

# Activity Descriptors of Mo<sub>2</sub>C-based Catalysts for C-OH Bond Activation

Raghavendra Meena,<sup>†,‡</sup> Michael J. Purcell,<sup>†</sup> Wouter Kluijtmans,<sup>†</sup> Han Zuilhof,<sup>‡,¶</sup>  
Johannes H. Bitter,<sup>†</sup> Runhai Ouyang,<sup>§</sup> and Guanna Li<sup>\*,†</sup>

<sup>†</sup>*Biobased Chemistry and Technology, Wageningen University & Research, Bornse Weilanden 9, 6708WG Wageningen, The Netherlands*

<sup>‡</sup>*Laboratory of Organic Chemistry, Wageningen University & Research, Stippeneng 4, 6708 WE Wageningen, The Netherlands*

<sup>¶</sup>*College of Biological and Chemical Engineering, Jiaxing University, Jiaxing 314001, China*

<sup>§</sup>*Materials Genome Institute, Shanghai University, No. 99 Shangda Road, Shanghai 200444, China*

E-mail: guanna.li@wur.nl

## Abstract

Deriving structure-activity relationships is crucial for designing efficient catalysts. To aid in this quest, data-driven methods such as machine learning (ML) are emerging. In this work, we incorporate ML tools with accurate density functional theory (DFT) energetics, electronic and geometric features, in the context of designing efficient molybdenum carbide ( $\text{Mo}_2\text{C}$ )-based catalysts for biomass conversion. Previously, it was shown that C-OH activation is the rate-determining step in the hydrodeoxygenation (HDO) reaction. Therefore, in this work, DFT was used to obtain accurate barriers ( $E_a$ ) and reaction energy ( $\Delta E$ ) for the C-OH activation in HDO over the most stable (111) and meta-stable facets (010, 101, and 110) of  $\text{Mo}_2\text{C}$  and transition metal doped  $\text{Mo}_2\text{C}$ . The 101 facet was identified as the most active facet for C-OH activation. While, doping of the active site with Zr and Nb was identified as promising strategy to improve the activity. Further, scikit-learn's ML models were used to obtain the best primary features correlating with the  $E_a$ . Ridge Regression (RR) gives the best ML model for predicting  $E_a$  with a test RMSE of 0.21 eV. SHapley Additive exPlanations (SHAP) analysis was then performed which reveals that  $\Delta E$  and  $d$ -band center are the most important features contributing to the activity. Finally, SISSO was used to validate RR model and SHAP results, and to obtain a 2-dimensional physically interpretable generic-descriptor comprising of dopant's local environment,  $d$ -band features, and  $\Delta E$ .

## Introduction

Understanding structure-activity relationships of catalysts is a fundamental challenge in heterogeneous catalysis, crucial for rational catalyst design. High-throughput computational chemistry and data-driven methods have been developed to address this challenge. Most notably, Nørskov et al. have pioneered the use of the Brønsted-Evans-Polanyi (BEP) principle and linear scaling (LS) relationships to identify critical electronic and energy features of metal catalysts that can correlate with external catalyst performance metrics such as reac-

tion rate and TOF.<sup>1–3</sup> The BEP relationship has since been applied extensively in de-signing metal catalysts with improved activity for various reactions.<sup>4,5</sup> However, with rapidly evolving catalyst space, complex catalysts such as transition metal oxides often deviate from the BEP principle.<sup>6,7</sup> LS relationships also break down for new and upcoming complex catalysts such as transition metal carbides (TMCs).<sup>8</sup> The BEP and LS relationships can break down because their linear, few-descriptor form is specific to site geometry and mechanism.<sup>9,10</sup> Further, almost intrinsic to heterogeneous catalysis, the nature of the active site is rather complex, and the complicated reaction mechanism can often not be mapped by only using one or two descriptors and a simple linear regression method.

To overcome the above-mentioned shortcomings, recently, exponential growth in computational capacity, coupled with efficient machine learning (ML) algorithms, has enabled the derivation of complex expressions that capture the nature of active sites and their structure–activity relationships.<sup>11,12</sup> Various ML methods have been developed and multi-dimensional descriptors beyond linear scaling relationships have been identified to improve structure–activity correlation and catalytic understanding.<sup>11</sup> Among these, the Sure Independence Screening and Sparsifying Operator (SISSO) has recently emerged as a particularly powerful approach for constructing interpretable descriptors and capturing non-linear correlations, as discussed below.<sup>13,14</sup> SISSO enables complex feature construction and estimation of catalyst activity by finding the best physically interpretable descriptors, even when only a small training set is available. For example, Wang et al.<sup>15</sup> established a general theory of metal-support interactions (MSIs) for metal catalysts on oxide supports, grounded in both metal-metal interactions (MMIs) and metal-oxygen interactions (MOIs). Using a large dataset of experimentally measured adhesion energies, interpretable machine learning, and theoretical derivation, these authors derived a predictive and interpretable formula for MSIs. Based on this they showed that for late transition metal catalysts, MMIs dominate the support effects and encapsulation behavior, and formulated a principle that strong

MMIs, rather than strong oxophilic MSIs, determine encapsulation occurrence. This theory is validated by extensive experiments and simulations, providing a comprehensive framework for understanding and designing supported metal catalysts. Building upon this, Shu et al.<sup>12</sup> identified that the topologically undercoordinated number, valence electron count, lattice constant, and the reaction energy can be combined to form a 2-dimensional descriptor that serves as the best interpretable descriptor for structure sensitivity and reaction barriers. Similarly, Xu et al.<sup>16</sup> applied SISSO for predicting the adsorption enthalpies of the oxygen evolution reaction (OER) in-termediates over doped RuO<sub>2</sub> and IrO<sub>2</sub> surfaces. These authors reported a test RMSE as low as 0.18 eV using clean-surface primary features (electronegativity, *d*-band center, and *d*-band kurtosis), and 0.12 eV when including the O\* adsorption enthalpy as an additional feature. Overall, the SISSO-obtained de-scriptor revealed local charge-transfer-related features as critical in predicting accurate adsorption enthalpies. More recently, Lin et al.<sup>7</sup> identified an optimal SISSO-derived descriptor for CO dissociation barriers in iron-based Fischer–Tropsch catalysis. These authors derived a 4-dimensional descriptor combining five features: work function, C-vacancy formation energy, CO adsorption energy, coordination number, and active-site size, with the dominant contribution arising from the C-vacancy formation energy. In combination with the reaction energy term, this 4-dimensional descriptor effectively captured structure sensitivity and reaction barriers. Collectively, these studies demonstrate the versatility and robustness of SISSO, also show how it enables the discovery of material-specific, interpretable descriptors that go beyond simple linear scaling relations and extend the applicability of ML models in catalysis and materials science.

Many reactions on the orthorhombic molybdenum carbide (Mo<sub>2</sub>C) 101 surface cannot be accurately simulated by one- or two-dimensional parametrizations, given the structural complexity of both the catalytically active surface and the reactants.<sup>17,18</sup> A good case in point is the hydrodeoxygenation (HDO) of butyric acid over orthorhombic molybdenum car-

bide ( $\text{Mo}_2\text{C}$ ) 101 surface. In previous work butanol dissociation, specifically C–OH bond cleavage, was determined as the rate-determining step (RDS),<sup>17,19</sup> More recently, we demonstrated that C–OH bond dissociation plays a similarly crucial role in the case of  $\text{W}_2\text{C}$ -based catalysts as well.<sup>20</sup> Subsequently, we modified the surfaces of  $\text{Mo}_2\text{C}$  and tungsten carbide ( $\text{W}_2\text{C}$ ) by introducing oxygen to model *in situ* formed oxycarbide-like species ( $\text{Mo}_2\text{CO}_x$  and  $\text{W}_2\text{CO}_x$ ).<sup>20</sup> Across both carbide and oxycarbide surfaces, microkinetic modeling (MKM) results consistently revealed that C–OH bond dissociation in butanol is among the most kinetically challenging steps. Based on these studies, we concluded that the C–OH activation barrier serves as a reliable descriptor for the overall catalytic activity of Mo and W carbide-based systems.

However, accurately estimating this barrier is highly sensitive to the choice of the initial surface model, and locating the transition state is particularly challenging due to aforementioned structural complexity of both the catalyst surfaces and the reactants. Additionally, transition state (TS) calculations are computationally expensive and time-consuming. In our previous work, we performed the *d*-band analysis in an effort to correlate the C–OH activation barrier,  $E_a(\text{C-OH})$ , with various electronic and geometric descriptors. We observed a notable correlation between  $E_a(\text{C-OH})$  and properties such as the atomic radius of the doped active metal site and its *d*-band filling. Additionally, a BEP relationship between C–OH activation barrier and  $\Delta E$  was obtained ( $R^2 = 0.79$ ). While these correlations provided qualitative insights, they were insufficient for quantitatively predicting activation barriers as these correlations were based on a very small dataset (1 phase, 1 facet, 11 dopants).

Therefore in the current work we aim to obtain a generic descriptor correlating structure and activity in hydrodeoxygenation reactions on  $\text{Mo}_2\text{C}$  catalysts. In this, we cover both the most stable facet 111 (with 3 unique facet terminations) and three meta-stable facets (101, 010, and 110), and for each facet, 11 dopants (Ti, V, Cr, Fe, Co, Ni, Zr, Nb, W, Pt,

and Au) were taken into consideration to evaluate their influence on the C-OH activation barrier. To achieve this goal, we evaluated scikit-learn ML models,<sup>21</sup> specifically including the Ridge Regressor, to predict the activation barrier, and subsequently applied SHapley Additive exPlanations (SHAP) - a unique and powerful analysis tool traditionally used and designed specifically for biomedical applications to understand the output of a ML model - for a quantitative evaluation.<sup>22,23</sup> Based on these outcomes, we then used SISSO<sup>14</sup> to derive an interpretable mathematical model that retains the underlying physics. Finally, the important features identified using the scikit-learn ML model and SHAP analysis were compared with the best SISSO descriptors, which clearly points to the significant potential of this approach. Overall, we successfully obtained a 2-dimensional descriptor engineered using primary energy ( $\Delta E$ ), electronic ( $\varepsilon_{d_{sk\downarrow}}$ ,  $\varepsilon_{d_{f\downarrow}}$ , BE<sub>M-O</sub>), and geometry (R<sub>d</sub>, CN<sub>C</sub>) features. These results further validated the SHAP analysis as all of these features, except BE<sub>M-O</sub>, were appeared as the most important features.

## Results and Discussion

### Data Generation

The C-OH bond breaking barriers of *n*-butanol (C<sub>4</sub>H<sub>8</sub>OH → C<sub>4</sub>H<sub>8</sub> + OH) were calculated over the stable orthorhombic Mo<sub>2</sub>C-based catalysts. For this purpose, the most stable facets of the orthorhombic Mo<sub>2</sub>C catalyst reported in the literature were used, as shown in Figure 1, i.e., 010, 101, 110, 111 (ter1), 111 (ter2), and 111 (ter3).<sup>24</sup> The Wulff construction in Figure 1 is obtained at a carbon chemical potential of -10.1 eV, i.e., at a carburization ability similar to CH<sub>4</sub>/H<sub>2</sub>, which is used for preparing Mo<sub>2</sub>C catalysts. All the facets used were Mo/C mixed-terminated to provide multiple active sites for the reaction to occur, making them more active surfaces in general.<sup>25</sup>

Table 1 (**Mo** column) shows that for pure Mo<sub>2</sub>C catalysts, 101 surface is the most active, followed by 010, 111-ter1, 111-ter3, 111-ter2, and 110 surfaces. Further, in all the facets, the

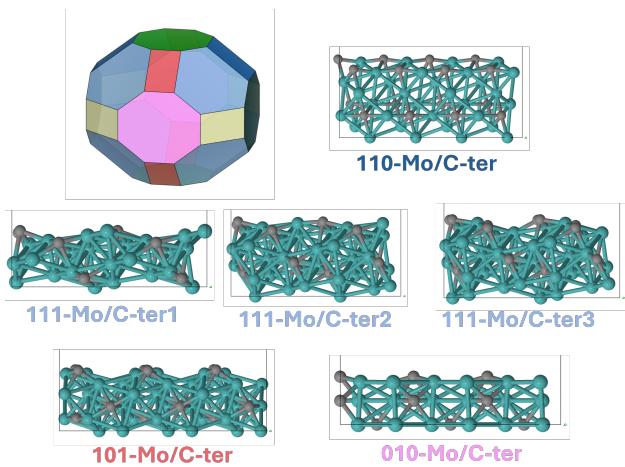


Figure 1: Wulff construction for orthorhombic  $\text{Mo}_2\text{C}$  catalyst and the facets chosen in this work based on the surfaces exposed.

active Mo site was doped with the following transition metals: Ti, V, Cr, Fe, Co, Ni, Zr, Nb, W, Pt, or Au. The C-OH activation barriers of butanol were evaluated upon heteroatom doping of the active metal site, and reported in Table 1. These results show that doping the active metal site with Zr and Nb reduces the C-OH activation barrier over all facets, irrespective of the specific surface structures. This further supports our previous work on butanol dissociation over  $\text{Mo}_2\text{C}$  (101) surface, where it was found that doping with Zr or Nb improves the activity for C-OH bond activation.<sup>19</sup> It is also evident that meta-stable facets (010, 110, and 101) are more active than the most stable 111 facet.

Table 1: The C-OH activation barriers as a function of facets and dopants (in eV).

Facet/dopant	Ti	V	Cr	Fe	Co	Mo	Ni	Zr	Nb	W	Pt	Au
010	0.80	0.88	0.99	1.16	0.97	0.94	1.01	0.71	0.71	0.78	1.28	NaN
110	0.70	1.26	1.60	NaN	NaN	1.58	1.07	1.01	0.91	1.01	1.64	2.34
101	0.95	0.77	1.00	1.32	1.63	0.83	1.59	0.54	0.58	0.78	1.60	NaN
111-ter1	1.02	1.04	1.34	NaN	NaN	1.25	1.53	0.85	0.82	0.81	1.70	1.88
111-ter2	1.2	1.24	1.44	NaN	1.94	1.54	NaN	1.05	0.93	0.96	2.35	NaN
111-ter3	1.23	1.21	1.46	1.90	NaN	1.26	NaN	0.98	0.88	0.97	2.35	2.50

\*NaN: the cases where a TS was not stabilized and hence excluded from this work.

Further, electronic structure analysis and geometry analysis were performed to obtain the potential primary features of activity in  $\text{Mo}_2\text{C}$ -based catalysts for C-OH activation. The data

Table 2: List of energy, electronic, and geometric features considered for ML models and SISSO analysis.

Electronic/Geometric feature	Symbol
Atomic radius	$R_d$
Coordination number (dopant with metal)	$CN_M$
Coordination number (dopant with carbon)	$CN_C$
Coordination number (dopant)	$CN_{total}$
Common coordination number (in stable oxides)	$CN_{oxides}$
Closest dopant-metal bond distance	$d_{MM}$
Closest dopant-carbon bond distance	$d_{MC}$
$d$ -band centre (up spin/down spin)	$\varepsilon_{d_{c\uparrow}}/\varepsilon_{d_{c\downarrow}}$
$d$ -band filling (up spin/down spin)	$\varepsilon_{d_{f\uparrow}}/\varepsilon_{d_{f\downarrow}}$
$d$ -band skewness (up spin/down spin)	$\varepsilon_{d_{sk\uparrow}}/\varepsilon_{d_{sk\downarrow}}$
$d$ -band kurtosis (up spin/down spin)	$\varepsilon_{d_{k\uparrow}}/\varepsilon_{d_{k\downarrow}}$
$d$ -band width (up spin/down spin)	$\varepsilon_{d_{w\uparrow}}/\varepsilon_{d_{w\downarrow}}$
Number of valence electrons	$V_e$
Pauling's electronegativity (dopant)	$\chi_P$
Electron affinity	$EA$
1st ionization energy (dopant)	$IP^{1st}$
Binding energy (dopant-oxygen) in stable oxides	$BE_{M-O}$
Reaction energy (C-OH activation)	$\Delta E$

we obtained contain the C-OH activation barrier ( $E_a$ ), reaction energy ( $\Delta E$ ), bond distance between the dopant (M) and the closest Carbon atom ( $d_{MC}$ ), bond distance between the dopant and the closest molybdenum atom ( $d_{MM}$ ), coordination number of the active metal site with metal or carbon in the first-shell ( $CN_{M/C}$ ), and the atomic radius of the active metal site ( $R_d$ ). Furthermore, clean surface properties such as active metal site's  $d$ -band filling,  $d$ -band center,  $d$ -band kurtosis,  $d$ -band skewness,  $d$ -band width, electronegativity, first-ionization energy, electron affinity, and binding energy with Oxygen (BEM-O) are evaluated. In total, this implies that a set containing 61 samples, each containing information about 23 primary features, were investigated (Table 2). The correlation matrix for the different primary features is shown in Figure 2. These correlations show that there are some features correlated with each other; however, no single rigorous descriptor has a significant correlation with the targeted C-OH activation barrier ( $E_a$ ).

From the correlation matrix and linear regression,<sup>26</sup> we derive the BEP relationship

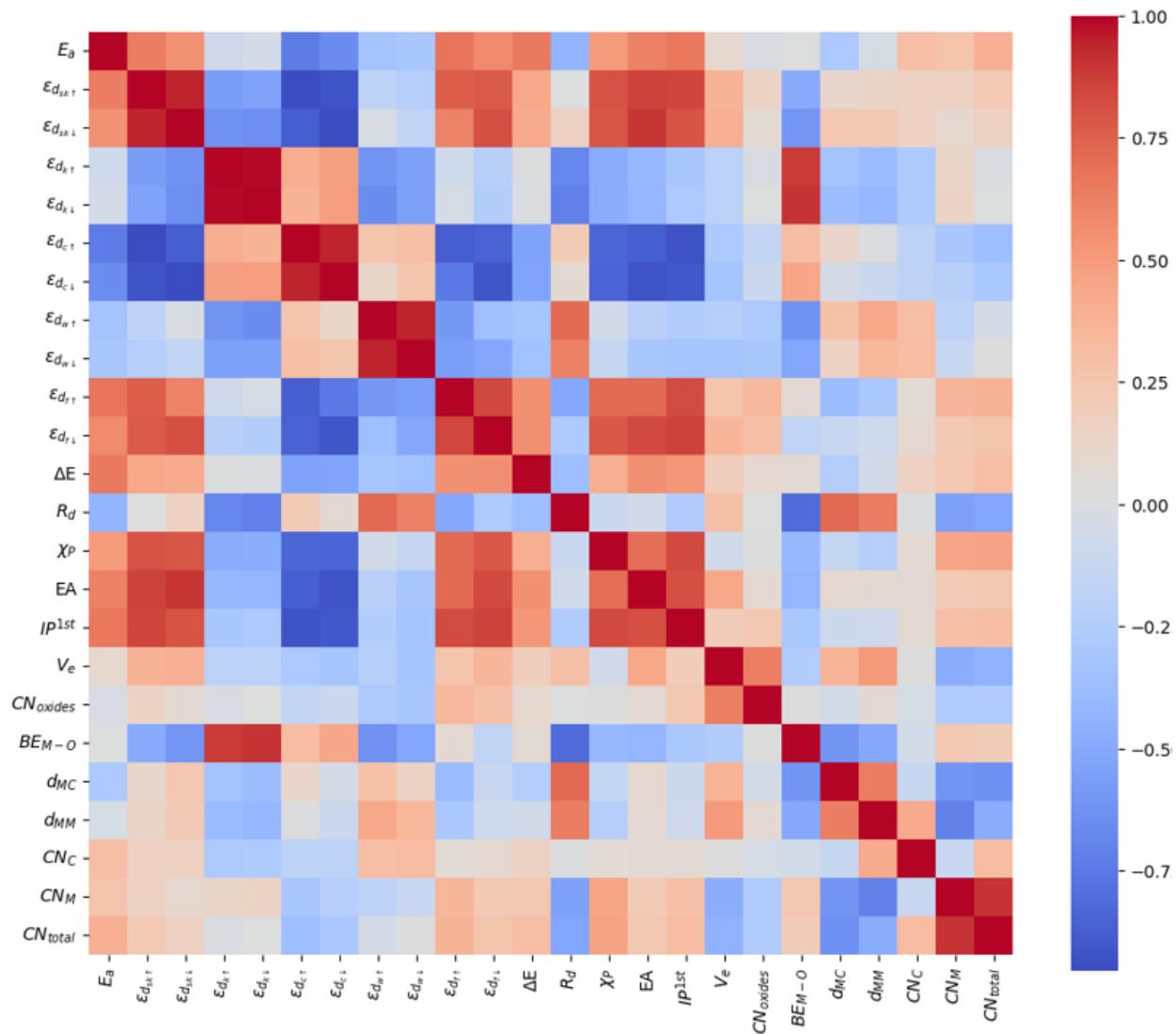


Figure 2: Correlation matrix between the activation barrier of C-OH bond activation  $E_a$  and primary energy, electronic, and geometric features.

between the activation energy  $E_a$  and reaction energy  $\Delta E$  for all the facets and dopants combined. The  $R^2$  score of 0.42 indicates that the different facets of Mo<sub>2</sub>C behave differently, and confirm that the BEP relationship is broken, as only facets 111-ter3 ( $R = 0.99$ ), 101 ( $R = 0.79$ ), and 111-ter1 ( $R = 0.77$ ) show a strong BEP relationship, while others facets such as 111-ter2 ( $R = 0.48$ ) and 110 ( $R = 0.37$ ) show a very weak linear BEP relation-ship. In contrast, the 010 facet ( $R = 0.06$ ) does not show a BEP relationship at all. Hence, more rigorous scikit-learn-based regressors were applied.

## Prediction using scikit-learn based ML methods

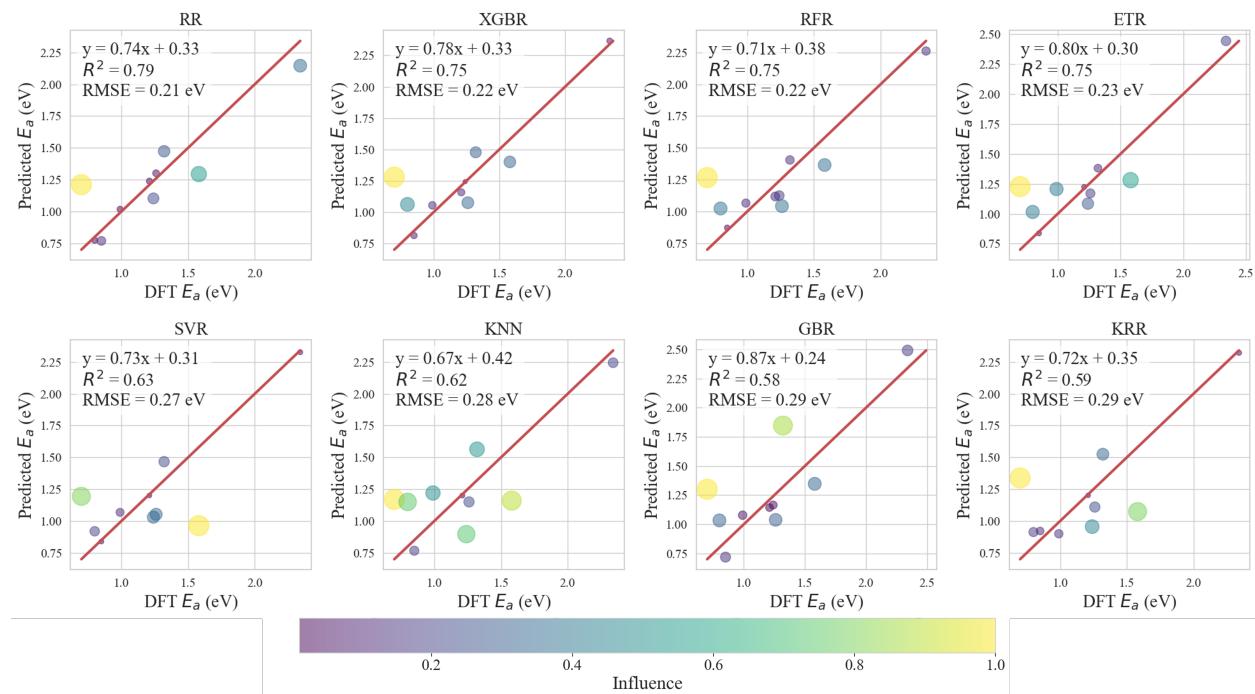


Figure 3: Test *RMSE* for the different scikit-learn ML regressors. Note: the influence is quantified by Cook's distance - a statistical measure used in regression analysis to identify influential data points that have a large impact on the regression model's coefficients.

The regression problem was first approached using the ML algorithms present within Python's scikit-learn package. The ML models used here are K-Nearest Neighbors Regressor (KNN), Kernel Ridge Regressor (KRR), Gradient Boosting Regressor (GBR), Ridge Regressor (RR), Extreme Gradient Boosting Regressor (XGBR), Support Vector Regressor (SVR), Random Forest Regressor (RFR), and Extra Trees Regressor (ETR). For all these regressors, the 85:15 training-to-testing ratio was used for splitting. The results of these algorithms for prediction of  $E_a(\text{C-OH})$  are in Figure 3. Ridge Regression (RR) emerges as the best performer ( $R^2=0.79$ , RMSE = 0.21 eV), indicating that the relationship between features and DFT activation energies is largely linear, so a simple but properly regularized linear model suffices. Further, tree-based methods (RFR, ETR, XGBR) also perform well, as they flexibly capture both linear and non-linear trends, though without outperforming RR. In contrast, kernel-based (SVR, KRR) and distance-based (KNN) models perform worse,

reflecting sensitivity to hyperparameters and limited data. On a par with KRR, Gradient Boosting (GBR) also shows a weak performance ( $R^2 = 0.59$ , RMSE = 0.29 eV), likely due to underfitting. Of course, all the ML models outperform a simple linear BEP relationship (Figure 3), due to the capture of any non-linear trends.

These results show that ML models do not perform well on test data as the highest  $R^2$  score achieved is 79% for RR model with a test RMSE of 0.21 eV, which does not breach the DFT accuracy, i.e.,  $\pm 0.20$  eV. Although the ML models are not good enough for descriptor-based predictions, they could be used for calculating the feature importance. Therefore, the best ML model, RR, was selected based on the lowest test RMSE (0.21 eV). While the scikit-learn-based machine learning models predict the activation barrier with reasonable accuracy, there are several limitations. First, these models typically do not provide an explicit mathematical expression for the prediction model, as they are often based on decision trees and hence are often referred to as *blackbox* models. Second, the analysis does not incorporate dimensional considerations, making the resulting models purely mathematical rather than physically meaningful. Finally, the limited size of the dataset weakens the predictive performance, as such ML models generally require large amounts of data to achieve robustness. Therefore, the SHAP analysis was performed to understand the output of the best ML model, RR, and to calculate the feature importance of the primary descriptors.

## Feature Importance – SHAP analysis

Next we applied SHAP analysis – as indicated above: developed for the biomedical field, but now applied to catalysis – to gain insights into the nature of the active site by understanding the synergistic interactions of different electronic, energy, and geometric features to render catalytic activity.<sup>27–31</sup> The SHAP summary plots (Figure 4) generated from the RR model provide a detailed view of how different descriptors contribute to the model's predictions,

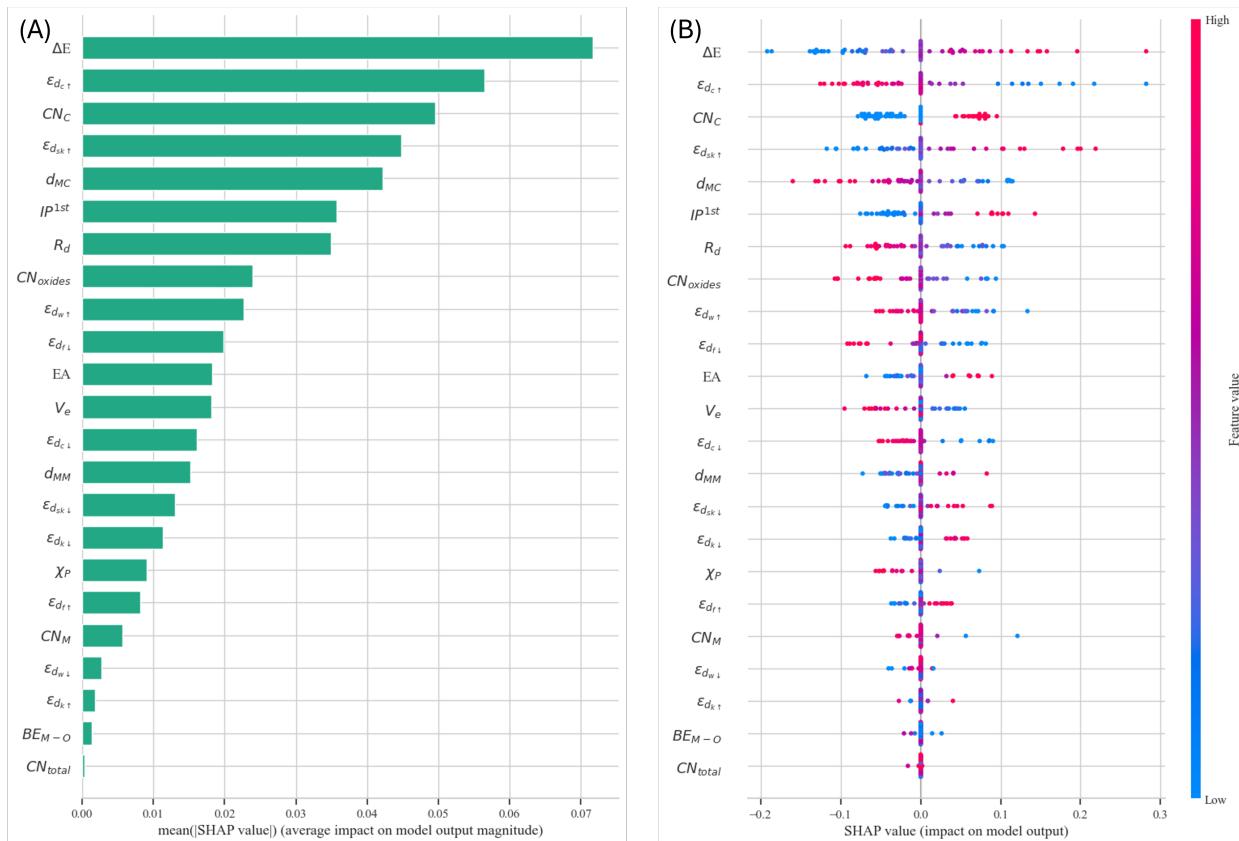


Figure 4: SHAP analysis. (A) Feature importance based on the RR model. (B) Spread of the feature values and its influence on the activation barrier prediction.

i.e., C-OH activation barrier ( $E_a$ ). Each feature is represented on the  $y$ -axis and is ranked by its average absolute SHAP value, which reflects its importance in the model (Figure 4A). In Figure 4B, the  $x$ -axis shows the SHAP values for each feature across the dataset, indicating whether a particular feature increases or decreases the model's output  $E_a$  for a given sample. Each point represents a single observation, and the color corresponds to the actual value of the feature — red for high, blue for low.

The SHAP analysis highlights the relative importance of different descriptors in predicting the C-OH activation energy ( $E_a$ ). Among all descriptors, the most impactful features (with mean  $|SHAP| < 0.04$ ) are the reaction energy ( $\Delta E$ ), the up-spin  $d$ -band center ( $\varepsilon_{d_{c\uparrow}}$ ), the dopant's coordination number with nearby Carbon atoms ( $CN_C$ ), the up-spin  $d$ -band

skewness ( $\varepsilon_{d_{sk\uparrow}}$ ), and the closest dopant-carbon distance ( $d_{MC}$ ). These descriptors show a strong positive contribution to the model's predictions. High values of  $\Delta E$ ,  $CN_C$ , and  $\varepsilon_{d_{sk\uparrow}}$  (red points) increase  $E_a$ , whereas low values (blue points) reduce  $E_a$ , confirming a strong relationship between the electronic structure and catalytic performance, whereas this trend is exactly opposite for  $\varepsilon_{d_{c\uparrow}}$  and  $d_{MC}$ . Moderately important features ( $0.03 > \text{mean } |\text{SHAP}| > 0.04$ ) include the dopant's atomic radius ( $R_d$ ) and , and first ionization potential ( $IP^{1st}$ ). For example, larger  $R_d$  values are associated with lower  $E_a$ , consistent with our earlier findings.<sup>19</sup> Similarly, lower  $CN_C$  tends to reduce  $E_a$  by creating a softer metal coordination environment. These top 7 features (Figure 4A) represent a mix of electronic and geometric effects that largely control the  $E_a$ . Additional descriptors such as dopant's CN in its most stable oxide form ( $CN_{oxides}$ ), up-spin  $d$ -band width ( $\varepsilon_{d_{w\uparrow}}$ ), down-spin  $d$ -band center ( $\varepsilon_{d_{c\downarrow}}$ ), down-spin  $d$ -band filling ( $\varepsilon_{d_{f\downarrow}}$ ), distance with the closest metal ( $d_{MM}$ ), electron affinity (EA), valence electron count ( $V_e$ ), and down-spin  $d$ -band skewness ( $\varepsilon_{d_{sk\downarrow}}$ ) also show noticeable influence, but markedly smaller ( $0.01 > \text{mean } |\text{SHAP}| > 0.03$ ). Many of these features are also directly positively or negatively correlated with the top features. Then some of the least influential descriptors include the dopant's electronegativity ( $\chi_P$ ), binding energy with oxygen ( $BE_{M-O}$ ), CN with nearby metals ( $CN_M$ ), total coordination number ( $CN_{total}$ ), and some  $d$ -band properties (e.g.,  $\varepsilon_{d_{k\uparrow}}$ ,  $\varepsilon_{d_{k\downarrow}}$ ,  $\varepsilon_{d_{f\uparrow}}$ ,  $\varepsilon_{d_{w\downarrow}}$ ). These features contribute minimally, as their SHAP values cluster around zero. These features are likely redundant or correlated with stronger descriptors, and are candidates for feature reduction.

In conclusion, the SHAP analysis not only ranked the importance of various molecular features, but also revealed how specific values of each descriptor alter the prediction of  $E_a$  and effectively highlighted the most sensitive and important features in our dataset. Furthermore, SHAP analysis demonstrates that both electronic descriptors ( $\varepsilon_{d_{c\uparrow}}$ ,  $\varepsilon_{d_{sk\uparrow}}$ ,  $IP^{1st}$ ) and geometric descriptors ( $CN_C$ ,  $d_{MC}$ ,  $R_d$ ), combined with energy descriptor  $\Delta E$ , significantly influence  $E_a$  predictions, with  $\Delta E$  and  $\varepsilon_{d_{c\uparrow}}$  emerging as the dominant factors. It

is important to note that while SHAP analysis provides a quantitative evaluation, it provides little under-standing and does not account for the synergistic effects between different primary features. Therefore, in the following section we use SISSO to obtain the complex interpretable descriptor of activity.

## Prediction using SISSO method

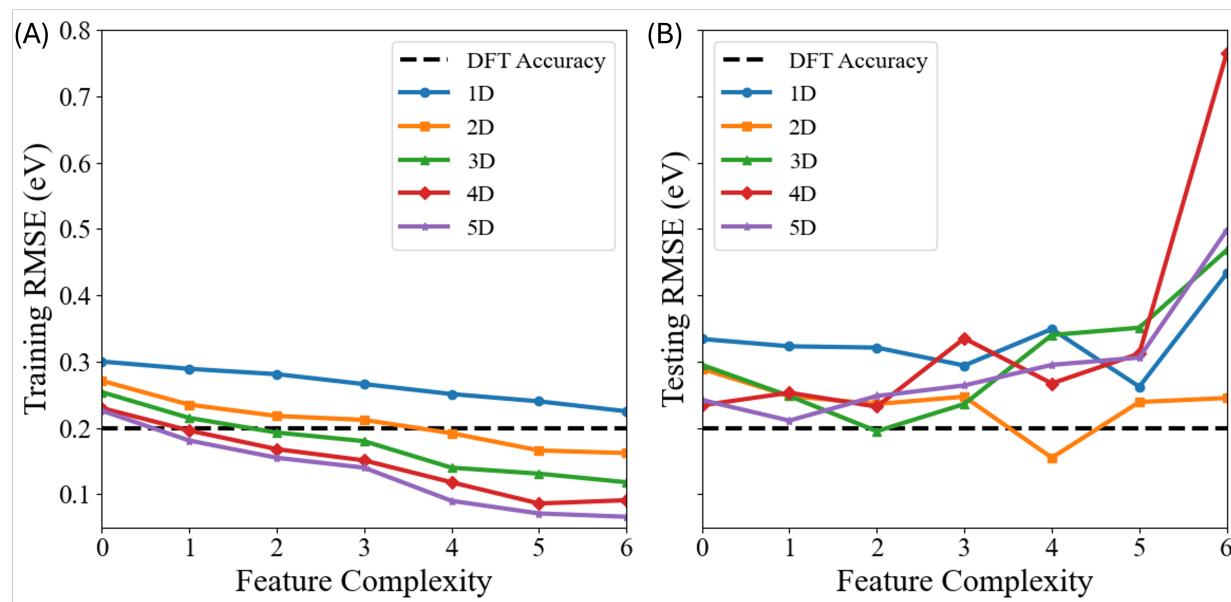


Figure 5: Training and testing RMSE using SISSO as a function of feature complexity per dimension.

The SISSO method is employed to capture the synergistic effects of electronic and geometric features, and to develop a quantitative model for predicting the C-OH activation barriers. SISSO also offers dimensional analysis functionality, ensuring that the engineered features are physically interpretable. The results obtained using SISSO are presented in Figure 5. Figure 5(A) shows that a SISSO model becomes better as a function of dimensions. Naturally and fundamentally, the more terms (primary features) there are in the linear model, the better the model can capture the behavior of the catalyst. Similarly, each model improves as a function of feature complexity, as the interaction between different electronic and geometric parameters is accounted for. In these models, values as low as 0.07 eV of

RMSE were obtained on the training set, which is way below the DFT-level error margin of  $\pm 0.20$  eV. It is evident that more terms in the model and more feature complexity typically render a better model for predicting the C-OH activation barriers.

Although SISSO models perform better on seen (training) data, they do not necessarily guarantee a good prediction over unseen (test) data. This can be seen from Figure 5B, in which there is no linear/monotonous trend in test RMSE's as a function of dimension and feature complexity. In fact, in some cases, like for 4-dimension and 6-feature complexity, it yields a very poor test RMSE of 0.76 eV.

The best SISSO models, based on test RMSE, we obtain which breach the DFT-level error margin are:

Best 2D model:

$$E_a = 1.12 \cdot \left( e^{\varepsilon_{d_{sk\downarrow}} - \Delta E} - \log(R_d) \right) - 0.0004 \cdot \left( \frac{BE_{M-O}}{(\varepsilon_{d_{f\downarrow}} \cdot CN_C) - e^{-\Delta E}} \right) + 7.04 \quad (1)$$

Best 3D model:

$$E_a = -0.02 \cdot (CN_{oxides} * CN_C) * \varepsilon_{d_{c\uparrow}} + 1.64 \cdot \left( \frac{\Delta E}{\sqrt[3]{BE_{M-O}}} \right) + 4.78 \cdot \left( \frac{|CN_{oxides} - CN_M|}{\Delta E} \right) + 1.28 \quad (2)$$

Best 5D model:

$$E_a = 0.007 \cdot \frac{IP^{1st}}{d_{MC}} + 22.44 \cdot \frac{\varepsilon_{d_{sk\uparrow}}}{BE_{M-O}} + 0.56 \cdot \frac{\Delta E}{CN_C} + 0.06 \cdot \frac{CN_{oxides}}{\Delta E} - 0.002 \cdot \frac{R_d}{\Delta E} - 0.66 \quad (3)$$

It is evident from the three mathematical expressions mentioned above that any generic descriptor derived using SISSO contains contributions combinedly from the primary electronic features and geometric features. More specifically, the dopant's *d*-band features, contributing explicitly to  $\Delta E$ , and dopant's local environment (coordination number and

atomic radius) dominate these generic descriptors. Interestingly, these primary descriptors dominate the feature importance in Figure 4A as well. Together, these two analysis (SHAP and SISSO) complement each other. Therefore, from these results, it can be established that the reaction energy ( $\Delta E$ ) combined with the dopant's *d*-band features ( $\varepsilon_{d_{sk\downarrow}}$ ,  $\varepsilon_{d_{f\downarrow}}$ ) and its local environment ( $R_d$ ,  $CN_C$ ) strongly influence the activity of a given Mo<sub>2</sub>C-based catalyst.

As highlighted in the introduction, the purpose of using SISSO and SISSO-like multidimensional analysis tools is to obtain the least complex descriptor of activity that is physically interpretable as well. Hence, the 2-dimensional with 4-feature complexity descriptor obtained from SISSO is selected as the best generic descriptor. It is very important to mention here that although these results offer some physical insight, the underlying model remains challenging to interpret. Nevertheless, the obtained 2D descriptor was then used to recalculate the activation barriers for all the training and testing data, as in Figure 6A, and then compared with the popular BEP relationship (Figure 6B). It is clear that SISSO-obtained 2D descriptor performs significantly better in predicting the activation barrier ( $R^2 = 0.83$ , RMSE = 0.17 eV) compared to the BEP relationship ( $R^2 = 0.42$ , RMSE = 0.34 eV). Figure 6 also shows the influence of each value on *x*-axis on the model's prediction of  $E_a$ . This influence is quantified by Cook's distance - a statistical measure used in regression analysis to identify influential data points, or outliers, that have a large impact on the regression model's coefficients. The greater number of such larger points (or outliers) heavily influences the model's prediction, rendering a bad BEP relationship, as is the case in Figure 6B.

On the other hand, comparing the SISSO-2D results with the scikit-learn's ML models shows that the best scikit-learn model (RR) does an okay job in predicting the activation barrier ( $R^2 = 0.79$ , RMSE = 0.21 eV). However, it fails to breach the DFT-level error margin of  $\pm 0.20$  eV. Additionally, RR does not provide an explicit mathematical expression for the underlying model, unlike SISSO. Overall, it is shown in this work that despite the

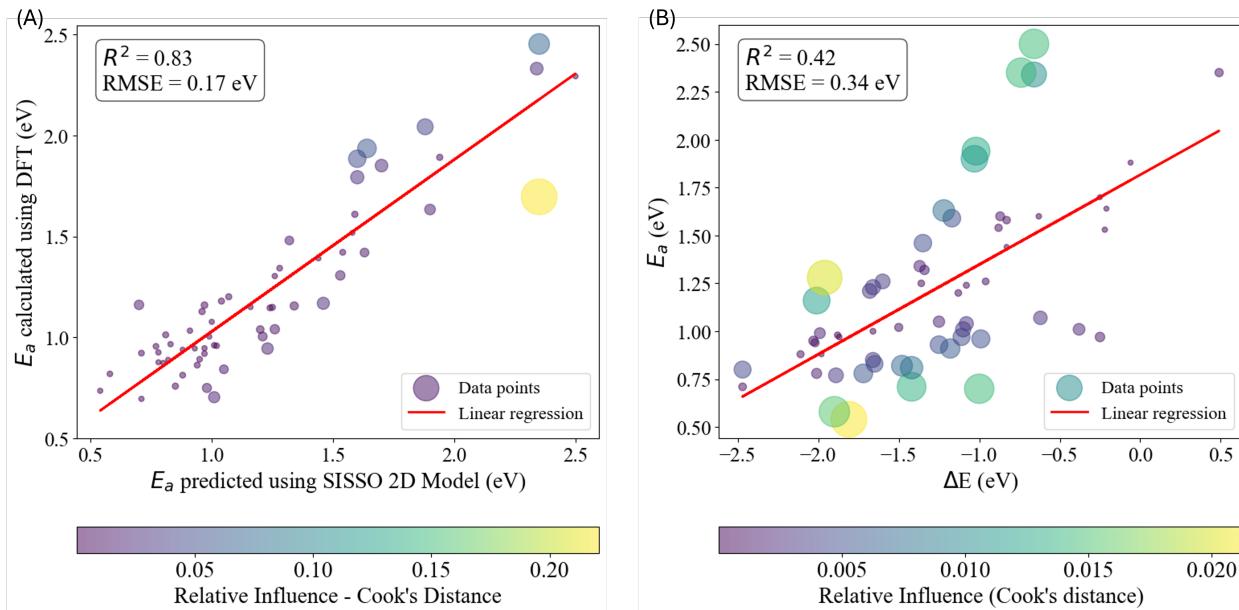


Figure 6: (A) Correlation of the SISSO-predicted  $E_a$  with the DFT-calculated  $E_a$ . (B) BEP relationship using DFT-calculated  $E_a$ .

small dataset, containing 61 samples, SISSO is able to produce a generic and, to some degree, physically interpretable descriptor for evaluating the C-OH activation barrier for Mo<sub>2</sub>C-based catalysts. Ultimately, Figure 7 highlights the most important local electronic and geometric primary features, as identified with SHAP and SISSO analyses, governing the catalytic activity of Mo<sub>2</sub>C-based catalysts for C-OH bond activation. These features essentially reflect the distinct local environment of the active site.

## Conclusions

This study demonstrates that for complex surface catalysts (Mo<sub>2</sub>C) the linear scaling BEP relationship does not hold, and that machine learning (ML) tools can be effectively used to extract physically meaningful descriptors that allow for an accurate prediction of the transition state energy. Scikit-learn's ML models generate a descriptor-based prediction of  $E_a$  that outperforms the traditional BEP relationship, but not with sufficient precision and without providing interpretable descriptors. Analysis of the output of the best ML model (RR) by

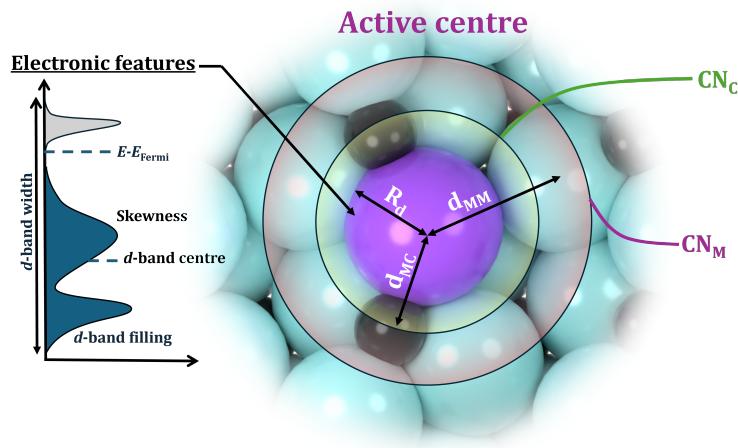


Figure 7: Scheme to highlight the majorly contributing primary descriptors controlling activity of  $\text{Mo}_2\text{C}$ -based catalysts.

SHAP reveals that the reaction energy ( $\Delta E$ ) and the dopant's  $d$ -band center contribute most in predicting the  $E_a$ . Finally, SISSO is used to obtain a low-dimensional physically interpretable descriptor, and to validate the findings from scikit-learn's ML model and SHAP analysis. It is found that a 2D descriptor containing contributions from the electronic features ( $d$ -band filling,  $d$ -band skewness, and  $\Delta E$ ) and geometric features (atomic radius and dopant's CN with C atoms) can predict the activity of  $\text{Mo}_2\text{C}$ -based catalysts ( $R^2 = 0.83$ , RMSE = 0.17 eV), which is within DFT-level error margin of  $\pm 0.20$  eV and better than any of the above. Our results predict that the local environment of active metal sites plays a key role in C-OH bond activation. Meeta-stable 101 and 010 surfaces are the most active, and surface heteroatom doping with Zr and Nb is a promising strategy to improve the performance of  $\text{Mo}_2\text{C}$ -based catalysts. This descriptor-activity relationship needs to be further validated for other transition metal carbide-based system for its generalized applicability.

# Computational Details

## DFT

All DFT calculations have been performed using the Vienna Ab Initio Simulation Package (VASP) and the Perdew-Burke-Ernzerhof functional as implemented in there,<sup>32,33</sup> with Grimme's DFT-D3BJ dispersion corrections.<sup>34,35</sup> The kinetic energy cut-off of the plane wave basis set was set to 500 eV. The convergence criterion for energy calculation and structure relaxation was set to a self-consistent field threshold of  $10^{-5}$  eV, and a maximum force threshold of 0.05 eV/Å.  $\Gamma$ -centered k-meshes of the size of  $6 \times 6 \times 6$  and  $2 \times 2 \times 1$  were used for sampling the Brillouin zone in the case of bulk and slab models, respectively. Gaussian-type smearing with a width of 0.05 eV was applied for the electronic energy density of states. For identifying the transition states, the climbing-image nudged elastic band (CI-NEB) method was used, and frequency analysis is done on the obtained transition state to confirm that there was only one imaginary frequency along the reaction coordinate.<sup>36</sup> For CI-NEB calculations, the maximum force threshold of 0.10 eV/Å was implemented. Dipole corrections were applied in the vacuum (z) direction. The bulk structure of orthorhombic Mo<sub>2</sub>C (mp-1552) was obtained from the Materials Project website and was fully relaxed. The obtained lattice parameters for Mo<sub>2</sub>C:  $a = 4.75$  Å,  $b = 5.23$  Å,  $c = 6.05$  Å (from experiments:  $a = 4.74$  Å,  $b = 5.21$  Å,  $c = 6.03$  Å),<sup>37</sup> are in good agreement with the experimentally re-reported values. From the optimized bulk(s), we cleaved the most stable 111 surface and metastable 010, 110, and 101 surfaces. Depending on the chosen facet we built slab models, deemed to be a big enough surface for the butanol C-OH activation reaction, with two or three stoichiometric layers of Mo<sub>2</sub>C. For all the slab models, a vacuum distance of 15 Å was introduced in the z-direction to minimize interaction with the periodic images. The bottom one or two stoichiometric layers, depending on the chosen facet, of the supercell were fixed to reduce the computational cost of the calculations and to mimic the bulk. Further, the metal active site in each of these facets was doped with relevant metals listed in Section 3.1. The *cif* files

for all the slab models used in this work is provided in the SI.

The adsorption energies ( $E_{\text{ads}}$ ), reaction energies ( $\Delta E$ ), and activation barriers ( $E_a$ ) were calculated as follows:

$$E_{\text{ads}} = E_{\text{slab+reactant}} - E_{\text{slab}} - E_{\text{reactant}} \quad (4)$$

$$\Delta E = E_{\text{product}} - E_{\text{reactant}} \quad (5)$$

$$E_a = E_{\text{transition state}} - E_{\text{reactant}} \quad (6)$$

Here,  $E_{\text{slab+reactant}}$  is the total energy of the slab with a reactant adsorbed on it,  $E_{\text{slab}}$  is the total energy of the clean slab,  $E_{\text{reactant}}$  and  $E_{\text{product}}$  are the total energies of the reactants and products of each elementary reaction step, and  $E_{\text{transition state}}$  is the total energy of the transition state (TS). The electronic structure parameters, such as dopant's  $d$ -band center and  $d$ -band filling, are derived from the density of states (DOS) and were calculated using Python's *pymatgen* package.<sup>38</sup>

X-band filling ( $X = s, p, d$ ) was calculated as:

$$f_x = \frac{\int_{-\infty}^{Fermi} \rho(\varepsilon)}{\int_{-\infty}^{\infty} \rho(\varepsilon)} \quad (7)$$

Here,  $\varepsilon$  means energy and  $\rho(\varepsilon)$  means the density of states.

The unoccupied  $d$ -band center was calculated as:

$$\varepsilon_{d-un} = \frac{\int_{Fermi}^{\infty} \varepsilon \rho(\varepsilon)}{\int_{-\infty}^{\infty} \rho(\varepsilon)} \quad (8)$$

## SISSO

The sure independence screening and sparsifying operator (SISSO)<sup>14</sup> in a Fortran-based code was used to efficiently extract relevant material descriptors from huge and strongly correlated feature spaces, even when only small training sets are available. The SISSO ap-

proach finds the best descriptor by combining material's features (electronic, geometric, and energy properties), identifies the most correlated features and discards the irrelevant ones, and expresses the descriptor-property relationship in the form of a mathematical function, also known as the SISSO-derived model. A generic SISSO model to predict the material's property ( $P^{SISSO}$ ) can be expressed as a linear combination of  $N$ -dimensional descriptors ( $\Phi_i$ ):

$$P^{SISSO} = \sum_{i=0}^N c_i \Phi_i \quad (9)$$

Here  $c_i$ 's are the fitting coefficients.  $\mathbf{n} \in (1,6)$ , and  $\mathbf{i} \in (1,5)$ .

While,

$$\Phi_i = \cup_{i=1}^n \hat{H}^{(m)}[\phi_1, \phi_2], \forall \phi_1, \phi_2 \in \Phi_{i-1} \quad (10)$$

Here, the  $\hat{H}$  is a set of mathematical operators considered for constructing complex features by combining primary features, e.g.  $\phi_1$  and  $\phi_2$  in eq. 10. The  $\hat{H}$  in this work contains the following operators:

$$\hat{H}^{(m)} = I, +, -, *, /, exp, log, | - |, \sqrt{\phantom{x}}^{-1}, ^2, ^3 \quad (11)$$

The superscript  $m$  indicates that when applying  $\hat{H}^{(m)}$  to primary features  $\phi_1$  and  $\phi_2$  a dimensional analysis is performed, which ensures that only physically meaningful combinations are retained, i.e., only primary features with the same unit are added or subtracted. Therefore, the complexity of a SISSO model depends on i) dimensionality: the number of linear terms in the model  $P^{SISSO}$  (eqn. 9), and ii) feature complexity: the number of operators included in  $\Phi_i$  (eqn. 10).

## Scikit-learn ML-Model's Hyperparameters

Following are the scikit-learn<sup>21</sup> ML models used in this study, and their corresponding hyperparameters as optimized using the grid method: RR (alpha = 10), XGBR (learning rate = 0.01, max depth = 7, n estimators = 500, subsample = 0.8), RFR (max depth = None,

min sample split = 2, n estimators = 100), ETR (max depth = 10, min sample split = 2, n estimators = 500), SVR (C = 100, gamma = 0.01, kernal = poly), KNN (n neighbors = 3, p = 1, weights = uniform), GBR (learning rate = 0.2, max depth = 3, n estimators = 300), and KRR (alpha = 1.0, gamma = 0.1, kernal = polynomial).

## SHAP

Built on the foundational work by Kononenko and Štrumbelj,<sup>39,40</sup> the SHAP (SHapley Additive exPlanations)<sup>22,41</sup> is a widely used tool for interpreting ML predictions. The Shapley value for a feature  $i$  represents its contribution to the prediction, calculated as:

$$\phi_i = \sum_{S \subseteq N \setminus \{i\}} \frac{|S|!(|N| - |S| - 1)!}{|N|!} [f(S \cup \{i\}) - f(S)] \quad (12)$$

where,  $\mathbf{N}$  is the set of all features,  $\mathbf{S}$  is a subset of features that does not include  $i$ ,  $f(S)$  is the model's prediction using only features in  $\mathbf{S}$ , and the fraction is a weighting term ensuring fairness across all possible feature orders. The term  $[f(S \cup \{i\}) - f(S)]$  is the marginal contribution of feature  $i$  when added to the subset  $\mathbf{S}$ . The sum averages this marginal contribution over all possible subsets of features, so no feature is favored just because it's considered earlier or later. For any single instance  $x$ , SHAP produces an additive explanation:

$$f(x) = \phi_0 + \sum_{i=1}^M \phi_i \quad (13)$$

where  $\phi_0$  is the baseline value (average model output over the dataset) and each  $\phi_i$  is the contribution of feature  $i$  for that specific prediction. By simulating the effect of removing each feature (and considering all possible combinations), SHAP assigns each feature a value representing its contribution to pushing the prediction up or down. This allows to explain both individual predictions (local explanations) and the model's general behavior across the dataset (global explanations). Finally, summary plots were obtained which include ranking

features by importance and showing the direction of their effect.

## Acknowledgement

This research was part of the Sector Plan Engineering II, funded by the Dutch Ministry of Education, Culture, and Science (OCW). The authors acknowledge the Dutch Organization for Scientific Research (NWO) for access to the Dutch national e-infrastructure (NWO-2025.019/L1, EINF-11982/L1, and EINF-7987).

## Supporting Information Available

The supporting information is available free of charge. The *cif* files for all the slab models used in this work are provided in **slabs.zip**.

## References

- (1) Bligaard, T.; Nørskov, J.; Dahl, S.; Matthiesen, J.; Christensen, C.; Sehested, J. The Brønsted–Evans–Polanyi relation and the volcano curve in heterogeneous catalysis. *Journal of Catalysis* **2004**, *224*, 206–217.
- (2) Abild-Pedersen, F.; Greeley, J.; Studt, F.; Rossmeisl, J.; Munter, T. R.; Moses, P. G.; Skúlason, E.; Bligaard, T.; Nørskov, J. K. Scaling Properties of Adsorption Energies for Hydrogen-Containing Molecules on Transition-Metal Surfaces. *Physical Review Letters* **2007**, *99*, 016105.
- (3) Nørskov, J. K.; Bligaard, T.; Rossmeisl, J.; Christensen, C. H. Towards the computational design of solid catalysts. *Nature Chemistry* **2009**, *1*, 37–46.
- (4) Michaelides, A.; Liu, Z.-P.; Zhang, C. J.; Alavi, A.; King, D. A.; Hu, P. Identification of General Linear Relationships between Activation Energies and Enthalpy Changes for

Dissociation Reactions at Surfaces. *Journal of the American Chemical Society* **2003**, *125*, 3704–3705.

- (5) Göltl, F.; Mavrikakis, M. Generalized Brønsted-Evans-Polanyi Relationships for Reactions on Metal Surfaces from Machine Learning. *ChemCatChem* **2022**, *14*.
- (6) Vojvodic, A.; Calle-Vallejo, F.; Guo, W.; Wang, S.; Toftelund, A.; Studt, F.; Martínez, J. I.; Shen, J.; Man, I. C.; Rossmeisl, J.; Bligaard, T.; Nørskov, J. K.; Abild-Pedersen, F. On the behavior of Brønsted-Evans-Polanyi relations for transition metal oxides. *The Journal of Chemical Physics* **2011**, *134*, 244509.
- (7) Lin, Y.; Ushna; Lin, Q.; Wei, C.; Wang, Y.; Huang, S.; Chen, X.; Ma, X. Machine learning descriptors for CO activation on iron-based FischerTropsch catalysts. *Journal of Catalysis* **2025**, *442*, 115921.
- (8) Calle-Vallejo, F. What we talk about when we talk about breaking scaling relations. *Applied Physics Reviews* **2024**, *11*, 021305.
- (9) Logadottir, A.; Rod, T.; Nørskov, J.; Hammer, B.; Dahl, S.; Jacobsen, C. The Brønsted–Evans–Polanyi Relation and the Volcano Plot for Ammonia Synthesis over Transition Metal Catalysts. *Journal of Catalysis* **2001**, *197*, 229–231.
- (10) Nørskov, J.; Bligaard, T.; Logadottir, A.; Bahn, S.; Hansen, L.; Bollinger, M.; Ben-gaard, H.; Hammer, B.; Sljivancanin, Z.; Mavrikakis, M.; Xu, Y.; Dahl, S.; Jacobsen, C. Universality in Heterogeneous Catalysis. *Journal of Catalysis* **2002**, *209*, 275–278.
- (11) Xin, H.; Mou, T.; Pillai, H. S.; Wang, S.-H.; Huang, Y. Interpretable Machine Learning for Catalytic Materials Design toward Sustainability. *Accounts of Materials Research* **2024**, *5*, 22–34.
- (12) Shu, W.; Li, J.; Liu, J.-X.; Zhu, C.; Wang, T.; Feng, L.; Ouyang, R.; Li, W.-X. Structure Sensitivity of Metal Catalysts Revealed by Interpretable Machine Learning and First-

Principles Calculations. *Journal of the American Chemical Society* **2024**, *146*, 8737–8745.

- (13) Ghiringhelli, L. M.; Vybiral, J.; Ahmetcik, E.; Ouyang, R.; Levchenko, S. V.; Draxl, C.; Scheffler, M. Learning physical descriptors for materials science by compressed sensing. *New Journal of Physics* **2017**, *19*, 023017.
- (14) Ouyang, R.; Curtarolo, S.; Ahmetcik, E.; Scheffler, M.; Ghiringhelli, L. M. SISSO: A compressed-sensing method for identifying the best low-dimensional descriptor in an immensity of offered candidates. *Physical Review Materials* **2018**, *2*, 083802.
- (15) Wang, T.; Hu, J.; Ouyang, R.; Wang, Y.; Huang, Y.; Hu, S.; Li, W.-X. Nature of metal-support interaction for metal catalysts on oxide supports. *Science* **2024**, *386*, 915–920.
- (16) Xu, W.; Andersen, M.; Reuter, K. Data-Driven Descriptor Engineering and Refined Scaling Relations for Predicting Transition Metal Oxide Reactivity. *ACS Catalysis* **2021**, *11*, 734–742.
- (17) Shi, Y.; Yang, Y.; Li, Y.-W.; Jiao, H. Theoretical study about Mo<sub>2</sub>C(101)-catalyzed hydrodeoxygenation of butyric acid to butane for biomass conversion. *Catalysis Science & Technology* **2016**, *6*, 4923–4936.
- (18) Hyeong, S.; Moon, B. J.; Lee, A.; Im, M. J.; Yang, H. Y.; Choi, J.; Kim, S.; Moon, J.; Park, S.; Jang, S. K.; Kim, T.; Lee, J.; Bae, S.; Lee, S. Artificial Modulation of the Hydrogen Evolution Reaction Kinetics via Control of Grain Boundaries Density in Mo<sub>2</sub>C Through Laser Processing. *Advanced Functional Materials* **2025**, *35*.
- (19) Meena, R.; Bitter, J. H.; Zuilhof, H.; Li, G. Toward the Rational Design of More Efficient Mo<sub>2</sub>C Catalysts for Hydrodeoxygenation—Mechanism and Descriptor Identification. *ACS Catalysis* **2023**, *13*, 13446–13455.

- (20) Meena, R.; Draijer, K. M.; Dam, B. v.; Zuilhof, H.; Bitter, J. H.; Li, G. Rationalizing Catalytic Performances of Mo/W-(Oxy)Carbides for Hydrodeoxygenation Reaction. *ChemCatChem* **2025**,
- (21) Bisong, E. Building Machine Learning and Deep Learning Models on Google Cloud Platform, A Comprehensive Guide for Beginners. **2019**, 215–229.
- (22) Lundberg, S. M.; Nair, B.; Vavilala, M. S.; Horibe, M.; Eisses, M. J.; Adams, T.; Liston, D. E.; Low, D. K.-W.; Newman, S.-F.; Kim, J.; Lee, S.-I. Explainable machine-learning predictions for the prevention of hypoxaemia during surgery. *Nature Biomedical Engineering* **2018**, *2*, 749–760.
- (23) Lundberg, S. M.; Erion, G.; Chen, H.; DeGrave, A.; Prutkin, J. M.; Nair, B.; Katz, R.; Himmelfarb, J.; Bansal, N.; Lee, S.-I. From local explanations to global understanding with explainable AI for trees. *Nature Machine Intelligence* **2020**, *2*, 56–67.
- (24) Wang, T.; Luo, Q.; Li, Y.-W.; Wang, J.; Beller, M.; Jiao, H. Stable surface terminations of orthorhombic Mo<sub>2</sub>C catalysts and their CO activation mechanisms. *Applied Catalysis A: General* **2014**, *478*, 146–156.
- (25) Tacey, S. A.; Jankousky, M.; Farberow, C. A. Assessing the role of surface carbon on the surface stability and reactivity of -Mo<sub>2</sub>C catalysts. *Applied Surface Science* **2022**, *593*, 153415.
- (26) Marill, K. A. Advanced Statistics: Linear Regression, Part II: Multiple Linear Regression. *Academic Emergency Medicine* **2004**, *11*, 94–102.
- (27) Mine, S.; Takao, M.; Yamaguchi, T.; Toyao, T.; Maeno, Z.; Siddiki, S. M. A. H.; Takakusagi, S.; Shimizu, K.; Takigawa, I. Analysis of Updated Literature Data up to 2019 on the Oxidative Coupling of Methane Using an Extrapolative Machine-Learning Method to Identify Novel Catalysts. *ChemCatChem* **2021**, *13*, 3636–3655.

- (28) Li, Z. et al. Machine learning-assisted Ru-N bond regulation for ammonia synthesis. *Nature Communications* **2025**, *16*, 7818.
- (29) Vidal-Lopez, A.; Mahringer, J.; Comas-Vives, A. Key Descriptors of Single-Atom Catalysts Supported on MXenes (Mo<sub>2</sub>C, Ti<sub>2</sub>C) Determining CO<sub>2</sub> Activation. *The Journal of Physical Chemistry C* **2025**, *129*, 8556–8569.
- (30) Praveen, C. S.; Comas-Vives, A. Design of an Accurate Machine Learning Algorithm to Predict the Binding Energies of Several Adsorbates on Multiple Sites of Metal Surfaces. *ChemCatChem* **2020**, *12*, 4611–4617.
- (31) Zhao, S.; Mine, S.; Wang, G.; Zhang, W.; Fakir, A. A. E.; Yang, B.; Qin, Z.; Dostagir, N. H. M.; Matsushita, K.; Takigawa, I.; Shimizu, K.-i.; Toyao, T. Development of highly active catalysts for low-temperature CO<sub>2</sub> hydrogenation to methanol using machine learning approach. **2025**,
- (32) Kresse, G.; Furthmüller, J. Efficiency of ab-initio total energy calculations for metals and semiconductors using a plane-wave basis set. *Computational Materials Science* **1996**, *6*, 15–50.
- (33) Kresse, G.; Furthmüller, J. Efficient iterative schemes for ab initio total-energy calculations using a plane-wave basis set. *Physical Review B* **1996**, *54*, 11169–11186.
- (34) Grimme, S.; Antony, J.; Ehrlich, S.; Krieg, H. A consistent and accurate ab initio parametrization of density functional dispersion correction (DFT-D) for the 94 elements H-Pu. *The Journal of Chemical Physics* **2010**, *132*, 154104.
- (35) Grimme, S.; Ehrlich, S.; Goerigk, L. Effect of the damping function in dispersion corrected density functional theory. *Journal of Computational Chemistry* **2011**, *32*, 1456–1465.

- (36) Henkelman, G.; Uberuaga, B. P.; Jónsson, H. A climbing image nudged elastic band method for finding saddle points and minimum energy paths. *The Journal of Chemical Physics* **2000**, *113*, 9901–9904.
- (37) Haines, J.; Léger, J. M.; Chateau, C.; E Lowther, J. Experimental and theoretical investigation of Mo<sub>2</sub>C at high pressure. *Journal of Physics: Condensed Matter* **2001**, *13*, 2447.
- (38) Ong, S. P.; Richards, W. D.; Jain, A.; Hautier, G.; Kocher, M.; Cholia, S.; Gunter, D.; Chevrier, V. L.; Persson, K. A.; Ceder, G. Python Materials Genomics ( pymatgen): A robust, open-source python library for materials analysis. *Computational Materials Science* **2013**, *68*, 314–319.
- (39) Štrumbelj, E.; Kononenko, I. Adaptive and Natural Computing Algorithms, 10th International Conference, ICANNGA 2011, Ljubljana, Slovenia, April 14-16, 2011, Proceedings, Part II. **2011**, 21–30.
- (40) Štrumbelj, E.; Kononenko, I. Explaining prediction models and individual predictions with feature contributions. *Knowledge and Information Systems* **2014**, *41*, 647–665.
- (41) Lundberg, S.; Lee, S.-I. A Unified Approach to Interpreting Model Predictions. *arXiv* **2017**,