# Sim2Real transfer for catalyst activity prediction

**Yuta YAHAGI**[1,2]
yuta-yahagi@nec.com

**Kiichi OBUCHI**[1,2]
kiichi-obuchi@nec.com

**Fumihiko KOSAKA**[2]
f.kosaka@aist.go.jp

**Kota MATSUI**[3]
matsui.kota.x3@f.mail.nagoya-u-ac.jp

[1]NEC Corporation, Minato-ku, Tokyo, Japan, 211-8666
[2]National Institute of Advanced Industrial Science and Technology, Tsukuba, Japan, 305-8568
[3]Nagoya University, Nagoya, Japan, 466-8550

## Abstract

Sim2real transfer, knowledge transfer from computational data to experimental data, has received increasing attention as a promising solution to small data problems in materials. We proposed a sim2real transfer method that significantly enhances catalyst activity predictions by harnessing the knowledge of catalyst chemistry. The proposed method transforms the feature space of source computational data into that of target experimental data, and then solves the problem as a homogeneous transfer learning. Through the demonstration, we confirmed that transfer learning model exhibits positive transfer on accuracy and robustness. Notably, significantly high accuracy was achieved despite using a few (less than 10) target data, whose accuracy is compatible with a full scratch model with more than 70 target data. This result indicates that the proposed method leverages the prediction performance with few target data, which helps saving the number of trials in real laboratories.

## 1 Introduction

In recent years, the rapid development of machine learning technology has paved the way for data-driven approaches to materials design and has attracted considerable attention. However, the data-driven approach for materials currently faces the problems of small data due to the high cost of data preparation, the difficulty of formatting data structure, and the exploration bias. The situation is particularly serious for catalyst chemistry, as both preparation and evaluation require significant time and human resources. Although there are attempts to obtain relatively large datasets by running high-throughput experiments [1, 2, 3], data sizes are typically on the order of O(100), which is still below the requirements of advanced methods such as neural networks (NNs).

Instead of conducting experiments in the real world, there are also attempts to produce data employing first-principles calculations [4, 5, 6]. In general, such simulations can generate data at reasonable cost and enable to generate relatively large data sets. Nevertheless, exact reproduction of real materials is unfeasible, thus we have to solve idealised models by sacrificing the fidelity. Therefore, computation data and experimental data are complementary, and this is the motivation for introducing transfer learning; Computational data transfer their abundant knowledge to fill a lack of experimental data, while experimental data complement the low fidelity of computational data. Such a transfer scheme is called as the sim2real (simulation-to-real) transfer.
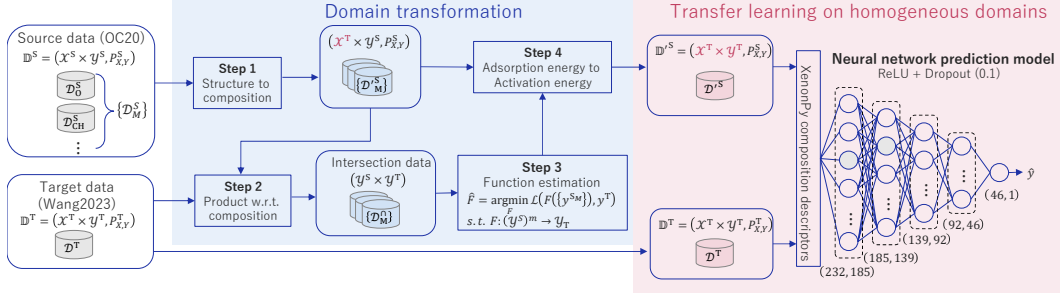
Figure 1: Schematics of the sim2real transfer model.

Despite of those advantages, applying the sim2real transfer approach to materials discovery is not straightforward due to the domain differences between computational and experimental data. Most critical is the difference in scale: whereas a first-principles calculation gives a microscopic picture on a single (often simple) surface structure at a single elementary process of reaction, an experimental measurement gives a macroscopic picture of an entire reaction path on complex surfaces involving various facets, reconstructions, and support interactions. The central challenge addressed in this study is to bridge this gap by harnessing the domain knowledge of catalyst chemistry and enabling sim2real transfer.

This study proposes a transfer learning scheme from calculation to experiment and assesses its feasibility through catalyst activity prediction for the reverse water gas shift reaction (RWGS) as an example. This paper first gives an overview of the proposed method and then presents experimental results of the sim2real transfer.

## 2 Method

### 2.1 Datasets

Throughout this study, we use the *Open Catalyst 2020* (OC20) [4], a simulation dataset based on the density functional theory (DFT), for a source computational dataset. It presents a large collection of pairs of a slab structure $x^{S_M}$ and an adsorption energy $y^{S_M}$ per adsorbate $M$, including approximately $10^6$ entries ($\sim$2000 structures $\times$ 48 adsorbates).

Regarding a target dataset, we refer a high-throughput experimental dataset opened by Wang *et al.* [3]. Their dataset provides 300 pairs of a catalyst composition $x^T$ and an activation energy $y^T$. Here, we disregard the compositions including lanthanoids as they exceed the coverage of the OC20, reducing the data size from 300 to 141.

Let the domain be $\mathbb{D}^r \equiv (\mathcal{X}^r \times \mathcal{Y}^r, P_{X,Y}^r)$ with the feature space, $\mathcal{X}^r$, and the output space, $\mathcal{Y}^r$, and the joint distribution of probability, $P_{X,Y}^r$ where $r = $ S or T for source or target. We express the source dataset as $\{\mathcal{D}^{S_M}\}_M$ where $\mathcal{D}^{S_M} \equiv \{(x_i^{S_M}, y_i^{S_M})\}_{i=1}^{\sim 2000}$ with $x_i^{S_M} \in \mathcal{X}^{S_M}, y_i^{S_M} \in \mathcal{Y}^{S_M}$ and the target dataset as $\mathcal{D}^T \equiv \{(x_i^T, y_i^T)\}_{i=1}^{141}$ with $x_i^T \in \mathcal{X}^T, y_i^T \in \mathcal{Y}^T$.

### 2.2 Model

The proposed method consists of two stages, a domain transformation part and a transfer learning part, as shown in Figure 1. The first part, domain transformation, maps the source domain to the target domain as $\mathbb{D}^S = (\mathcal{X}^S \times \mathcal{Y}^S, P_{X,Y}^S) \rightarrow \mathbb{D}'^S = (\mathcal{X}^T \times \mathcal{Y}^T, P_{X,Y}^S)$ so that $P_{Y|X}^S \approx P_{Y|X}^T$. This procedure reduces the problem to a homogeneous domain shift, enabling standard transfer learning methods to solve.

The transformation proceeds as follows:

**Step 1: conversion from structures to a composition** We first convert structures $\{x^{S_M}\} \in \{\mathcal{X}^S\}^n$ to the corresponding composition $x'^{S_M} \in \mathcal{X}^T$ by marginalizing their atomic positions. Since structures and compositions have a many-to-one relationship with the multiplicity $n$, a corresponding

adsorption energy cannot be determined uniquely. In this work, we simply adopt the lowest adsorption energy among the same compositions as a representative value, that is, $\bar{y}_j^{S_M} = \min\{y_i^{S_M} \mid x_i^{S_M} = x'^{S_M}_j\}_i \in \mathcal{Y}^S$.

**Step 2: extraction of the product set**   Next, we arrange intersection datasets with respect to compositions, that is, $\{\mathcal{D}^{\cap_M}\}_M \equiv \{(\bar{y}_i^{S_M}, y_i^T) \mid x'^{S_M}_i = x_i^T\}_i$. This datasets will be used in the following step to determine a map of $\bar{y}^{S_M}$ onto $\mathcal{Y}^T$.

**Step 3: function estimation for linear scaling**   Let $F : (\mathcal{Y}^S)^m \to \mathcal{Y}^T$ be a projection function from $\{\bar{y}^{S_M}\}_M$ to $y'^S \in \mathcal{Y}^T$. Using $\{\mathcal{D}_M^{\cap}\}_M$, we estimate $F$ by solving a regression problem,

$$\hat{F} = \underset{F}{\operatorname{argmin}} \frac{1}{N} \sum_i^N \mathcal{L}_{MSE}(F(\{\bar{y}_i^{S_M}\}_M), y_i^T), \tag{1}$$

where $\mathcal{L}_{MSE}$ is the mean squared error (MSE) loss function. The function form of $F$ is assumed as a linear model,

$$F(\{\bar{y}^{S_M}\}_M) = f_{M_*}(\bar{y}^{S_{M_*}}), \qquad F_M(y) = a_M y + b_M, \tag{2}$$

where $M_*$ is the adsorbate that best fits $y^T$. This hypothesis relies on the established fact that an activation energy is linearly scaled with the adsorption energy of the rate-determining step[7]. At this process, unanticipated effects involved in the experiments, such as support interactions, may be integrated into the source model, thereby complementing the low-fidelity of simulations.

**Step 4: conversion from adsorption energies to an activation energy**   We convert $\{\bar{y}_i^{S_M}\}_M$ to $y'^S$ with the estimated function $\hat{F}$, resulting a source data on the target space, $\mathcal{D}'^S = \{(x'^S, y'^S)\}$, with the transformed source domain, $\mathbb{D}'^S = (\mathcal{X}^T \times \mathcal{Y}^T, P_{X,Y}^S)$.

These processes correspond to the transformation of microscopic instances into macroscopic ones. Further details are discussed in Appendix B.

The transfer learning pert performs domain adaptation based on a NN model. The model is designed as a fully connected network with five layers. The number of units in each layer decreases progressively, starting from the input dimension, $d$, in the first layer, followed by $0.8d$, $0.6d$, $0.4d$, and $0.2d$, with the final output layer having one unit. All units are activated by ReLU (rectified linear unit) and applied dropout at a rate of $0.1$.

For the model input, the composition is converted into the (standardized) chemical descriptors generated by XenonPy code [8]. While XenonPy originally provides 290 features, we omit the 'sum' features, that is identical to the 'average' features on a compositional data. It reduces the number of features from 290 to 232.

To construct a prediction model, we first prepare a source model from $\mathcal{D}'^S$. With the target training, we compare two different method: the fine-tuning (FT) that re-optimises all layers, and the transfer learning (TL) that fixes the middle layers and retrain only the final layer. During both training phases, the parameters are optimized so that minimizing the mean squared error loss by the stochastic gradient descent algorithm with 16 batches, which is implemented in the PyTorch framework [9].

## 3   Results

### 3.1   Function estimation for linear scaling

The linear scaling function $\hat{F}$ was estimated by the linear regression in Eqs. (1) and (2). To determine $M_*$, the regression performance for each $M$ was scored with the coefficient of determination, $R^2$. Figure 2 shows the results of the best 5 adsorbates in order of $R^2$. As the result, we eventually obtained $\hat{F} = 0.136 y^{S_{OH}} + 0.663$ with $R^2 = 0.75$ from $\mathcal{D}^{\cap_{OH}}$ that has 9 compositions.

### 3.2   Training and catalyst activity prediction

We constructed a prediction model according to the method shown in Section 2.2. After domain transformation, the source model was trained using $\mathcal{D}'^S$, then the prediction models were prepared by
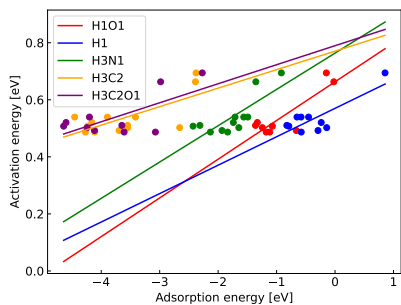
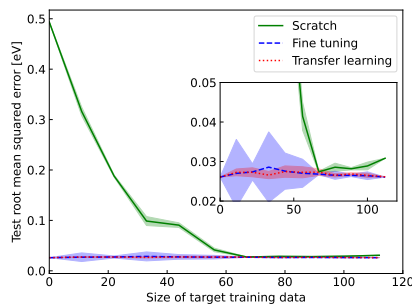Figure 2: Linear regression of activation energy by adsorption energy.



Figure 3: Test root mean squared error as a function of the size of training data.

retraining the source model on $\mathcal{D}^T$. Here, we use 80% of the target data for retraining and separate the remaining 20% as the test data. To prevent leaks, test data were sampled avoiding $\mathcal{D}^{\cap OH}$ whose data have already been used. Further details of training are presented in Appendix. C.

The effects of sim2real transfer were measured by evaluating the test loss as a function of the target training data size, as shown in Figure 3. Here, we evaluate the mean loss and its standard deviation over five trials where the data were randomly selected from the training data pool in each trial.

As the result, we observed a sign of positive transfer; Both the FT and the TL model show lower losses than the full scratch model. In comparison with the FT and the TL, although the later is more stable as showing the smaller standard deviations than the former, there are little difference. It is noteworthy that the losses of transferred models have minimized even without the target training steps, indicating that the source model has almost converged on optimal parameters in the target space. Although the source model knows only 9 data used in estimating $F$, it achieves an accuracy of 0.026eV, which is comparable to that of full scratch model with $> 70$ target data. In addition, the transferred models are more robust than the full scratch model that shows a tendency of over-fitting as decreasing accuracy with adding training data. We also performed some additional examinations, shown in Appendix. D

## 4 Conclusion

In this work, we proposed a machine learning method for catalyst activity prediction based on transfer learning from computational data to experimental data across their gap harnessing the knowledge of catalyst chemistry.

The proposed method transforms the feature space of microscopic computational data into that of macroscopic experimental data by using the linear scaling relation in catalysis, and then solves the problem as a homogeneous transfer learning.

Using the *Open Catalyst 2020* as the source dataset and the high-throughput experimental dataset from Wang *et al.* as the target dataset, we construct a neural network prediction model for activation energy of the reverse water gas shift reaction. Our demonstration confirmed that the proposed method shows a positive transfer of increasing accuracy and robustness compared to the model trained only with target data. Moreover, the source model has already been optimised in the target space, despite only using few target data during the domain transformation. This result indicates the importance of appropriate domain transformation, and suggests that the prediction performance can be improved efficiently by expanding the data so that there is more overlap between calculations and experiments.

In conclusion, we believe that our proposed method helps saving the number of trials in real laboratories significantly, leading drastic reduction of the cost and time of exploring novel catalysts. Since this work is still at the stage of proof of concept, there is a number of implications for further improvement. Future works should include the improvement of each procedure, for instance, designing the mapping functions $F$ and $G$, or tailoring the architecture of prediction models.

## Acknowledgement

## References

[1] Thanh Nhat Nguyen, Thuy Tran Phuong Nhat, Ken Takimoto, Ashutosh Thakur, Shun Nishimura, Junya Ohyama, Itsuki Miyazato, Lauren Takahashi, Jun Fujima, Keisuke Takahashi, and Toshiaki Taniike. High-throughput experimentation and catalyst informatics for oxidative coupling of methane. *ACS Catal.*, 10(2):921–932, January 2020.

[2] Sunao Nakanowatari, Thanh Nhat Nguyen, Hiroki Chikuma, Aya Fujiwara, Kalaivani Seenivasan, Ashutosh Thakur, Lauren Takahashi, Keisuke Takahashi, and Toshiaki Taniike. Extraction of catalyst design heuristics from random catalyst dataset and their utilization in catalyst development for oxidative coupling of methane. *ChemCatChem*, 13(14):3262–3269, July 2021.

[3] Gang Wang, Shinya Mine, Duotian Chen, Yuan Jing, Kah Wei Ting, Taichi Yamaguchi, Motoshi Takao, Zen Maeno, Ichigaku Takigawa, Koichi Matsushita, Ken-Ichi Shimizu, and Takashi Toyao. Accelerated discovery of multi-elemental reverse water-gas shift catalysts using extrapolative machine learning approach. *Nat. Commun.*, 14(1):5861, September 2023.

[4] Lowik Chanussot, Abhishek Das, Siddharth Goyal, Thibaut Lavril, Muhammed Shuaibi, Morgane Riviere, Kevin Tran, Javier Heras-Domingo, Caleb Ho, Weihua Hu, Aini Palizhati, Anuroop Sriram, Brandon Wood, Junwoong Yoon, Devi Parikh, C Lawrence Zitnick, and Zachary Ulissi. Open catalyst 2020 (OC20) dataset and community challenges. *ACS Catal.*, 11(10):6059–6072, May 2021.

[5] Keisuke Takahashi, Lauren Takahashi, Son Dinh Le, Takaaki Kinoshita, Shun Nishimura, and Junya Ohyama. Synthesis of heterogeneous catalysts in catalyst informatics to bridge experiment and high-throughput calculation. *J. Am. Chem. Soc.*, 144(34):15735–15744, August 2022.

[6] Richard Tran, Janice Lan, Muhammed Shuaibi, Brandon M Wood, Siddharth Goyal, Abhishek Das, Javier Heras-Domingo, Adeesh Kolluru, Ammar Rizvi, Nima Shoghi, Anuroop Sriram, Félix Therrien, Jehad Abed, Oleksandr Voznyy, Edward H Sargent, Zachary Ulissi, and C Lawrence Zitnick. The open catalyst 2022 (OC22) dataset and challenges for oxide electrocatalysts. *ACS Catal.*, 13(5):3066–3084, March 2023.

[7] J K Nørskov, T Bligaard, A Logadottir, S Bahn, L B Hansen, M Bollinger, H Bengaard, B Hammer, Z Sljivancanin, M Mavrikakis, Y Xu, S Dahl, and C J H Jacobsen. Universality in heterogeneous catalysis. *J. Catal.*, 209(2):275–278, July 2002.

[8] Chang Liu, Erina Fujita, Yukari Katsura, Yuki Inada, Asuka Ishikawa, Ryuji Tamura, Kaoru Kimura, and Ryo Yoshida. Machine learning to predict quasicrystals from chemical compositions. *Adv. Mater.*, 33(36):e2102507, September 2021.

[9] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An imperative style, high-performance deep learning library. *arXiv [cs.LG]*, December 2019.

## A Dataset preparation and preprocessing

### A.1 Computational dataset

*Open Catalyst 2020* (OC20) [4] is a set of simulated data on molecular adsorption processes at metal surfaces, obtained by first-principles calculations based on the density functional theory (DFT). While

it consists of several tasks, we particularly focus on the initial structure to relaxed energy (IS2IR) dataset because we are interested in the adsorption energy. IS2IR is a large collection of pairs of slab structures $S$ and adsorption energies $E^M$ per adsorbate $M$, including approximately $10^6$ entries ($\sim$2000 catalysts $\times$ 48 molecules).

Let the source domain be $\mathbb{D}^S \equiv (\mathcal{X}^S \times \mathcal{Y}^S, P^S_{X,Y})$ with the source input (output) space, $\mathcal{X}^S(\mathcal{Y}^S)$, and the joint distribution on source space, $P^S_{X,Y}$. We express the source dataset as $\{\mathcal{D}^S_M\}$ with $\mathcal{D}^S_M \equiv \{(x^{S_M}_i, y^{S_M}_i)\}^{\sim 2000}_{i=1}$, that is, $x^{S_M}_i \in \mathcal{X}^S$ and $y^{S_M}_i \in \mathcal{Y}^S$.

### A.2 Experimental dataset

Wang *et al.* presents measured data on the activity of RWGS on Pt-based catalysts synthesised on TiO2 supports [3]. The catalysts are prepared such that their compositions are Pt(3)/$M1(X1)$-$M2(X2)$-$M3(X3)$-$M4(X4)$-$M5(X5)$/TiO$_2$ where $M_i$ is an element contained in a loading amount $X_i$ (wt%) with the 3 wt% Pt, on the TiO$_2$ support. This dataset consists of pairs of loading amount and CO formation rate, $r_{CO}$, and is one of the largest experimental dataset of catalyst with a total of 300 entries (45 initial compositions and 255 compositions obtained by Bayesian search). Here, we disregard the compositions including lanthanoids as they exceed the coverage of OC20, reducing the entries to 141 (38+103).

For simplicity, we assume that each catalyst is characterised by a single compositional data among $\{M_i\}$ plus Pt as $x^T$. We derive the activation energy $y^T$ from the CO formation rate using the Arrhenius equation,

$$r_{CO} = A\exp\left(\frac{y^T}{kT}\right),\tag{3}$$

where $k$ is the Boltzmann constant, $T$ is the experimental temperature, $A$ is called the pre-exponential factor. $y^T$ is measured in the unit of eV.

In this work, $A$ is represented by the value of Pt(3)/Rb(1)-Ba(1)-Mo(0.6)-Nb(0.2)/TiO$_2$, irrespective of composition. Eventually, we consider the target dataset as $\mathcal{D}^T \equiv \{(x^T_i, y^T_i)\}^{141}_{i=1}$ with $x^T_i \in \mathcal{X}^T$ and $y^T_i \in \mathcal{Y}^T$ on the target domain $\mathbb{D}^T \equiv (\mathcal{X}^T \times \mathcal{Y}^T, P^T_{X,Y})$.

## B  Detailed discussions on transformation functions

Physically, Steps 2 and 3 in the transformation part are the projections from microscopic to macroscopic and may be defined according to the theory of thermodynamics.

Step 2 applies a function that discards the information on atomic positions from the structure and extracts only its composition. In general, structures and compositions have the many-to-one relationship, and two or more structures can give the same composition. Ideally, it should be a good approximation of adsorption on real surfaces by taking thermostatistical averages for possible (meta-)stable surfaces and adsorption sites. However, as the OC20 is based on random sample from low-index surfaces of relatively simple alloys, it does not necessarily include the active sites of real catalysts. Since exploration of stable surfaces is beyond the scope of this study, we adopt the structure that gives the lowest adsorption energy among duplicates as the representative sample instead.

$F$ mediates between the adsorption energies and the activation energy, and is related to microkinetic analysis of the reaction path, assumed as Eqs. (1) and (2). This assumption is based on the following empirical row, which is well-established in catalyst chemistry:

1. In many cases, catalytic activity is determined by a rate-determining step (RDS) in the reaction path. Thus, an adsorbate related to the RDS is sufficient for consideration.

2. There is a linear scaling between the activation energy and the adsorption energy in RDS [7], known as the Brønsted–Evans–Polanyi relation.

With regard to the latter, if we consider the Sabatier principle, which states that adsorption energy should be neither too strong nor too weak as it follows a volcano curve, we can improve the prediction by extending a piecewise linear model. However, in this study, a simple linear model was employed due to the limited quantity of training data.
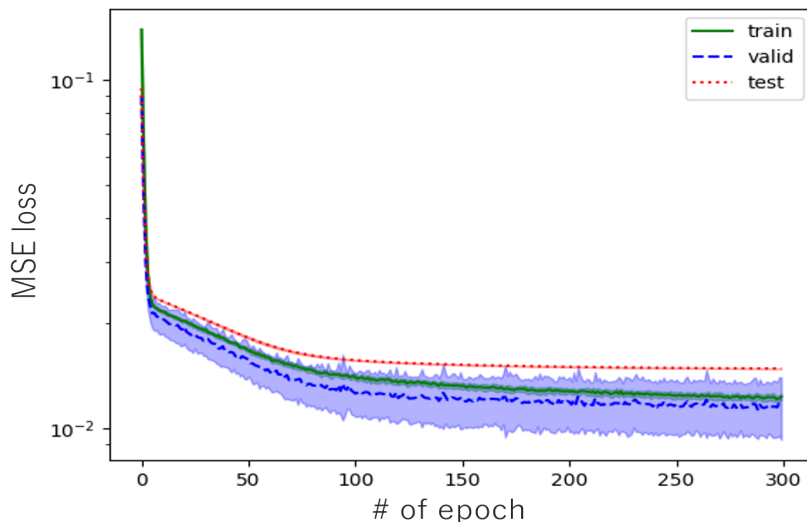
Figure 4: Source training curve with respect to the losses of training, validation, and test, as functions of the number of epochs.

## C    Experimental details

### C.1    Source training

Figure 4 shows the training curves with source data, displaying the mean and the standard deviation of losses obtained by 10-fold cross validation. Here, we separated the 20% of source data for testing, and performed cross validation with the remaining 80%. Apart from the validation, the source model was trained using all the source data, and carried out early stopping on 100 epochs.

### C.2    Target training

Figure 4 shows the target training curves for each retraining methods of the full scratch, the fine tuning, and the transfer learning, explained in Section 2.2. The mean and the standard deviation of validation losses obtained by 10-fold cross validation are displayed. With the prediction models, we stopped training after 100 epochs.

## D    Supporting experiments

We conducted some additional experiments to analyse the results discussed in Section 3.2. Since the fine tuning and the transfer learning exhibit almost identical performance, we only show the results with the later. First, to investigate the impact of choosing $M_*$, we compare the case if we differentiate $M_* = \text{OH}$ to, for instance, $M_* = \text{CH}_3\text{CO}$, that gives $\hat{F} = 0.0667 E^{\text{CH}_3\text{CO}} + 0.790$ with $R^2 = 0.47$. Figure 6 shows the comparison, suggesting that a wrong choice of $M_*$ degrades accuracy. Interestingly, even in such case, the transferred model starts with better initial conditions and gets optimised earlier than the full scratch model.

In order to check whether the parameters fall into a local minimum or not, we apply a white noise to the parameters before starting target training. Figure 7 shows the of retraining the model with white noise of different amplitudes, 0.025, 0.05, and 0.1. We can see the performance degradation due to noises, but if it is small, the performance will be recovered while learning the target data. Otherwise, the test error converges on a worse value. This result confirms that the pre-training model reached an optimal solution, at least under these conditions.
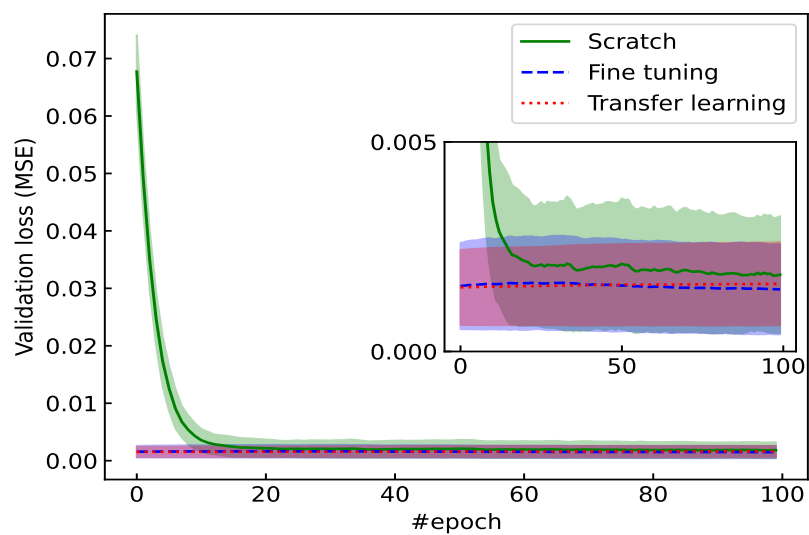
Figure 5: Target training curve with respect to the validation losses for each training methods, as functions of the number of epochs.
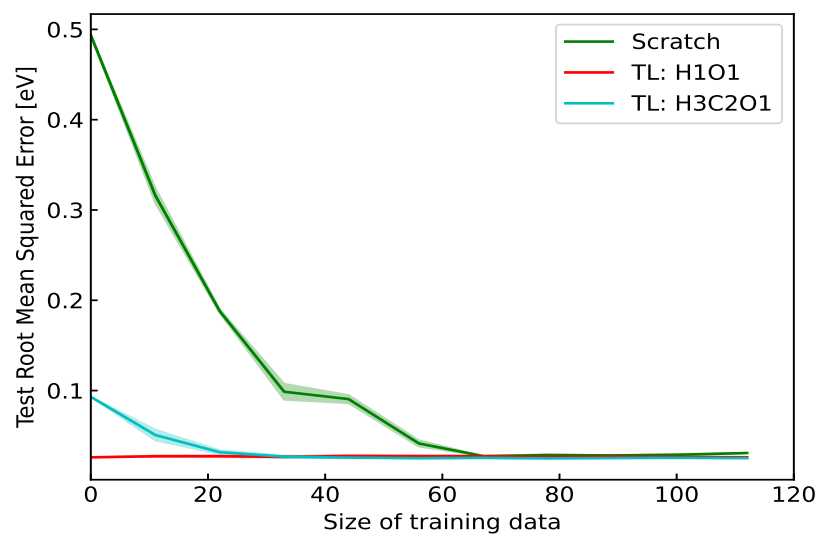


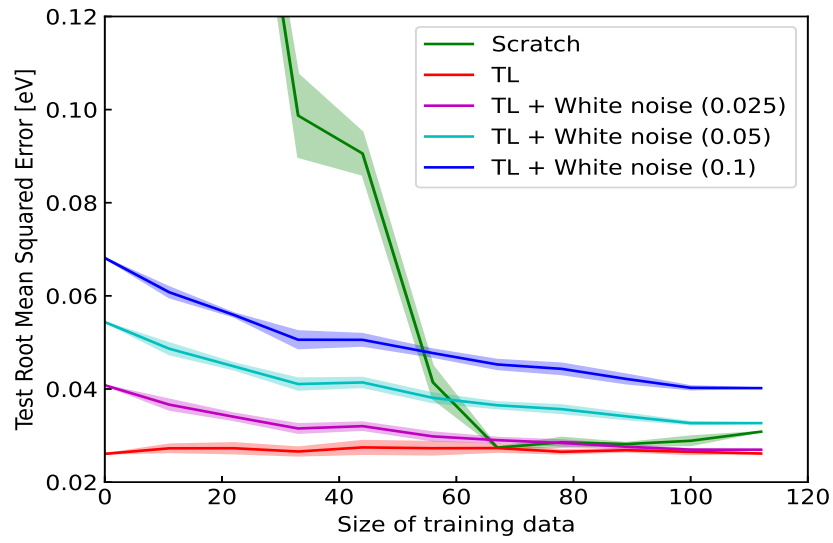Figure 6: Performance evaluation for each adsorbed molecule.

8

Figure 7: Performance evaluation for each amplitude of white noise.