,3-phenylene)dicarbamate



1,3-bis(3-amino-4-methylphenyl)urea



methyl (3-amino-4-methylphenyl)carbamate

Surfaces



Figure S11: BM-dataset: Polyurethane molecules and metals.

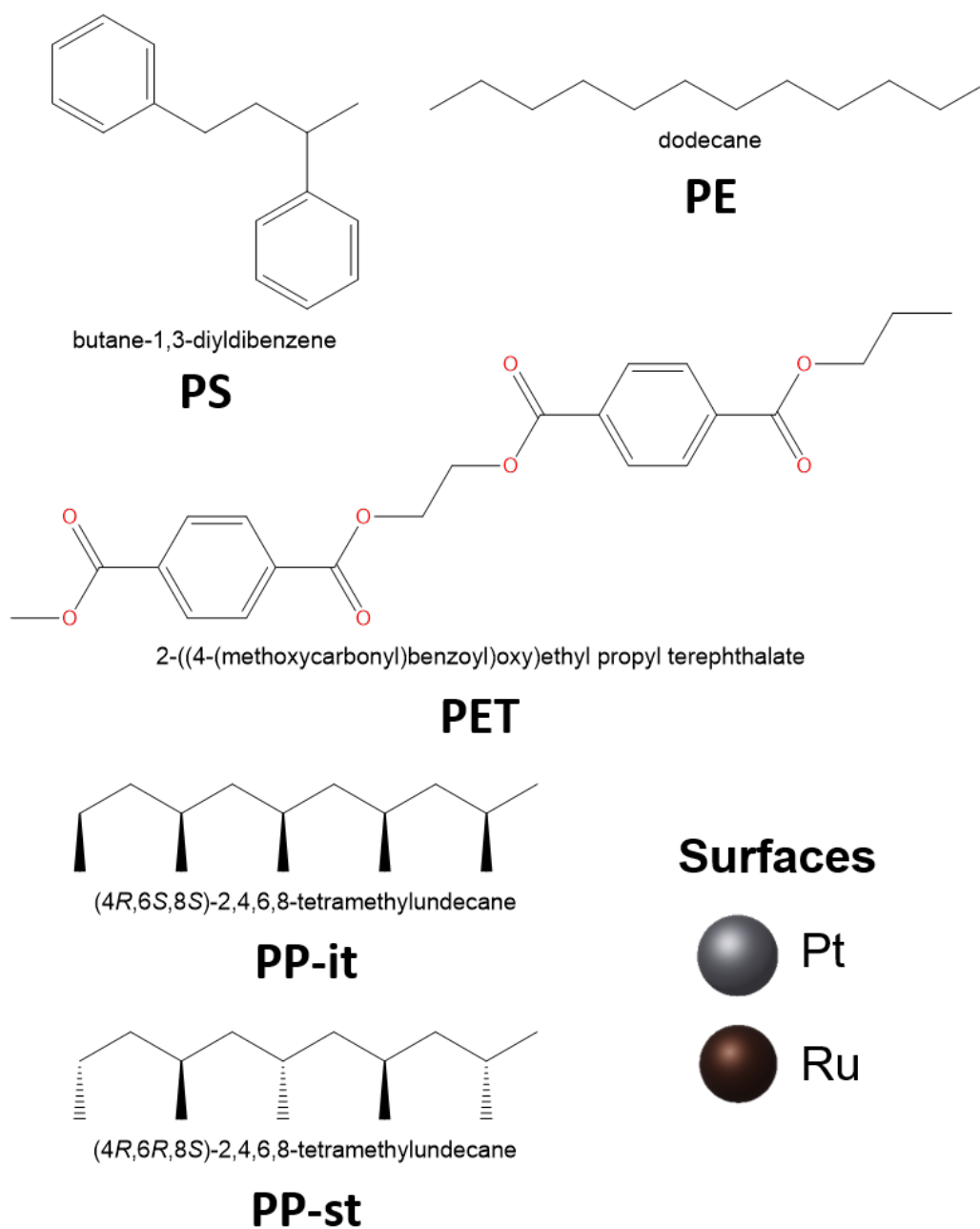


Figure S12: BM-dataset: Plastic molecules and metals.

Tables

Table S1: FG-dataset: Hydrocarbons, alcohols, aldehydes, ketones and ethers.

IUPAC name	SMILES
Formaldehyde	<chem>C=O</chem>
Acetylene	<chem>C#C</chem>
Ethylene	<chem>C=C</chem>
Acetaldehyde	<chem>CC=O</chem>
Ethane	<chem>CC</chem>
Dimethylether	<chem>COC</chem>
Ethanol	<chem>CCO</chem>
Propyne	<chem>C#CC</chem>
Propylene	<chem>C=CC</chem>
Acetone	<chem>CC(=O)C</chem>
Propane	<chem>CCC</chem>
Ethylmethylether	<chem>CCOC</chem>
N-propanol	<chem>CCCO</chem>
I-propanol	<chem>CC(C)O</chem>
2-butyne	<chem>CC#CC</chem>
1-butyne	<chem>C#CCC</chem>
1-butene	<chem>C=CCC</chem>
2-butene-cis	<chem>C/C=C\C</chem>
2-butene-trans	<chem>C/C=C/C</chem>
I-butene	<chem>CC(=C)C</chem>
Butanone	<chem>CCC(=O)C</chem>
N-butane	<chem>CCCC</chem>
I-butane	<chem>CC(C)C</chem>
N-butanol	<chem>CCCCO</chem>
2-butanol-R	<chem>C[C@H](CC)O</chem>
2-butanol-S	<chem>C[C@@H](CC)O</chem>
T-butanol	<chem>CC(C)(C)O</chem>
Propionaldehyde	<chem>CCC=O</chem>
Butyraldehyde	<chem>CCCC=O</chem>
Isobutyraldehyde	<chem>CC(C=O)C</chem>
Ethoxyethane	<chem>CCOCC</chem>
1-methoxypropane	<chem>CCCOCC</chem>

Table S2: FG-dataset: Carbonates, carboxylic acids and esters.

IUPAC name	SMILES
Formic acid	<chem>C(=O)O</chem>
Carbonic acid	<chem>C(=O)(O)O</chem>
Acetic acid	<chem>C(=O)(C)O</chem>
Methyl formate	<chem>C(=O)OC</chem>
Methyl hydrogen carbonate	<chem>C(=O)(O)OC</chem>
Propionic acid	<chem>CCC(=O)O</chem>
Ethyl formate	<chem>C(=O)OCC</chem>
Methyl acetate	<chem>C(=O)(C)OC</chem>
Ethyl hydrogen carbonate	<chem>C(=O)(OCC)O</chem>
Dimethyl carbonate	<chem>C(=O)(OC)OC</chem>
Butyric acid	<chem>CCCC(=O)O</chem>
Isobutyric acid	<chem>CC(C(=O)O)C</chem>
Propyl formate	<chem>C(=O)OCCC</chem>
Isopropyl formate	<chem>C(=O)OC(C)C</chem>
Ethyl acetate	<chem>CC(=O)OCC</chem>
Methyl propionate	<chem>CCC(=O)OC</chem>
Propyl hydrogen carbonate	<chem>C(=O)(OCCC)O</chem>
Isopropyl hydrogen carbonate	<chem>C(=O)(O)OC(C)C</chem>
Ethyl methyl carbonate	<chem>C(=O)(OC)OCC</chem>

Table S3: FG-dataset: Amines and imines.

IUPAC name	SMILES
Methanimine	<chem>C=N</chem>
Methanamine	<chem>CN</chem>
Ethanimine	<chem>CC=N</chem>
Methylmethanimine	<chem>C=NC</chem>
Ethanamine	<chem>CCN</chem>
Dimethylamine	<chem>CNC</chem>
Propan-1-imine	<chem>CCC=N</chem>
N-ethylmethanimine	<chem>C=NCC</chem>
(E)-N-methylethanamine	<chem>C/C=N/C</chem>
(Z)-N-methylethanamine	<chem>C/C=N\C</chem>
Propan-2-imine	<chem>CC(=N)C</chem>
Propan-1-amine	<chem>CCCN</chem>
Propan-2-amine	<chem>CC(C)N</chem>
N-methylethanamine	<chem>CCNC</chem>
Trimethylamine	<chem>CN(C)C</chem>
Butan-1-imine	<chem>CCCC=N</chem>
2-methylpropan-1-imine	<chem>CC(C)C=N</chem>
N-propylmethanimine	<chem>C=NCCC</chem>
N-isopropylmethanimine	<chem>C=NC(C)C</chem>
(E)-N-ethylethanamine	<chem>C/C=N/CC</chem>
(Z)-N-ethylethanamine	<chem>C/C=N\CC</chem>
(E)-N-methylpropan-1-imine	<chem>CC/C=N/C</chem>
(Z)-N-methylpropan-1-imine	<chem>CC/C=N\C</chem>
Butan-2-imine	<chem>CCC(=N)C</chem>
N-methylpropan-2-imine	<chem>CC(=NC)C</chem>
Butan-1-amine	<chem>CCCCN</chem>
(R)-butan-2-amine	<chem>C[C@H](CC)N</chem>
(S)-butan-2-amine	<chem>C[C@@H](CC)N</chem>
2-methylpropan-2-amine	<chem>CC(C)(C)N</chem>
Diethylamine	<chem>CCNCC</chem>
N-methylpropan-1-amine	<chem>CCCNC</chem>
N,N-dimethylethanamine	<chem>CCN(C)C</chem>

Table S4: FG-dataset: Amidines.

IUPAC name	SMILES
Formimidamide	<chem>C(=N)N</chem>
N-methylformimidamide	<chem>C(=N)NC</chem>
acetimidamide	<chem>C(=N)(C)N</chem>
(E)-N ¹ -methylformimidamide	<chem>C(=[NH2])[N]C</chem>
(Z)-N ¹ -methylformimidamide	<chem>C(=[NH2])[N]C</chem>
N-ethylformimidamide	<chem>C(=N)NCC</chem>
Propionimidamide	<chem>C(=N)(CC)N</chem>
(E)-N ¹ -ethylformimidamide	<chem>C(=[NH2])[N]CC</chem>
(Z)-N ¹ -ethylformimidamide	<chem>C(=[NH2])[N]CC</chem>
(E)-N ¹ -methylacetimidamide	<chem>C(=[NH2])(C)[N]C</chem>
(Z)-N ¹ -methylacetimidamide	<chem>C(=[NH2])(C)[N]C</chem>
N,N-dimethylformimidamide	<chem>C(=N)N(C)C</chem>
N-methylacetimidamide	<chem>C(=N)(C)NC</chem>
(E)-N,N ¹ -dimethylformimidamide	<chem>C(=N\C)/NC</chem>
(Z)-N,N ¹ -dimethylformimidamide	<chem>C(=N\C)\NC</chem>
N-isopropylformimidamide	<chem>C(=N)NC(C)C</chem>
N-propylformimidamide	<chem>C(=N)NCCC</chem>
Isobutyrimidamide	<chem>C(=N)(C(C)C)N</chem>
Butyrimidamide	<chem>C(=N)(CCC)N</chem>
(E)-N ¹ -isopropylformimidamide	<chem>C(=[NH2])[N]C(C)C</chem>
(Z)-N ¹ -isopropylformimidamide	<chem>C(=[NH2])[N]C(C)C</chem>
(E)-N ¹ -propylformimidamide	<chem>C(=[NH2])[N]CCC</chem>
(Z)-N ¹ -propylformimidamide	<chem>C(=[NH2])[N]CCC</chem>
(E)-N ¹ -ethylacetimidamide	<chem>C(=[NH2])(C)[N]CC</chem>
(Z)-N ¹ -ethylacetimidamide	<chem>C(=[NH2])(C)[N]CC</chem>
(E)-N ¹ -methylpropionimidamide	<chem>C(=[NH2])(CC)[N]C</chem>
(Z)-N ¹ -methylpropionimidamide	<chem>C(=[NH2])(CC)[N]C</chem>
N-ethyl-N-methylformimidamide	<chem>C(=N)N(C)CC</chem>
N-methylpropionimidamide	<chem>C(=N)(CC)NC</chem>
N-ethylacetimidamide	<chem>C(=N)(C)NCC</chem>
N,N-dimethylacetimidamide	<chem>C(=N)(C)N(C)C</chem>
(E)-N ¹ -ethyl-N-methylformimidamide	<chem>C(=N\CC)/NC</chem>
(Z)-N ¹ -ethyl-N-methylformimidamide	<chem>C(=N\CC)\NC</chem>
(E)-N-ethyl-N ¹ -methylformimidamide	<chem>C(=N\C)/NCC</chem>
(Z)-N-ethyl-N ¹ -methylformimidamide	<chem>C(=N\C)\NCC</chem>
(E)-N,N,N ¹ -trimethylformimidamide	<chem>C(=N\C)/N(C)C</chem>
(Z)-N,N,N ¹ -trimethylformimidamide	<chem>C(=N\C)\N(C)C</chem>
(E)-N,N ¹ -dimethylacetimidamide	<chem>C(=N\C)(\C)/NC</chem>
(Z)-N,N ¹ -dimethylacetimidamide	<chem>C(=N\C)(/C)\NC</chem>

Table S5: FG-dataset: Thiols, thials, thioketones and thioethers.

IUPAC name	SMILES
Methanethial	<chem>C=S</chem>
Methanethiol	<chem>CS</chem>
Ethanethial	<chem>CC=S</chem>
Ethanethiol	<chem>CCS</chem>
Dimethylsulfane	<chem>CSC</chem>
Propanethial	<chem>CCC=S</chem>
Propane-2-thione	<chem>CC(=S)C</chem>
Propane-1-thiol	<chem>CCCS</chem>
Propane-2-thiol	<chem>CC(C)S</chem>
Ethyl(methyl)sulfane	<chem>CCSC</chem>
Butanethial	<chem>CCCC=S</chem>
2-methylpropanethial	<chem>CC(C)C=S</chem>
Butane-2-thione	<chem>CCC(=S)C</chem>
Butane-1-thiol	<chem>CCCCS</chem>
(R)-butane-2-thiol	<chem>C[C@H](CC)S</chem>
(S)-butane-2-thiol	<chem>C[C@@H](CC)S</chem>
2-methylpropane-2-thiol	<chem>CC(C)(C)S</chem>
Diethylsulfane	<chem>CCSCC</chem>
Methyl(propyl)sulfane	<chem>CCCSC</chem>

Table S6: FG-dataset: Amides.

IUPAC name	SMILES
Formamide	<chem>C(=O)N</chem>
N-methylformamide	<chem>C(=O)NC</chem>
Acetamide	<chem>C(=O)(C)N</chem>
N-ethylformamide	<chem>C(=O)NCC</chem>
Propionamide	<chem>C(=O)(CC)N</chem>
N,N-dimethylformamide	<chem>C(=O)N(C)C</chem>
N-methylacetamide	<chem>C(=O)(C)NC</chem>
N-isopropylformamide	<chem>C(=O)NC(C)C</chem>
N-propylformamide	<chem>C(=O)NCCC</chem>
Isobutyramide	<chem>C(=O)(C(C)C)N</chem>
Butyramide	<chem>C(=O)(CCC)N</chem>
N-ethyl-N-methylformamide	<chem>C(=O)N(C)CC</chem>
N-ethylacetamide	<chem>C(=O)(C)NCC</chem>
N-methylpropionamide	<chem>C(=O)(CC)NC</chem>
N,N-dimethylacetamide	<chem>C(=O)(C)N(C)C</chem>

Table S7: FG-dataset: Oximes.

IUPAC name	SMILES
Formaldehyde oxime	<chem>C=NO</chem>
(E)-acetaldehyde oxime	<chem>C(=N\O)/C</chem>
(Z)-acetaldehyde oxime	<chem>C(=N\O)\C</chem>
(E)-propionaldehyde oxime	<chem>C(=N\O)/CC</chem>
(Z)-propionaldehyde oxime	<chem>C(=N\O)\CC</chem>
Propan-2-one oxime	<chem>C(=NO)(C)C</chem>
(E)-butyraldehyde oxime	<chem>C(=N\O)/CCC</chem>
(Z)-butyraldehyde oxime	<chem>C(=N\O)\CCC</chem>
(E)-isobutyraldehyde oxime	<chem>C(=N\O)/C(C)C</chem>
(Z)-isobutyraldehyde oxime	<chem>C(=N\O)\C(C)C</chem>
(E)-butan-2-one oxime	<chem>C(=N\O)(\C)/CC</chem>
(Z)-butan-2-one oxime	<chem>C(=N\O)(\CC)/C</chem>

Table S8: FG-dataset: Carbamate esters.

IUPAC name	SMILES
carbamic acid	<chem>C(=O)(N)O</chem>
methylcarbamic acid	<chem>C(=O)(NC)O</chem>
methyl carbamate	<chem>C(=O)(N)OC</chem>
ethylcarbamic acid	<chem>C(=O)(NCC)O</chem>
ethyl carbamate	<chem>C(=O)(N)OCC</chem>
dimethylcarbamic acid	<chem>C(=O)(N(C)C)O</chem>
methyl methylcarbamate	<chem>C(=O)(NC)OC</chem>
isopropylcarbamic acid	<chem>C(=O)(NC(C)C)O</chem>
propylcarbamic acid	<chem>C(=O)(NCCC)O</chem>
isopropyl carbamate	<chem>C(=O)(N)OC(C)C</chem>
propyl carbamate	<chem>C(=O)(N)OCCC</chem>
ethyl(methyl)carbamic acid	<chem>C(=O)(N(CC)C)O</chem>
methyl ethylcarbamate	<chem>C(=O)(NCC)OC</chem>
ethyl methylcarbamate	<chem>C(=O)(NC)OCC</chem>
methyl dimethylcarbamate	<chem>C(=O)(N(C)C)OC</chem>

Table S9: FG-dataset: Aromatic molecules.

IUPAC name	SMILES
Furan	<chem>o1cccc1</chem>
Thiophene	<chem>s1cccc1</chem>
Pyrrole	<chem>[nH]1cccc1</chem>
Pyridine	<chem>c1ccncc1</chem>
Cyclopentadiene	<chem>C1C=CC=C1</chem>
Benzene	<chem>c1ccccc1</chem>
Phenol	<chem>Oc1ccccc1</chem>
Thiophenol	<chem>Sc1ccccc1</chem>
Aniline	<chem>Nc1ccccc1</chem>
Toluene	<chem>Cc1ccccc1</chem>
Para-xylene	<chem>Cc1ccc(C)cc1</chem>
Meta-xylene	<chem>Cc1cccc(C)c1</chem>
Ortho-xylene	<chem>Cc1ccccc1C</chem>
Benzofuran	<chem>o1ccc2ccccc12</chem>
Isobenzofuran	<chem>o1cc2ccccc2c1</chem>
Benzo[b]thiophene	<chem>s1ccc2ccccc12</chem>
Benzo[c]thiophene	<chem>s1cc2ccccc2c1</chem>
1H-indole	<chem>[nH]1ccc2ccccc12</chem>
2H-indole	<chem>C1C=C2C=CC=CC2=N1</chem>
Quinoline	<chem>c1ccc2ncccc2c1</chem>
Isoquinoline	<chem>c1ccc2ncccc2c1</chem>
1H-indene	<chem>C1C=Cc2ccccc12</chem>
2H-indene	<chem>C1C=C2C=CC=CC2=C1</chem>
Naphthalene	<chem>c1ccc2ccccc2c1</chem>

Table S10: Standard deviation and standard error of the mean of the prediction error of the models obtained by cross validation sorted by chemical family. Values are reported in eV.

Chemical Family	Standard Deviation	Standard Error
$C_xH_yO_{(0,1)}$	0.26	0.0073
$C_xH_yO_{(2,3)}$	0.13	0.0024
C_xH_yN	0.23	0.0027
C_xH_yS	0.26	0.0044
Amidines	0.23	0.0039
Amides	0.15	0.0046
Oximes	0.19	0.0043
Carbamates	0.16	0.0035
Aromatics	0.48	0.0102

Table S11: Standard deviation and standard error of the mean of the prediction error of the models obtained by cross validation sorted by metal (the last row considers the gas-phase subset). Values are reported in eV.

Metal	Standard Deviation	Standard Error
Ag	0.23	0.0104
Au	0.16	0.0080
Cd	0.10	0.0041
Cu	0.12	0.0076
Ir	0.17	0.0065
Ni	0.14	0.0038
Os	0.20	0.0051
Pd	0.15	0.0051
Pt	0.19	0.0055
Rh	0.15	0.0080
Ru	0.22	0.0106
Zn	0.11	0.0045
Gas	0.22	0.0078

Table S12: Summary of the performed hyperparameter optimization. All the pooling-wise options refer to the GMT implementation in Pytorch Geometric (PyG).

Hyperparameter	Type	Search space	Optimum
Batch size	int	16, 32, 64, 128	16
Loss function	func	MAE, MSE	MAE
Initial lr	float	1e-1, 1e-2, 1e-3, 1e-4	1e-3
Lr-patience	int	5, 7, 9	5
Lr-factor	float	0.5, 0.7, 0.9	0.7
Minimum lr	float	1e-7, 1e-8, 1e-9	1e-8
Amsgrad	bool	True, False	True
Layers depth	int	64, 128, 256	256
Bias inclusion	bool	True, False	False
Linear layers	int	0, 1, 2, 3, 4	1
Convolutional layers	int	1, 2, 3, 4, 5	3
Convolution type	func	GraphSAGE, GATv2	GraphSAGE
Normalized conv.	bool	True, False	False
Root-weighted conv.	bool	True, False	True
Pool ratio	float	0.25, 0.50, 0.75	0.25
Pool heads	int	1, 2, 4	1
Pool sequence	list[func]	(see PyG Docs)	[GMPool_I]
Pool normalization	bool	True, False	False

Table S13: Adsorption energy of BM-dataset samples obtained with the optimized GNN and compared to DFT. Values are reported in eV.

Molecule	Family	Metal	E_{ads}^{DFT}	E_{ads}^{GNN}	Absolute error
mol1	Polyurethanes	Ag	-0.75	0.12	0.87
mol2	Polyurethanes	Ag	-0.81	-0.77	0.04
mol3	Polyurethanes	Ag	-0.77	-0.40	0.38
mol4	Polyurethanes	Ag	-0.84	-1.06	0.22
mol5	Polyurethanes	Ag	-1.11	-1.20	0.09
mol1	Polyurethanes	Au	-0.50	-0.12	0.38
mol2	Polyurethanes	Au	-0.40	-0.84	0.44
mol3	Polyurethanes	Au	-0.36	-0.19	0.17
mol4	Polyurethanes	Au	-0.39	0.00	0.39
mol5	Polyurethanes	Au	-0.84	-0.60	0.24
PE	Plastics	Pt	-1.17	-0.72	0.45
PP-it	Plastics	Pt	-0.93	0.00	0.93
PP-st	Plastics	Pt	-0.86	0.00	0.86
PS	Plastics	Pt	-2.95	-3.41	0.45
PET	Plastics	Pt	-2.86	-1.89	0.97
PE	Plastics	Ru	-0.18	-0.04	0.15
PP-it	Plastics	Ru	0.02	0.00	0.02
PP-st	Plastics	Ru	0.08	0.00	0.08
PS	Plastics	Ru	-3.10	-3.47	0.36
PET	Plastics	Ru	-4.78	-4.16	0.63
mol1	Biomass	Ni	-2.86	-2.51	0.36
mol2	Biomass	Ni	-2.39	-2.18	0.21
mol3	Biomass	Ni	-2.79	-1.64	1.15
mol4	Biomass	Ni	-1.92	-1.34	0.58
mol5	Biomass	Ni	-1.60	-1.51	0.09
mol1	Biomass	Ru	-3.65	-3.58	0.07
mol2	Biomass	Ru	-3.15	-2.95	0.20
mol3	Biomass	Ru	-3.41	-2.41	1.00
mol4	Biomass	Ru	-2.50	-2.16	0.33
mol5	Biomass	Ru	-1.99	-1.44	0.55

References

- (1) García-Muelas, R.; López, N. *J. Chem. Phys. C* **2014**, *118*, 17531–17537, DOI: [10.1021/jp502819s](https://doi.org/10.1021/jp502819s).
- (2) García-Muelas, R.; Li, Q.; López, N. *ACS Catalysis* **2015**, *5*, 1027–1036, DOI: [10.1021/cs501698w](https://doi.org/10.1021/cs501698w).
- (3) Weininger, D. *Journal of Chemical Information and Modeling* **1988**, *28*, 31–36, DOI: [10.1021/ci00057a005](https://doi.org/10.1021/ci00057a005).
- (4) O’Boyle, N. M.; Morley, C.; Hutchison, G. R. *Chemistry Central Journal* **2008**, *2*, 5, DOI: [10.1186/1752-153x-2-5](https://doi.org/10.1186/1752-153x-2-5).
- (5) Halgren, T. A. *Journal of Computational Chemistry* **1996**, *17*, 490–519, DOI: [https://doi.org/10.1002/\(SICI\)1096-987X\(199604\)17:5/6<490::AID-JCC1>3.0.CO;2-P](https://doi.org/10.1002/(SICI)1096-987X(199604)17:5/6<490::AID-JCC1>3.0.CO;2-P).
- (6) Ong, S. P.; Richards, W. D.; Jain, A.; Hautier, G.; Kocher, M.; Cholia, S.; Gunter, D.; Chevrier, V. L.; Persson, K. A.; Ceder, G. *Computational Materials Science* **2013**, *68*, 314–319, DOI: [10.1016/j.commatsci.2012.10.028](https://doi.org/10.1016/j.commatsci.2012.10.028).
- (7) Morandi, S. DFT data repository, DOI: [10.19061/iochem-bd-1-257](https://doi.org/10.19061/iochem-bd-1-257).
- (8) Álvarez-Moreno, M.; de Graaf, C.; López, N.; Maseras, F.; Poblet, J. M.; Bo, C. *Journal of Chemical Information and Modeling* **2014**, *55*, 95–103, DOI: [10.1021/ci500593j](https://doi.org/10.1021/ci500593j).
- (9) Isayev, O.; Oses, C.; Toher, C.; Gossett, E.; Curtarolo, S.; Tropsha, A. *Nat. Commun.* **2017**, *8*, 15679, DOI: [10.1038/ncomms15679](https://doi.org/10.1038/ncomms15679).

- (10) Cordero, B.; Gómez, V.; Platero-Prats, A. E.; Revés, M.; Echeverría, J.; Cremades, E.; Barragán, F.; Alvarez, S. *Dalton Transactions* **2008**, *nil*, 2832, DOI: [10.1039/b801115j](https://doi.org/10.1039/b801115j).
- (11) Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V., et al. *Journal of machine learning research* **2011**, *12*, 2825–2830.
- (12) Hamilton, W. L.; Ying, R.; Leskovec, J. Inductive Representation Learning on Large Graphs, 2017, DOI: [10.48550/ARXIV.1706.02216](https://doi.org/10.48550/ARXIV.1706.02216).
- (13) Baek, J.; Kang, M.; Hwang, S. J. Accurate Learning of Graph Representations with Graph Multiset Pooling, 2021, DOI: [10.48550/ARXIV.2102.11533](https://doi.org/10.48550/ARXIV.2102.11533).
- (14) Agarap, A. F. *CoRR* **2018**, DOI: <https://doi.org/10.48550/arXiv.1803.08375>.
- (15) Kingma, D. P.; Ba, J. *CoRR* **2014**, DOI: <https://doi.org/10.48550/arXiv.1412.6980>.
- (16) Li, Q.; López, N. *ACS Catalysis* **2018**, *8*, 4230–4240, DOI: [10.1021/acscatal.8b00067](https://doi.org/10.1021/acscatal.8b00067).
- (17) Puértolas, B.; Rellán-Piñeiro, M.; Núñez-Rico, J. L.; Amrute, A. P.; Vidal-Ferran, A.; López, N.; Pérez-Ramírez, J.; Wershofen, S. *ACS Catalysis* **2019**, *9*, 7708–7720, DOI: [10.1021/acscatal.9b02086](https://doi.org/10.1021/acscatal.9b02086).
- (18) Ding, S.; Hülsey, M. J.; Pérez-Ramírez, J.; Yan, N. *Joule* **2019**, *3*, 2897–2929, DOI: [10.1016/j.joule.2019.09.015](https://doi.org/10.1016/j.joule.2019.09.015).

- (19) Li, L.; Jamieson, K.; Rostamizadeh, A.; Gonina, E.; Hardt, M.; Recht, B.; Talwalkar, A. **2018**, DOI: [10.48550/ARXIV.1810.05934](https://doi.org/10.48550/ARXIV.1810.05934).
- (20) Liaw, R.; Liang, E.; Nishihara, R.; Moritz, P.; Gonzalez, J. E.; Stoica, I. *arXiv preprint arXiv:1807.05118* **2018**.