

Bachelor Thesis

Uncertainty in Graph Neural Networks, applied within Quantum Chemistry

Name: **Student number:**

Lars Lohmann s173875

Supervisor: **Title:**

Mikkel Nørgaard Schmidt Associate Professor

Department:

Department of Applied Mathematics and Computer Science, Technical University of Denmark



1 Abstract

Contents

1 Abstract	2
2 Acknowledgements	4
Preface	5
3 Introduction	6
4 Definition of Concepts	8
4.1 Model Terminology	8
5 Theory	9
5.1 Message Passing Neural Networks	9
5.2 Ensembles	10
6 Methods	12
6.1 Experimental Setup	12
6.2 Dataset	12
6.3 Model	13
6.4 Training- and Validation-loop	17
6.5 Test Loop	18
6.6 Modelling Hyperparameters	18
6.7 Experimental Hyperparameters	19
6.8 Software tools and Hardware	19
6.9 Reproduceability	20
7 Data	21
7.1 Cross Coupling Catalysts	21
7.2 Data Privacy- and Quality-issues	26
8 Results	28
9 Discussion	29
10 Conclusion	29
References	29
11 Appendix	31
11.1 Appendix A	31
11.2 Appendix B	31
11.3 Appendix C	31

2 Acknowledgements

The completion of this project would not have been possible without the dedicated supervisor efforts from Mikkel Nørgaard Schmidt.

Mikkel Nørgaard Schmidts help has manifested through all of the stages of the project, from planning, to scoping, to executing and finalizing. His inviting presence allows for curiosity to an extend very few prior of prior supervisors have managed. His own curiosity towards the topic is contagious. His direct and fair questions and advice, has contributed to the bettering of the project, and a deepening of the knowledge around the subject matter.

When things get tough, his calming manner gave the mental surplus that was needed to finalize the work, and deliver on this project.

I would recommend his counselling efforts to any student at DTU, whether it being in a project or teaching capacity.

Thank you!

Lars Lohmann

Approval

This thesis has been prepared over 12 weeks with the Department of Applied Mathematics and Computer Science, Technical University of Denmark, DTU, in partial fulfilment for the degree Bachelor of Science in Engineering.

Lars Lohmann - s173875 _____

3 Introduction

The implications of predicting molecular properties and dynamics have both large industrial- and research implications within the natural- and life-sciences. The large amounts of compute, necessary for simulation, often makes desirable metrics intractable to estimate. The efforts of this study, will help provide tools for industry to reduce lead time and risk in development fields such as drug discovery and material sciences for energy storage[1].

These interests, have for many years been supported by research and development methods within the traditional field of chemistry, with emphasis on estimation techniques, that carry a heavy computational load like traditional molecular dynamics (MD)[2] or Density Functional Theory[1]. Therefore, despite the emergence of cloud-data-centers and power full algorithms available at lower cost to research, some problems still exist which are intractable with known methods.

That fact, has prompted the emergence of the computational fields within Artificial Intelligence (AI), such as Machine Learning (ML), Deep Learning (DL) and Geometric Deep Learning (GDL), to enrich the solution space. These methods produce models which tries to decrease the lead time for predictions, while maintaining an acceptable error-rate[3]. Specifically the field of Geometric Deep Learning, has emerged with a viable set of data-driven methods, for predicting molecular properties and dynamics[4]. Geometric Deep Learning comprises a set of methods, which tries to enrich the initial data structures with a geometric prior through the structuring of data in i.e graphs, which is comprised of nodes and edges between nodes[4]. Examples of encoded information could be distances, angles or directions between nodes, as well as symmetry properties namely equivariance or invariance.

Since significant risk is embodied in the task of predicting chemical and biochemical dynamics, due to end-uses often being utilized in pharmaceutical products, proper mechanisms for estimating the the uncertainty around prediction methods, and mitigating the impact of high uncertainty through calibration, are needed. Traditionally, data-driven methods within machine learning and deep learning, have problems with providing the necessary uncertainties, on the metrics being predicted, making calibration for more robust overall performance hard[5][1][6]. Providing methods for estimating uncertainty, and calibrating, would allow for defaulting to more traditional methods for estimating chemical properties, when uncertainty is high in the data-driven models[1].

In support of producing such methods, uncertainty, can effectively be distinguished between two concepts, namely epistemic- and aleatoric- uncertainty[1][6]. Aleatoric, also known as statistical uncertainty, refers to the uncertainty inherent in experimental outcomes, due to random effects. This class of uncertainty is classified as irreducible, due to no amount of additional information in the modelling effort, being able to reduce this type of uncertainty[6]. Epistemic, also called

systemic uncertainty, refers to the uncertainty in a model, produced by lack of information, and can therefore be reduced. These two elements of uncertainty are considered by researchers to be the sole components of total uncertainty, their sum being equal to total uncertainty in an experiment[6].

This paper is concerned with the task of molecular estimation of a single property, namely the binding energy in a catalyst process between a transition metal, and a ligand[7]. The objective is to compare two ensembles of machine learning models that try to predict the binding energy, with similar architecture, but different hyperparameters in their internal representations of molecules, namely the state representation of scalar- and vector-properties of graphs. The main contribution of this paper, will be to assess the influence of this hyperparameter on modelling performance and uncertainty, through the comparison of the ensembles. This modelling performance will be measured in both a traditional metrics namely the Mean Squared Error between the regression metric and the target, but also the uncertainty in the ensembles, produced by changing the hyperparameters to a different value.

Specifically, this paper will utilize a Message Passing Neural Network (MPNN) architecture inspired by the PAINN architecture[8], to ingest graph representations of molecules, with Cartesian coordinates, atom-types and lengths of the bonds in molecules, in order to predict the binding energy in kcal per mol. This model structure, will be provided with a different set of internal state complexity, and placed in two separate uniformly weighted mixture models, one with higher sophistication and one with lower. The MPNN models will provide individual estimates of the target value in the test-set, where the mean and variance of the ensembles will be the bedrock for ensemble predictions and uncertainty estimations, inspired by[9], and [10].

Specifically, this paper asks:

Does the size of internal state representations, impact prediction error and uncertainty in ensembles of Message Passing Neural Networks

The following sections of this paper, namely[3,4,5], we will go through relevant concepts and theories, necessary for assessing the results and methods in this paper. This is followed by a description of the experimental methods, and a presentation of the data set, on which the methods have been utilized in sections[6,7]. Lastly, results followed by a conclusion and a discussion of the potential points of error, which influence the outcomes of the experiment in sections[8,9,10].

4 Definition of Concepts

4.1 Model Terminology

- **Sophisticated**

Sophisticated utilized in the scope of this paper, is used as a summary term to generally describe a models number of internal parameters. A model being more sophisticated, will have more parameters in its internal structure, than a less sophisticated model.

5 Theory

In this section, we will discuss the theoretical foundations of message passing neural networks (MPNNs) and ensembles for predicting the binding energy of molecules in a catalyst process.

5.1 Message Passing Neural Networks

The key idea behind MPNNs is the concept of message passing, where each node in the graph exchanges information with its neighbors in an iterative manner.

Let $G = (N, E)$ be a graph, where N is the set of nodes and E is the set of edges. Each node $n \in N$ is associated with an invariant representation vector $\mathbf{S}_n \in \mathbb{R}^{F \times 1}$, an equivariant representation vector $\vec{\mathbf{v}}_n \in \mathbb{R}^{F \times 3}$ and a position in 3D space $\vec{r}_i \in \mathbb{R}^3$. F denotes the number of features in the state representation. Each edge $(i, j) \in E$ is associated with a relative position $\vec{r}_{ij} = \vec{r}_i - \vec{r}_j$. A node $n_j \in N$ is said to be a neighbour of a node $n_i \in N \setminus \{n_j\} = \mathcal{N}(i)$ if n_j is within the cutoff distance of n_i : $\mathcal{N}(i) = \{j | \|\vec{r}_{ij}\| \leq r_{cut}\}$.

The message passing process in an MPNN can be described by the following equations 1 and 2 following the notation of [8].

$$\vec{\mathbf{m}}_i^{v,t+1} = \sum_{j \in \mathcal{N}(i)} \vec{M}_t(\mathbf{S}_i^t, \mathbf{S}_j^t, \vec{\mathbf{v}}_i^t, \vec{\mathbf{v}}_j^t, \vec{r}_{ij}) \quad (1)$$

$$\vec{\mathbf{s}}_i^{v,t+1} = \vec{U}_t(\mathbf{S}_i^t, \vec{\mathbf{v}}_i^t, \vec{\mathbf{m}}_i^{v,t+1}) \quad (2)$$

The U and M functions respectively being the update, and message functions, taking into account scalar, vector and position features. These equations utilize both invariant and equivariant representations, doing operations and letting representations interact, in accordance with knowledge surrounding the chemical descriptor trying to be predicted. The underlying operations of the functions U and M , have a linearity constraint on equivariant representations, in order to maintain information, throughout the modelling process[4]. These operations need to be chosen, such that they fulfill either of the two equations (3 for invariant representations or 4 for equivariant representations) as shown below for any rotation matrix $R \in \mathbb{R}^{3 \times 3}$:

$$\mathbf{f}(\vec{\mathbf{x}}) = \mathbf{f}(R\vec{\mathbf{x}}) \quad (3)$$

$$R\vec{\mathbf{f}}(\vec{\mathbf{x}}) = \mathbf{f}(R\vec{\mathbf{x}}) \quad (4)$$

Equivariant representations of directions between nodes, have shown to be more expressive, at a similar computational cost. This effectively means a higher level of information can be maintained through the modelling phase, if the equivariance is maintained. This information is centered around directional information, which might be useful for predicting the chemical descriptor in question. Directional information is also expanded from the relative positions between edges \vec{r}_{ij} , via the radial basis function (RBF)[4]. The radial basis function outputs are

chosen such that it is centered around twenty points, $C = \{1, 2, 3, \dots, 20\}$. The function takes directional information $d = \|\vec{r}_{ij}\|$ and a cutoff r_{cut} , and produces expansions following the formula below.

$$\mathbf{RBF}(d) = \sum_{c \in C} \sin\left(\frac{c\pi}{r_{cut}} * d\right) / d \quad (5)$$

The expansions provides a representation of similarity or dissimilarity, between the input relative position \vec{r}_{ij} , and the chosen points C [11]. This information is further expanded by a neural network layer with a set of weights biases, seen in equation 6, before arriving at the cosine cutoff.

$$y_{rbf} = \mathbf{w} \cdot \mathbf{RBF}(d) + \mathbf{b} \quad (6)$$

The cosine-cutoff function is originally derived from [2] in the form of the equation below7.

$$\mathbf{f}_{cut}(x) = \begin{cases} 0.5 \cdot \cos\left(\frac{\pi d}{r_{cut}} + 1\right) & \text{if } d \leq r_{cut} \\ 0 & \text{if } d > r_{cut} \end{cases} \quad (7)$$

The cutoff function implemented in this project, utilized the equation below8:

$$\mathbf{f}_{cut}(x) = \begin{cases} 0.5 \cdot \cos\left(\frac{\pi d}{r_{cut}} + 1\right) \cdot y_{rbf} & \text{if } d \leq r_{cut} \\ 0 & \text{if } d > r_{cut} \end{cases} \quad (8)$$

This alteration allows for the flow of information from both the radial basis expansions and the initial vector differences, which from the process diagram in[8], is not happening, although this is believed by the project to be the intention. These three highlighted steps, the rbf, the set of weights and biases and the cosine cutoff, constitute the rotationally invariant filter \mathcal{W}_s , proposed by [3], and [2], and have been highlighted in this section, due to the difference in approaches from these propositions.

5.2 Ensembles

Ensembles is a technique for combining the predictive performance of individual modelling approaches and to as an extension measuring total uncertainty in the model. In the context of MPNNs, ensembles can be created by training multiple models with different initializations, and combining their predictions.

Let \mathbf{y}_i be the predicted binding energy for molecule i by model M_k , where i represents the molecule index and k represents the model index. The ensemble prediction $\hat{\mathbf{y}}_i$ can be calculated as the mean of the individual model predictions, as done in[10]:

$$\hat{\mathbf{y}}_i = \frac{1}{T} \sum_{k=1}^T \mathbf{y}_i^{(j)} = \hat{\mu}_i \quad (9)$$

where T is the number of models in the ensemble. This format of calculation weights the inputs from each model equally, similarly to if it was a uniformly weighted mixture model as in [1].

To estimate the uncertainty in the ensemble predictions, we can calculate the variance of the individual model predictions[10]:

$$\text{Var}(\hat{y}_i) = \frac{1}{T} \sum_{k=1}^T (y_i^{(k)} - \hat{y}_i)^2 = \hat{\sigma}_i^2 \quad (10)$$

Finally as the foundation of estimating uncertainty in the ensemble predictions, we can estimate the expected normal calibration error. This method relies on sorting predictions of targets, in K equal sized bins, computing the predicted root mean variance¹¹.

$$RMV_k = ? \quad (11)$$

Thereafter producing the empirical root mean squared error¹²:

$$RMSE_k = \sqrt{\frac{1}{k} \sum_{i=1}^K (y_i^{(k)} - \hat{y}_i)^2} \quad (12)$$

and finally utilizing the expression, following equation the equation in[1]:

$$ENCE = \frac{1}{k} \sum_{k=1}^K \frac{|RMV_k - RMSE_k|}{RMV_k} \quad (13)$$

The approach of [1] is based on models in an ensemble, outputting both a mean value and a variance metric, effectively predicting both the target, and predicting its own uncertainty. This is not the case in this project, since the model solely predicts the target.

6 Methods

This section will detail the implementation and necessary technological foundation, for applying the theories described, in order to produce the experimental results for answering the research question. The section will also detail the necessary information, for reproducing and understanding the experimental results and their becoming.

6.1 Experimental Setup

The experimental setup are comprised of the following components:

- Dataset
- Model
- Training- and Validation-Loop
- Test loop

These four components all together make up the main functionality of the experiment, and will be detailed individually below. The detailing will have its focus on how concepts and theories previously explained, are formatted and shaped, in order to be viable for practical modelling purposes.

6.2 Dataset

The data set is comprised of one row for each node for 695 molecules, with a the following fields for each node:

- Atom-type
- Overall binding energy (target property)
- x-coordinate
- y-coordinate
- z-coordinate
- Graph id

The atom type signifies the type of atom for a given node by name. The overall binding energy is the target property, for the full catalyst process. This means all nodes with identical graph id will have a identical overall binding energy. The x, y, and z coordinates are the coordinates of the nodes in the molecular graph, measured in Ångström.

These fields are then translated to the relevant datastructures, via the graph data set constructor. This constructor defines the following data structures, for utilization in the modelling task. Important to say is that the data set constructor, constructs a data set of potentially several graphs, and not only one. The data set consists of the following critical fields:

- num_nodes: Number of nodes in the data set of graphs
- node_from: The node from which the edge originates
- node_to: The node to which the edge points
- num_graphs: Number of graphs in the data set
- node_graph_index: The index of the graph the edge resides in
- unique_atoms: The number of unique atoms in the data set
- edge_lengths: The lengths of the edges
- edge_vector_diffs: The vector differences of the edges

The number of nodes in dataset, allows for the model to create initial state representation structures with proper dimensions for the modelling task. The node_from property, allows for summation of all neighbouring node state representations, within the cutoff limit, into the designated node, we are trying to model in the message block. The node_to property, allows for initialising the state representation of the node in the update block, namely the scalar properties and the vector properties. The scalar properties are concatenated tensors of embedded atom representations, linked by the index in the node_to property. For scalar properties it is zero-vectors, instead of embedded tensors. The num_graphs property is used to initialize the initial state representation of the graphs. The node_graph_index property is an index linking an edge, to a corresponding graph. When summing over all neighbour nodes occur, the full message state representation of all neighbours under the cutoff, is summed over the node_graph_index property, to produce a state representation of graphs. The unique_atoms property, allows for the embedding function to create a unique atom representation for each type of atom in the data set. The edge_lengths allow for a masking of the edges and thereby nodes that are not within the cutoff limit. The edge_vector_diffs are the vector differences of the edges, being an input variable in the message block.

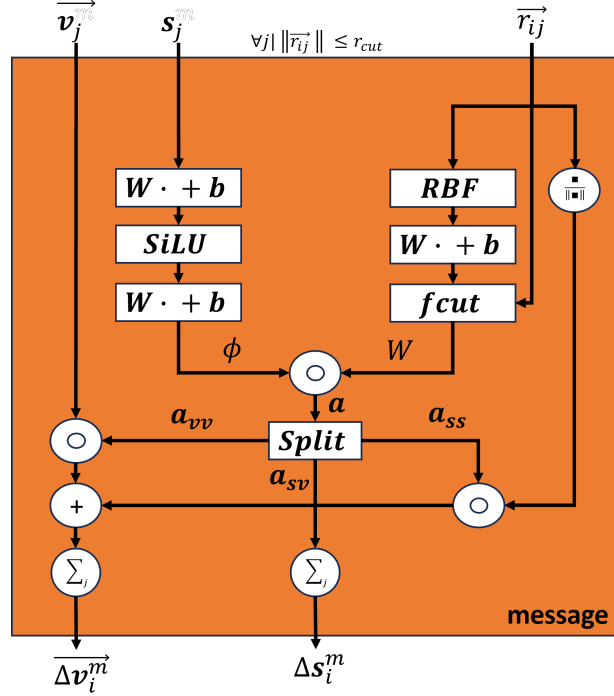
6.3 Model

The model architecture consists of two main components, namely the Message block, and the Update block as formerly introduced in the the theory section 1, and 2. The Message block is conceptually responsible for the message passing between the nodes in the graph, allowing for the flow of information between the nodes. This information is then summed up and stored in state variables in the form of a scalar and vector representation. This summation is done in a residual manner, following the implementation of the PAINN architecture[8]. An equation defining the residual calculation of the scalar-properties¹⁴ and vector-properties¹⁵ can be found below, as well as an illustration of the message block¹.

$$\Delta \mathbf{s}_i^m = (\phi_s(\mathbf{s}) * \mathcal{W}_s)_i = \sum_j \phi_s(\mathbf{s}) \circ \mathcal{W}_s (\|r_{ij}\|) \quad (14)$$

$$\Delta \vec{v}_i^m = \sum_j \vec{v}_j \circ \phi_{vv}(s_j) \circ \mathcal{W}_{vv}(\|r_{ij}\|) + \sum_j \phi_{vs}(s_j) \circ \mathcal{W}'_{vs}(\|r_{ij}\|) * \frac{r_{ij}}{\|r_{ij}\|} \quad (15)$$

Figure 1: Message Block



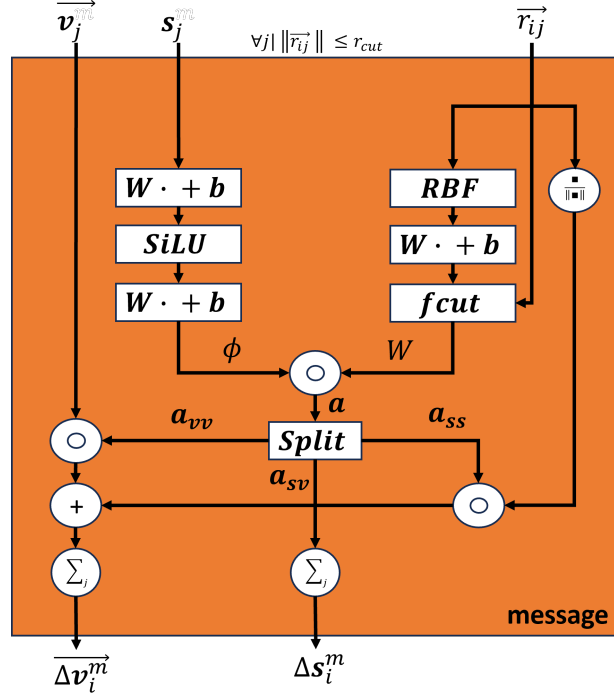
Worth noting about the figure, compared to the similar one in [8], is the use of the r_{ij} vector, that also flow into the 'fcut' function, as mentioned in 5.1. The model in this paper, does not highlight the dimension expansions or collisions from network element to network element, due to the model parameters varying in size, as it is the scope of this project to investigate similarities or dissimilarities in performance, with respect to the PAINN architecture and dissimilar state representation sizes. both sizes ϕ and \mathcal{W} are chosen sub-architectures based on the PAINN architecture from [8]. The image also contains hints of the \mathbf{a} -tensor, which is a product of the ϕ and \mathcal{W} tensors, and not highlighted in [8].

The residual values produced in the message block, then used in the update block, in order to update the node states, based on the summation over neighbouring states. This structure is also heavily inspired by the PAINN architecture [8]. The below equations represent the residual updates on both the scalar- 16 and the vector- 17 properties. Below is also a figure representing the update block2.

$$\Delta s_i^u = \mathbf{a}_{ss}(s_i, \|\mathbf{V}\vec{v}_i\|) + \mathbf{a}_{sv}(s_i, \|\mathbf{V}\vec{v}_i\|) \langle \mathbf{U}\vec{v}_i, \mathbf{V}\vec{v}_i \rangle \quad (16)$$

$$\Delta \mathbf{v}_i^u = \mathbf{a}_{vv}(s_i, \|\mathbf{V}\vec{v}_i\|) \mathbf{U}\vec{v}_i \quad (17)$$

Figure 2: Update block



Worth noting about the above figure 2, is that the indexes of the scalar and vector representations, have been altered, compared to [8]. This comes down to the project evaluating the indexing of the scalar and vector representations, in the paper, were faulty, and did not represent the PAINN architecture as described in prior sections of the paper.

The blocks are stacked in coupled sequences of five full rounds, like the illustration below highlights3, a corresponding model can be found in [8]. The illustration shows an architecture comprised of three sets of message- and update-blocks, the implementation in this project utilized five sets. The model representation property \vec{v}_i^0 is initialized to the zero vector, and the scalar property s_i^0 is initialized to a random embedding on atom level, for each node in the graph.

The diagram illustrates the proposed Gated Graph Neural Network (GGNN) architecture. The input is a node feature vector Z_i . It passes through an embedding block to produce v_i^0 and s_i^0 . The message blocks (orange) and update blocks (yellow) are stacked. The message blocks take v_i^0 and s_i^0 as inputs and produce v_i^1 and s_i^1 . The update blocks take v_i^1 and s_i^1 as inputs and produce v_i^2 and s_i^2 . The final output is E , which is the sum of the output of the last message block and the output of the last update block.

6.4 Training- and Validation-loop

Algorithm 1 MPNN Training Loop

```
1: Initialize size of ensemble:  $T$ 
2: Initialize size of batch:  $B$ 
3: Initialize data sets:  $Td, Vd$ 
4: initialize training data loader:  $TD = \text{Dataloader}(B, Td)$ 
5: initialize validation data loader:  $VD = \text{Dataloader}(B, Vd)$ 
6: Initialize number of epochs:  $E$ 
7: Initialize validation index:  $V$ 
8: for  $model = 1, 2, \dots, T$  do
9:   Initialize MPNN model instance:  $model = \text{PAINN}()$ 
10:  initialize MSE loss function:  $loss = \text{MSE}()$ 
11:  Initialize Adam optimizer:  $optimizer = \text{Adam}(model.parameters())$ 
12:  Initialize learning rate scheduler:  $scheduler = \text{StepLR}(optimizer)$ 
13:  for  $epoch = 1, 2, \dots, E$  do
14:    for  $batch = 1, 2, \dots, TD$  do
15:      Normalize targets:  $y = \text{normalize}(y)$ 
16:      Predict binding energy:  $\hat{y} = model(batch)$ 
17:      Calculate loss:  $l = loss(y, \hat{y})$ 
18:      Backward pass:  $l.backward()$ 
19:      Optimizer step:  $optimizer.step()$ 
20:      if  $epoch \% V == 0$  then
21:        Model in evaluation mode:  $model.eval()$ 
22:        for  $batch = 1, 2, \dots, VD$  do
23:          Normalize targets:  $y = \text{normalize}(y)$ 
24:          Predict binding energy:  $\hat{y} = model(batch)$ 
25:          Calculate loss:  $l = loss(y, \hat{y})$ 
26:          Apply exponential smoothing:  $l = l * 0.9 + l_{-1} * 0.1$ 
27:          Scheduler step:  $scheduler.step(l)$ 
28:        end for
29:      end if
30:    end for
31:  end for
32: end for
```

6.5 Test Loop

Algorithm 2 MPNN Testing Loop

```
Initialize size of ensemble:  $T$ 
2: Initialize size of batch:  $B$ 
   Initialize data sets:  $Testd$ ,
4: initialize training data loader:  $TestD = Dataloader(B, Testd)$ 
   for  $model = 1, 2, \dots, T$  do
6:   Initialize target list:  $y_l$ 
     Initialize prediction list:  $\hat{y}_l$ 
8:   for  $batch = 1, 2, \dots, TestD$  do
       Normalize targets:  $y = normalize(y)$ 
10:   Predict binding energy:  $\hat{y} = model(batch)$ 
       Calculate loss:  $l = loss(y, \hat{y})$ 
12:   Append target to list:  $y_l$ 
       Append prediction to list:  $\hat{y}_l$ 
14:   end for
       Save predictions and targets to file
16:   Save model to file
   end for
```

6.6 Modelling Hyperparameters

Generally, two subsets of hyperparameters exists, those relevant to the less sophisticated set of ensemble models (titled MPNN64), and those relevant to the more sophisticated ensemble models (titled MPNN128). The table below¹, details shared and non-shared model hyperparameters, as well as highlighting the parameters, which are consistent with the PAINN architecture described in [8] with a star (*).

Table 1: Modelling hyperparameters

Model Hyperparameters	MPNN64	MPNN128	Unit
Shared			
Number of physical dimensions*	3		Dimensions
Number of Message Passing Rounds*	5		Rounds
Patience*	5		Validation steps
Weight Decay*	0.01		N/A
Learning Rate Decay*	0.5		N/A
Exponential Smoothing for Validation*	0.9		N/A
Radial Basis Parameter*	20		N/A
Learning Rate	0.001		N/A
Cutoff Distance	4		Ångström
Not shared			
State dimension size	64	128	Dimensions

An explanation of the individual hyperparameters influence on the model, and a justification for selection of the above values, can be found in the appendix section 11.1.

6.7 Experimental Hyperparameters

The experimental hyperparameters defining the implementation of the above model setting, can be found in table below 2.

Table 2: Experimental Hyperparameters

Experimental Hyperparameters	Value	Unit
Number of molecules in Data Set	695	Molecules
Training Split	0.8	N/A
Validation Split	0.1	N/A
Test Split	0.1	N/A
Ensemble Size	5	Models
Initial number of Epochs	200	Epochs
Validation Index	5	Epochs
Model Saving Interval	0.01	N/A
Batch Size	5	Molecules

An in explanation of the selection of experiment hyperparameters and justification, be found in the appendix section 11.2.

6.8 Software tools and Hardware

The projects execution were developed in Python 3.10.9, with heavy emphasis on Pytorch, Numpy and Pandas, supported by few shell scripts for cloud job definitions. A full list of pack-

Table 4: Cloud Job Parameters

Cloud Job Parameters	<i>MPNN128</i>	<i>MPNN64</i>
Queue	gpub100	gpua100
Number of cores	8	8
Memory/Cores	5 gb	5 gb
Number of hosts	1	1

ages, and their versions utilized in the experiment, can be seen in appendix 11.3. The initial data-structures and pipeline were inspired by course-material made by Mikkel Nørgaard Schmidt. The hardware utilized for the final experimental setup, were two types of GPUs, belonging to the High Performance Computing (HPC) cluster at the Technical University of Denmark (DTU). The two queues utilized were 'gpua100' and 'gpub100'. For further details on the hardware, please refer to the following reference:[12]. Total computation time in the cloud was three days, in order to produce the trained ensemble of ten individual models.

6.9 Reproduceability

For reviewing and reproduceability purposes, the following github link stores the latest version of the code: <https://github.com/Lohmann94/head>, which was utilized in the experiment. The experiment can be reproduced with the following seed-values³, which can be found for individual models:

Table 3: Model seeds for reproduction

MPNN128 Ensemble:		MPNN64 Ensemble:	
<i>Model</i>	<i>Seed</i>	<i>Model</i>	<i>Seed</i>
MPNN128_1	49	MPNN64_1	90
MPNN128_2	10	MPNN64_2	52
MPNN128_3	100	MPNN64_3	7
MPNN128_4	20	MPNN64_4	9
MPNN128_5	30	MPNN64_5	87

The cloud-job configuration variables for reproducing with the described performance, can be found in the below table4:

7 Data

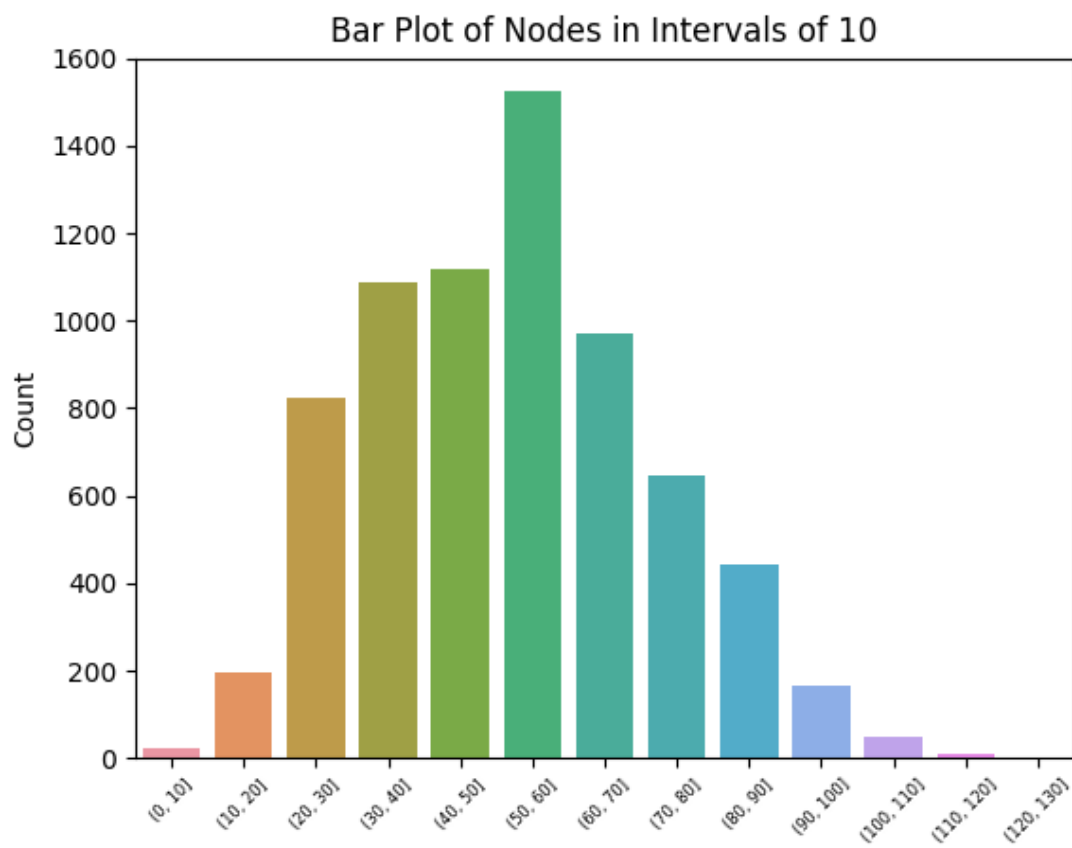
7.1 Cross Coupling Catalysts

The data set of interest consists of 7.054 optimized geometries of compounds being candidate structures in oxidative addition catalytic process. The optimized geometries will have a binding energy associated with it, that dictates the net energy produced by the oxidative addition of a specific transition metal[7], which is the target of our modelling task. The original papers set out to further select 557 catalyst candidates which are within a predefined thermodynamic window, and a further selection of candidates based on their pricepoint per mol.

The molecular compounds are represented as graph structures of nodes and edges. The nodes constitute atoms, and the edges molecular connections between atoms. The representation will further consist of Cartesian coordinates provided in the data set, fixing the individual nodes in three dimensional space. This also provides the data set with the possibility of deriving distances between the nodes, or the length of edges, which are crucial to the modelling efforts.

The number of number of nodes in the graphs, range from a minimum of 5, a mean of 52.51, to a maximum of 123 nodes. A bar plot of the number of graphs in various buckets of node-counts, can be seen below4.

Figure 4: Bar plot of number of nodes per graph



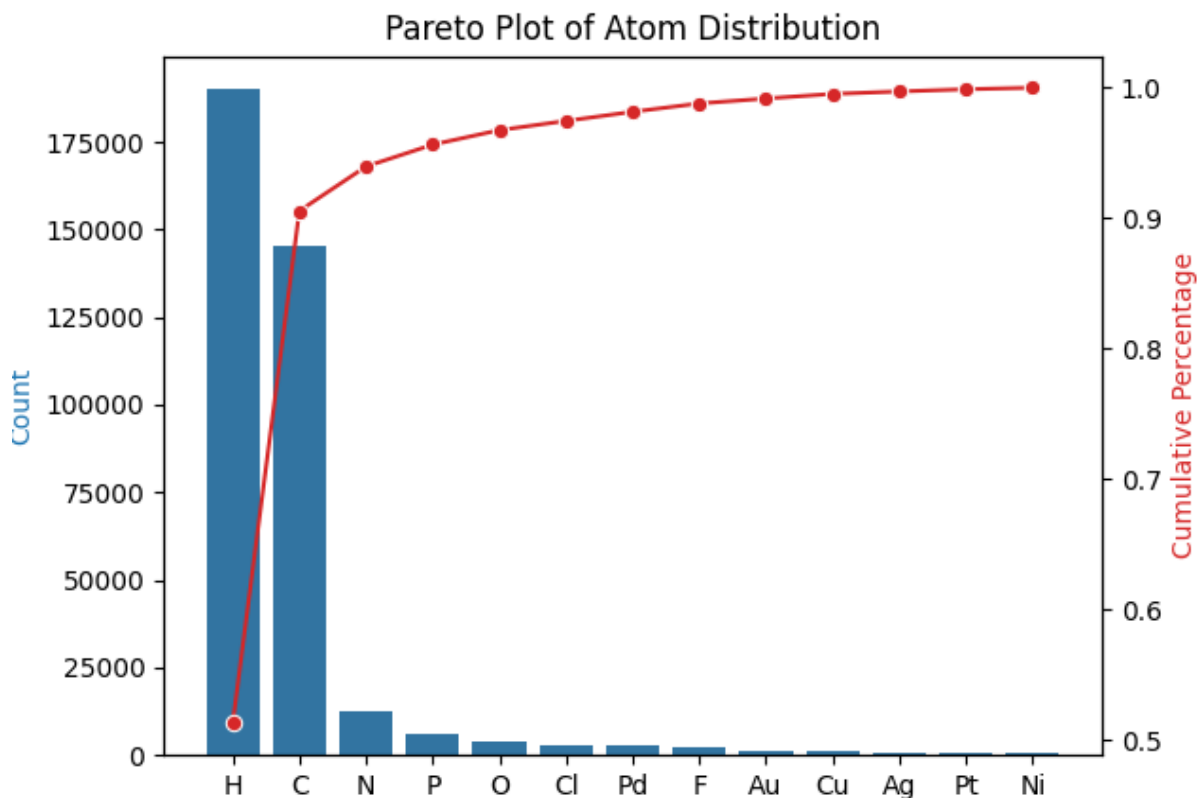
The atom distribution of the graphs a constituted by 13 unique atoms being:

- H: Hydrogen
- C: Carbon
- N: Nitrogen
- P: Phosphorus
- O: Oxygen
- Cl: Chlorine
- Pd: Palladium
- F: Fluorine
- Au: Gold
- Cu: Copper
- Ag: Silver
- Pt: Platinum

- Ni: Nickel

The distribution of the above atoms, can be seen in a Pareto graph below⁵. The distribution is mainly made up by hydrogen and carbon, those two being 90.6 percent of the atom-representation combined. The next highest contributor is nitrogen, only representing 3,4 percent of the atoms. The remaining atoms constitute 6 percent of the atom representation in the data set.

Figure 5: Pareto plot of atom distribution in data set

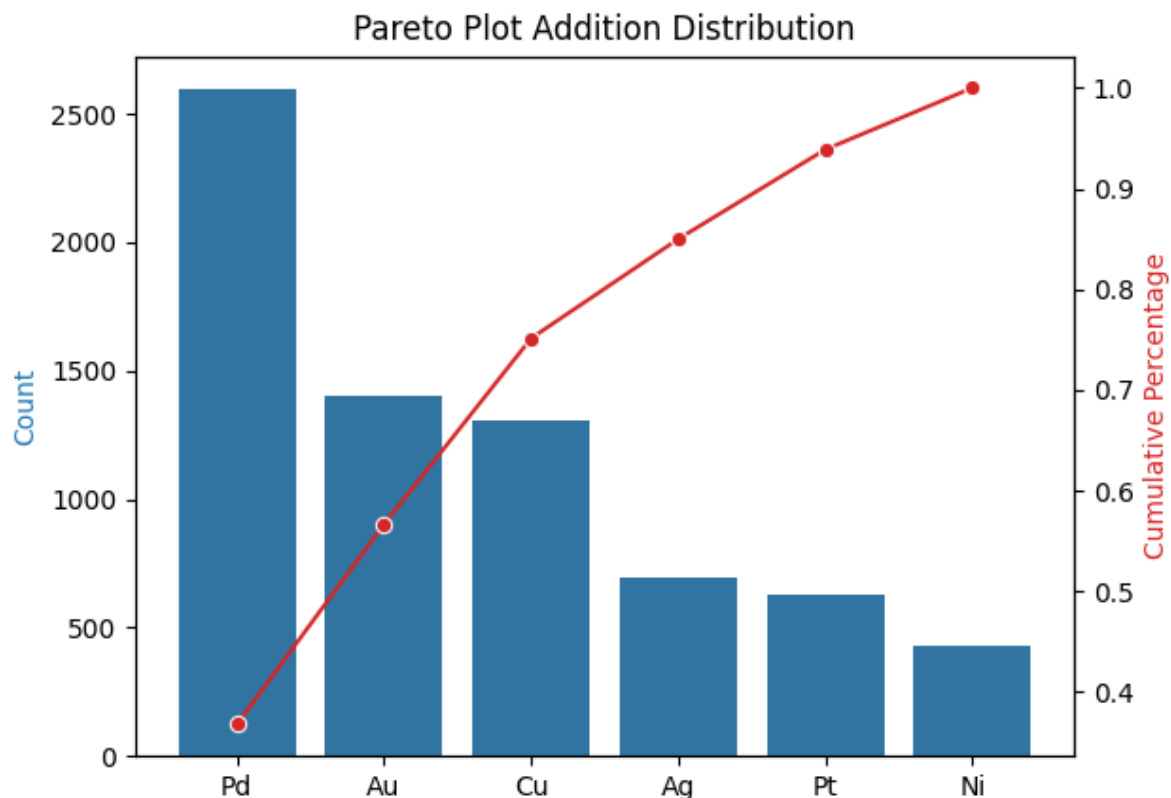


The additive part of the catalyst process, are selected transition metals, namely:

- Pd: Palladium
- Au: Gold
- Cu: Copper
- Ag: Silver
- Pt: Platinum
- Ni: Nickel

The data set consists of the following distribution of the above transition metals, below can be seen a pareto similar to further above, but now show the distribution of transition metals in the data set⁶.

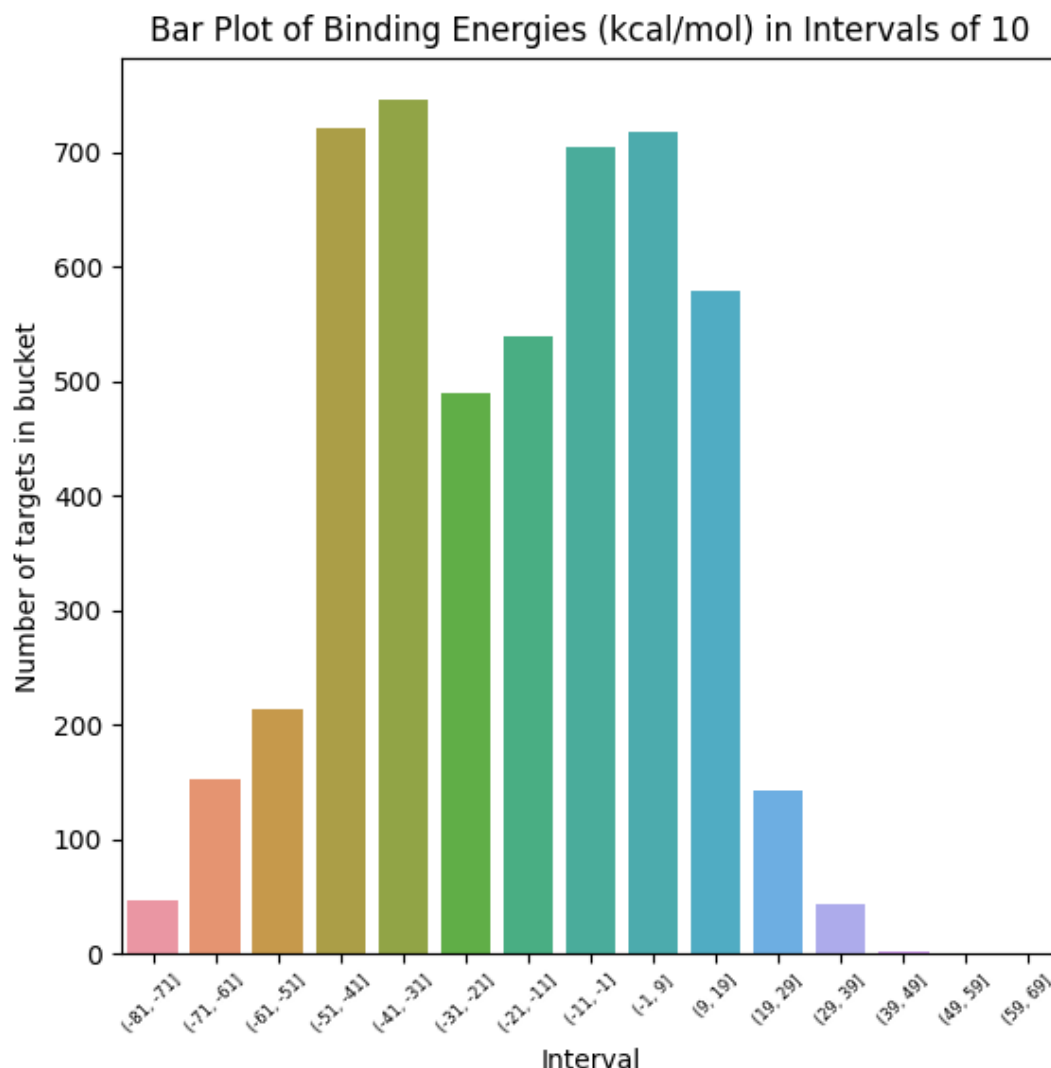
Figure 6: Pareto plot of transition metal distribution in data set



The distribution is mainly comprised of Palladium and Gold, standing for 57 percent of the transition metals in the data set.

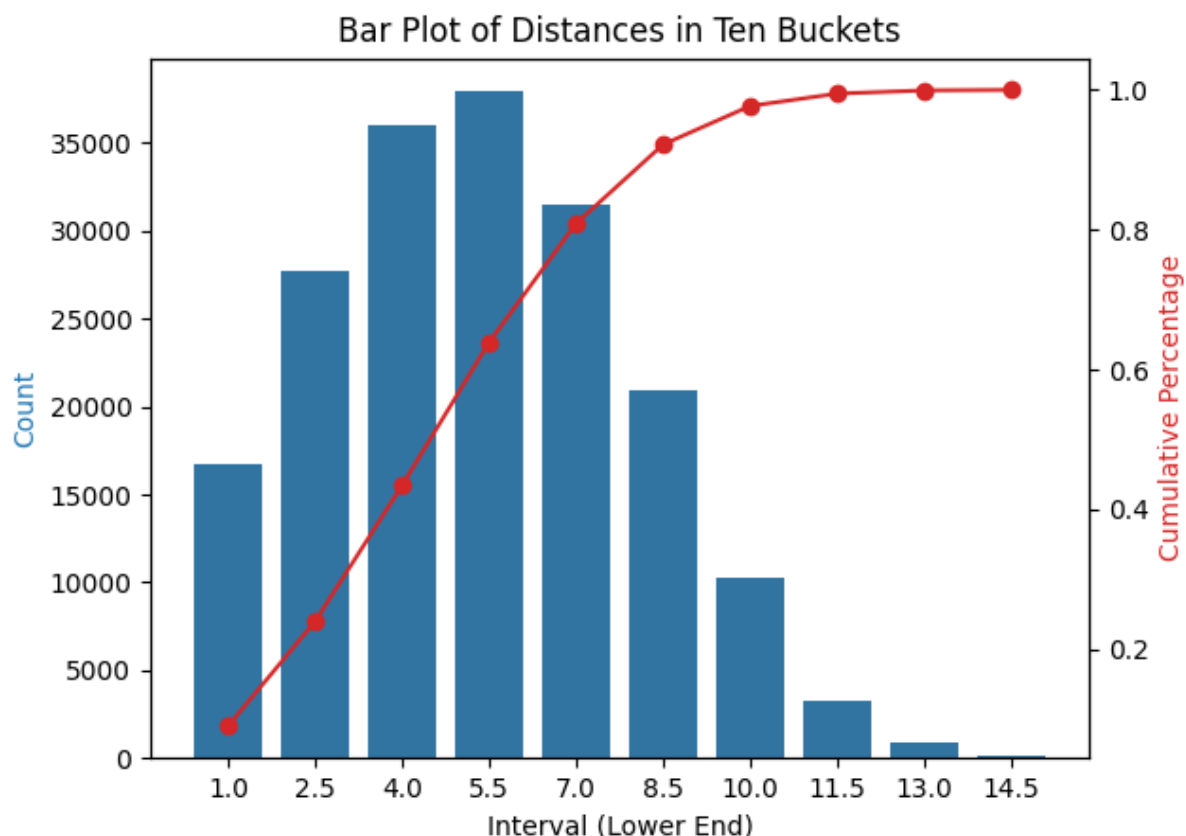
The catalyst process can be both exothermic and endothermic, represented in either a positive net binding energy as a target, or a negative net binding energy as a target for the regression task. The regression targets range from a minimum of -80.82 being endothermic, to 61.49 being exothermic, where the mean is placed at -18.62. A bar plot with buckets of ten, of the distribution of targets can be seen below⁷.

Figure 7: Bar plot of binding energy distribution in targets



The edge-wise distances between nodes, measured in ångstrom, are crucial to the modelling task, since a model-parameter called ‘r-cut’, defines a cutoff limit from which edges with a distance lower than the cutoff are defined as neighbours of a given node, and those above as non-neighbours. The neighbours goes are allowed a higher model influence, and therefore understanding where to set the cutoff, is a crucial hyperparameter. The distances of the full data set is intractable to compute locally, so a subset of 50 randomly chosen graphs and their edge-wise distances are plotted as a representation of the full data set⁸. As we can see in the plot, most distances are covered by the interval 0 to 7, at 64 percent, and only 20 percent of distances lie in the interval above 8.5.

Figure 8: Pareto plot of Distances between edges from 50 graphs



In a chemistry context, it is argued that a large portion in the variance in binding energy, can be explained by local interactions among atoms[8]. So even though a high degree of information lies beyond cutoff limits of i.e four or five, then information might not be crucial to the prediction in energy we are trying to make.

The modelling task will take in a given molecular structure, and produce a regression metric, trying to predict the associated binding energy from the oxidative addition in the unit kcal/mol. Further more, the regression efforts will be put in an ensemble context of a number of structurally similar models, producing a mean and a variance over the predicted regression metric from the ensemble.

7.2 Data Privacy- and Quality-issues

Due to the data sources being publicly available for download, under the creative commons attribution 4.0 license, the implementation and storage of the data, has been chosen to support all project activities in the most convenient way. Specifically that mean local and cloud storage, without password protection or encryption of the raw data.

The field of chemistry, that this study supports computationally, provide research which are used

within the life-science industry for drug-discovery among other fields.

With this in mind, sourcing of the data set needs to be carefully handled, taking into account potential end-user risks of promoting various molecules and processes through research, at some point in the process before consumption. The sourcing of the Cross-Coupling data set, is partly transparent. The initial data set is comprised of 25.116 graph structures, and was sourced through generating options based on 6 transition metals, mentioned above, and 91 ligands. This generation process was made by combining all possible combinations of the six transition metals and the 91 ligands, applying a filter for redundant chemical properties to reach the initial 25.116 size data set[7]. The reduction in data set size to the 7.054, was done by training a model to predict a chemical descriptor value relating to certain energy-levels, and then selecting these 7.054 molecules and transition metals based on that. No apparent precaution towards end-user complications are mentioned, and should therefore be applied further down stream of the research efforts within this field.

8 Results

9 Discussion

10 Conclusion

References

- [1] J. Busk, P. B. Jørgensen, A. Bhowmik, M. N. Schmidt, O. Winther, and T. Vegge, “Calibrated uncertainty for molecular property prediction using ensembles of message passing neural networks,” *Machine Learning: Science and Technology*, vol. 3, 7 2021. [Online]. Available: <https://arxiv.org/abs/2107.06068v2>
- [2] J. Behler, “Atom-centered symmetry functions for constructing high-dimensional neural network potentials,” *Journal of Chemical Physics*, vol. 134, p. 74106, 2 2011. [Online]. Available: [/aip/jcp/article/134/7/074106/954787/Atom-centered-symmetry-functions-for-constructing](http://aip/jcp/article/134/7/074106/954787/Atom-centered-symmetry-functions-for-constructing)
- [3] J. Gastegger, J. Groß, and S. Günnemann, “Directional message passing for molecular graphs,” *8th International Conference on Learning Representations, ICLR 2020*, 3 2020. [Online]. Available: <https://arxiv.org/abs/2003.03123v2>
- [4] K. Atz, F. Grisoni, and G. Schneider, “Geometric deep learning on molecular representations,” *Nature Machine Intelligence*, vol. 3, pp. 1023–1032, 12 2021.
- [5] H. Song, T. Diethe, M. Kull, and P. Flach, “Distribution calibration for regression,” *36th International Conference on Machine Learning, ICML 2019*, vol. 2019-June, pp. 10 347–10 356, 5 2019. [Online]. Available: <https://arxiv.org/abs/1905.06023v1>
- [6] E. Hüllermeier and W. Waegeman, “Aleatoric and epistemic uncertainty in machine learning: an introduction to concepts and methods,” *Machine Learning*, vol. 110, pp. 457–506, 3 2021.
- [7] B. Meyer, B. Sawatlon, S. Heinen, O. A. V. Lilienfeld, and C. Corminboeuf, “Machine learning meets volcano plots: Computational discovery of cross-coupling catalysts,” *Chemical Science*, vol. 9, pp. 7069–7077, 2018.
- [8] K. T. Schütt, O. T. Unke, and M. Gastegger, “Equivariant message passing for the prediction of tensorial properties and molecular spectra,” *Proceedings of Machine Learning Research*, vol. 139, pp. 9377–9388, 2 2021. [Online]. Available: <https://arxiv.org/abs/2102.03150v4>
- [9] B. Lakshminarayanan, A. Pritzel, and C. Blundell, “Simple and scalable predictive uncertainty estimation using deep ensembles,” *Advances in Neural Information Processing Systems*, vol. 2017-December, pp. 6403–6414, 12 2016. [Online]. Available: <https://arxiv.org/abs/1612.01474v3>
- [10] K. Tran, W. Neiswanger, J. Yoon, Q. Zhang, E. Xing, and Z. W. Ulissi, “Methods for comparing uncertainty quantifications for material property predictions,” *Machine*

Learning: Science and Technology, vol. 1, 12 2019. [Online]. Available: <https://arxiv.org/abs/1912.10066v2>

[11] "The radial basis function kernel."

[12] "Using gpus under lsf10." [Online]. Available: https://www.hpc.dtu.dk/?page_id=2759

11 Appendix

11.1 Appendix A

11.2 Appendix B

11.3 Appendix C

Two lists of main packages, and their versions vital to local development and cloud execution.
First the local development packages11.3:

- Python 3.10.9
- Pytorch 2.0.1+cu118
- numpy 1.23.5
- pandas 1.5.2
- matplotlib 3.7.0
- tqdm 4.65.0
- python-dotenv 1.0.0

And for cloud execution11.3:

- Python 3.9.14
- Pytorch 2.1.0+cu118
- numpy 1.23.5
- pandas 1.5.2
- scipy 1.9.1bstat
- tqdm 4.66.1