## **Supporting information**

# Elucidating the impact of oxygen functional groups on the catalytic activity of M-N<sub>4</sub>-C catalysts for oxygen reduction reaction: A density functional theory and machine learning approach

Liang Xie<sup>b</sup>, Wei Zhou<sup>a, b\*</sup>, Yuming Huang<sup>b</sup>, Zhibin Qu<sup>b</sup>, Longhao Li<sup>b</sup>, Chaowei Yang<sup>b</sup>, Yani Ding<sup>b</sup>, Junfeng Li<sup>b</sup>, Xiaoxiao Meng<sup>b</sup>, Fei Sun<sup>b</sup>, Jihui Gao<sup>b\*</sup>, Guangbo Zhao<sup>b</sup>, Yukun Qin<sup>b</sup>.

a State Key Laboratory of Urban Water Resource and Environment, Harbin Institute of Technology, Harbin, Heilongjiang, 150001 P. R. China

b School of Energy Science and Engineering, Harbin Institute of Technology, Harbin, Heilongjiang, 150001 P. R. China

\* Corresponding author:

E-mail addresses: hitzhouw@hit.edu.cn (Wei Zhou)

E-mail addresses: gaojh@hit.edu.cn (Jihui Gao)

#### 1. Method for calculating the electronegativity of OGs

In order to investigate the effect of OGs on  $MN_4$ , we consider the OGs as a whole and calculate their average electronegativity ( $\chi_{OGs}$ ). The formula for calculating by harmonic average method is as follows:

$$\chi_{OGS} = \frac{N_{OGS}}{\frac{N_C}{\chi_C} + \frac{N_H}{\chi_H} + \frac{N_O}{\chi_O}}$$

Where  $N_{OGs}$  represents the total number of atoms in OGs,  $N_C$ ,  $N_H$  and  $N_O$  respectively represent the number of atoms C, H and O in OGs,  $\chi_C$ ,  $\chi_H$  and  $\chi_O$  respectively represent the Pauli electronegativity of C, H and O.

#### 2. Machine learning model used in this work

The ML models used for training and prediction include RFR, KNN, SVR, NN, XGBR, LINER, GBR, GPR, and LASSO. Here is a brief introduction to them:

- RFR (Random Forest Regressor): RFR is an integrated learning model based on decision trees, which improves prediction accuracy by constructing multiple decision trees and taking their average values. Applicable to regression problems. (Type: integrated learning model)
- (2) KNN (K-Nearest Neighbors): Description: KNN performs classification or regression based on the nearest neighbor relationship of the sample. For a given sample, predictions are made by looking at its nearest k neighbors, either by majority voting or by average. (Type: Supervised learning model)
- (3) SVR (Support Vector Regressor):SVR is an application of support vector machines that focuses on regression problems. By finding the support vector in the data, an optimal hyperplane is constructed to minimize the prediction error. (Type: Support vector machine model)
- (4) NN (Neural Network): NN is a model that simulates the structure of human brain neural network. Composed of input layers, hidden layers, and output layers, it learns weights to achieve complex non-linear relationships and is suitable for a variety of tasks, including classification and regression. (Type: Deep learning model)
- (5) XGBR (XGBoost Regressor): XGBR is a gradient lifting algorithm that improves performance by integrating multiple decision trees. It is excellent at handling structured data and regression problems. (Type: integrated learning model)
- (6) LINER (Linear Regression): LINER is a simple but powerful linear model that performs regression by fitting linear relationships in the data. Suitable for problems where linear relationships are obvious. (Type: Linear regression model)
- (7) GBR (Gradient Boosting Regressor): GBR is a gradient lifting algorithm that improves prediction performance by iteratively training a weak model and correcting the errors of the previous model. Applicable to regression problems. (Type: integrated learning model)
- (8) GPR (Gaussian Process Regressor): GPR is based on Bayesian inference, which treats predictions as Gaussian distributions over the underlying functions. It is suitable for small sample regression problems and provides uncertainty estimation. (Type: Gaussian process model)
- (9) LASSO (Least Absolute Shrinkage and Selection Operator): LASSO is a linear regression model that achieves feature selection by adding L1 regularization to the coefficients. It is suitable for regression problems in high dimensional data sets. (Type: Linear regression model)

### **3.Supplementary Figures**



Figure S1. In OGs@MN<sub>4</sub> catalyst, oxygen functional groups (OH, COOH, CHO, COC, C-O-C, C=O, etc.) have electron-rich oxygen atom(s).



Figure S2. Optimized structures of single doped OGs@CoN<sub>4</sub>. (a) CoN<sub>4</sub>, (b) COC@CoN<sub>4</sub>, (c) C-O-C@CoN<sub>4</sub>, (d) C=O@CoN<sub>4</sub>, (e) OH@CoN<sub>4</sub>, (f) CHO@CoN<sub>4</sub> and (g) COOH@CoN<sub>4</sub>.



Figure S3. Optimized structures of single doped OGs@FeN<sub>4</sub>. (a) FeN<sub>4</sub>, (b) COC@FeN<sub>4</sub>, (c) C-O-C@FeN<sub>4</sub>, (d) C=O@FeN<sub>4</sub>, (e) OH@FeN<sub>4</sub>, (f) CHO@FeN<sub>4</sub> and (g) COOH@FeN<sub>4</sub>.



Figure S4. Four configurations of 2OH@CoN\_4 are considered in this work.



Figure S5. Optimized structures of double doped OGs@MN<sub>4</sub>. (a)  $2C=O@CoN_4$ , (b)  $2CHO@CoN_4$ , (c)  $2COC@CoN_4$ , (d)  $2C-O-C@CoN_4$ , (e)  $2COOH@CoN_4$ , (f)  $2OH@CoN_4$ , (g)  $2C=O@FeN_4$ , (h)  $2CHO@FeN_4$ , (i)  $2COC@FeN_4$ , (j)  $2C-O-C@FeN_4$ , (k)  $2COOH@FeN_4$ , (l)  $2OH@FeN_4$ ,



Figure S6. Optimized structures of double doped OGs@MN<sub>4</sub>. (a) COOH+OH@MN<sub>4</sub>, (b) COOH+C-O-C@MN<sub>4</sub>, (c) COOH+COC@MN<sub>4</sub>, (d) COOH+CHO@MN<sub>4</sub>, (e) COC+COOH@MN<sub>4</sub>, (f) COC+CHO@MN<sub>4</sub>, (g) CHO+COC@MN<sub>4</sub>, (h) C=O+COOH@MN<sub>4</sub>.



Figure S7. Optimized structures of multiple oxygen doped OGs@MN<sub>4</sub>. (a) OH+COOH+C-O-C@CoN<sub>4</sub>, (b) OH+COOH+2C-O-C@CoN<sub>4</sub>, (c) OH+COOH+2C-O-C+CHO@CoN<sub>4</sub>, (d) 2OH+2CHO@CoN<sub>4</sub>, (e) 2OH+2CHO+COOH@CoN<sub>4</sub>, (f) 2OH+2CHO+2COOH@CoN<sub>4</sub>, (g) 2OH+2CHO+2COOH+C-O-C@CoN<sub>4</sub>, (h) 2OH+2CHO+2COOH+2C-O-C@CoN<sub>4</sub>, (i) OH+COOH +C-O-C@FeN<sub>4</sub>, (j) OH+COOH+2C-O-C@FeN<sub>4</sub>, (k) OH+COOH+2C-O-C+CHO@FeN<sub>4</sub>, (l) 2OH+2CHO@FeN<sub>4</sub>, (m) 2OH+2CHO+COOH@FeN<sub>4</sub>, (n) 2OH+2CHO+2COOH@FeN<sub>4</sub>, (o) 2OH+2CHO+2COOH+C-O-C@FeN<sub>4</sub>, (p) 2OH+2CHO+2COOH+2C-O-C@FeN<sub>4</sub>, (b) 2OH+2CHO+2COOH@FeN<sub>4</sub>, (c) 2OH+2CHO+2COOH=2C-O-C@FeN<sub>4</sub>, (c) 2OH+2CHO+2C-O-C@FeN<sub>4</sub>, (c) 2OH+2CHO+2C-O-C@FeN<sub>4</sub>, (c) 2OH=2CHO+2C-O-C@FeN<sub>4</sub>, (



Figure S8. The variations of (a) temperature and (b) energy versus the AIMD simulation time for 2 ps of OGs@CoN<sub>4</sub> under 298.15 K.



Figure S9. The variations of (a) temperature and (b) energy versus the AIMD simulation time for 2 ps of OGs@CoN<sub>4</sub> under 1298.15 K.



Figure S10. The variations of (a) temperature and (b) energy versus the AIMD simulation time for 2 ps of OGs@FeN<sub>4</sub> under 298.15 K.



ure S11. The variations of (a) temperature and (b) energy versus the AIMD simulation time for 2 ps of OGs@FeN<sub>4</sub> under 1298.15 K.



Figure S12. The free energy diagram of 4eORR and 2eORR pathway on OGs@CoN<sub>4</sub> and OGs@FeN<sub>4</sub> in implicit water solvent.



Figure S13. The overpotential and free energy diagram of 2eORR pathway on double doped  $OGs@CoN_4$  and  $OGs@FeN_4$  in implicit water solvent.



Figure S14. The free energy diagram of 2eORR pathway on double doped OGs@CoN<sub>4</sub> and OGs@FeN<sub>4</sub> in implicit water solvent.



Figure S15. Linear correlation between  $\triangle G_{*OOH\_vac}$  and  $\triangle G_{*OH\_vac}$ .



Figure S16. Charge distribution of single doped OGs@CoN<sub>4</sub>. (a) CoN<sub>4</sub>, (b) COC@CoN<sub>4</sub>, (c) C-O-C@CoN<sub>4</sub>, (d) C=O@CoN<sub>4</sub>, (e) OH@CoN<sub>4</sub>, (f) CHO@CoN<sub>4</sub> and (g) COOH@CoN<sub>4</sub>.



Figure S17. Charge distribution of single doped OGs@FeN<sub>4</sub>. (a) FeN<sub>4</sub>, (b) COC@FeN<sub>4</sub>, (c) C-O-C@FeN<sub>4</sub>, (d) C=O@FeN<sub>4</sub>, (e) OH@FeN<sub>4</sub>, (f) CHO@FeN<sub>4</sub> and (g) COOH@FeN<sub>4</sub>.



Figure S18. D band center of single doped OGs@CoN<sub>4</sub>. (a) CoN<sub>4</sub>, (b) COC@CoN<sub>4</sub>, (c) C-O-C@CoN<sub>4</sub>, (d) C=O@CoN<sub>4</sub>, (e) OH@CoN<sub>4</sub>, (f) CHO@CoN<sub>4</sub> and (g) COOH@CoN<sub>4</sub>.



Figure S19. D band center of single doped OGs@FeN<sub>4</sub>. (a) FeN<sub>4</sub>, (b) COC@FeN<sub>4</sub>, (c) C-O-C@FeN<sub>4</sub>, (d) C=O@FeN<sub>4</sub>, (e) OH@FeN<sub>4</sub>, (f) CHO@FeN<sub>4</sub> and (g) COOH@FeN<sub>4</sub>.



Figure S20. Charge distribution of double doped OGs@CoN<sub>4</sub>. (a) 2C=O@ CoN<sub>4</sub>, (b) 2CHO@ CoN<sub>4</sub>, (c) 2COC@ CoN<sub>4</sub>, (d) 2C-O-C@ CoN<sub>4</sub>, (e) 2COOH@ CoN<sub>4</sub>, and (f) 2OH@ CoN<sub>4</sub>.



 $\begin{array}{l} \mbox{Figure S21. Charge distribution of double doped OGs@FeN_4. (a) 2C=O@ FeN_4, (b) 2CHO@ FeN_4, (c) 2COC@ FeN_4, (d) 2C-O-C@ FeN_4, (e) 2COOH@ FeN_4, and (f) 2OH@ FeN_4. \end{array}$ 



Figure S22. Charge of (a) Co in CoN<sub>4</sub>, single doped OGs@CoN<sub>4</sub>, double doped OGs@CoN<sub>4</sub>, and (b) Fe in FeN<sub>4</sub>, single doped OGs@FeN<sub>4</sub>, double doped OGs@FeN<sub>4</sub>.



Figure S23. The conventional charge based method and the  $E_g$  based method proposed in this study were used to predict \*OOH adsorption free energy.



Figure S24. The correlation between the calculated OOH adsorption free energy and the fundamental gap under implicit solvent conditions.