| 1 | Supplementary Information |
|--------|--|
| 2 | |
| 3 | An accurate and interpretable deep learning model for yield |
| 4 | prediction using hybrid molecular representations |
| 5 | Yingying Wang ^{abc} , Xinyi Sun ^{abc} , Yuanyuan Li* ^{abc} , Li Wang* ^{abc} and Jinglai Zhang* ^{abc} |
| 6 7 | ^a Henan Key Laboratory of Protection and Safety Energy Storage of Light Metal Materials, Henan University, Kaifeng, Henan 475004, China: |
| 8 | ^b Henan Province Engineering Research Center of Green Anticorrosion Technology for Magnesium Alloys |
| 9 | Henan University, Kaifeng, Henan 475004, China; |
| 10 | ^c College of Chemistry and Molecular Sciences, Henan University, Kaifeng, Henan 475004, China; |
| 11 | |

12 Feature selection process

Prior to model construction, the dataset's molecular fingerprints undergo initial processing to reduce the dimension of the descriptors. A large number of redundant features with identical numerical distributions are removed from the molecular fingerprints, as these features do not contribute meaningfully to machine learning training models. As a result, the features of these descriptors have been considerably simplified.

18

19 Model input layers

In this study, three widely used molecular fingerprints, Molecular Access System (MACCS), Extended-Connectivity Fingerprint (ECFP), and Functional Connectivity Fingerprint (FCFP), are utilized as descriptors. The distinct molecular fingerprints are utilized as separate input layers rather than being integrated into a unified input layer, primarily for

^{*}Corresponding authors

E-mail addresses: yylichem@henu.edu.cn_(Y. Li), chemwangl@henu.edu.cn (L. Wang), zhangjinglai@henu.edu.cn (J. Zhang).

24 the following reasons:

Feature diversity and uniqueness: Distinct molecular fingerprints capture different aspects of molecular structure information. For instance, MACCS captures functional group information within molecules; ECFP reflects the local environment and connectivity of each atom within a specified radius; FCFP, based on ECFP, disregards distinctions between similar functional groups and represents a more generalized functional molecular fingerprint. Separate input layers better preserve each fingerprint's unique characteristics and prevent potential information overlap or masking that could arise if combined into a single layer.

Model flexibility and adaptability: The separate input layer design allows the network architecture to be tailored for each type of fingerprint. For example, in this study, hidden layers with varying numbers of nodes are established for each distinct input layer. This flexibility enables the model to more effectively adapt to different types of input data, thereby enhancing the overall performance.

Avoiding information confusion: Combining all fingerprints into a unified input layer may cause information confusion. For example, high-dimensional features of certain fingerprints might overshadow critical information from other fingerprints. By separating the input layers, such mutual interference can be avoided, allowing the model to learn the characteristics of each fingerprint more distinctly.

42

43 Reaction condition optimization process

44 The reaction conditions are optimized through deep learning model. The CO_2 pressure is 45 always maintained at 0.1 MPa. The time range is set to [1 h, 2 h, 3 h, 4 h, 5 h], and the 46 temperature range is set to [50 °C, 60 °C, 70 °C, 80 °C, 90 °C]. The catalyst dosage range is 47 considered as [0.2 mol%, 0.4 mol%, 0.6 mol%, 0.8 mol%, 1.0 mol%, 1.2 mol%, 1.4 mol%]. 48 The initial reaction conditions are 4 h, 90 °C, 0.1 MPa, 0.8%. Not all combinations of 49 reaction conditions are traversed, but the control variable method from experiments is 50 employed to determine parameters one by one. With the aim of reducing reaction conditions 51 as much as possible while ensuring a predicted yield close to 0.95, the reaction conditions are 52 optimized by yield prediction through deep learning.

53

54 General applicability and potential limitations of the model

55 The deep learning approach presented in this study offers a novel and efficient alternative for predicting and designing IL catalysts for CO₂ cycloaddition reactions. While 56 57 the current study focuses on imidazolium-based and pyrazolium-based ILs, the approach can be extended to other types of molecules and reactions with appropriate modifications. The 58 use of hybrid fingerprint features allows the model to capture essential structural information 59 for predicting catalytic performance. Additionally, the model's ability to handle high-60 dimensional data and complex relationships, combined with SHAP analysis and a two-step 61 screening strategy, provides valuable insights and reduces experimental burden. 62

Despite promising results, the model's generalizability is limited by the dataset's focus on imidazolium-based and pyrazolium-based ILs. Extending the approach to other molecules or reactions will require more diverse training data. Additionally, the model's adaptability to more complex or less studied catalytic systems remains to be validated and may need further refinement. Future work should aim to expand the dataset, validate the model on additional 68 systems, and optimize computational efficiency to enhance its broader applicability.



71 Fig. S1. 27 types of imidazolium-based cation structures in the dataset.



Fig. S2. 13 types of anion structures associated with imidazolium-based ionic liquids in thedataset.



78 Fig. S3. 25 types of pyrazolium-based cation structures in the dataset.79



80

77

81 Fig. S4. 5 types of anion structures associated with pyrazolium-based ionic liquids in the

82 dataset.



Fig. S5. 14 types of reactant structures in the dataset.





85

86

Fig. S6. Number of data points for each of the 14 types of reactants in the dataset.

89

90 Table S1

91 Parameters of DNN models.

| Model | Random seed | Learning rate | Epochs |
|-------|-------------|---------------|--------|
| М | 28 | 0.006 | 310 |
| Е | 29 | 0.005 | 374 |
| F | 6 | 0.004 | 486 |
| ME | 3 | 0.007 | 715 |
| MF | 25 | 0.007 | 390 |
| EF | 2 | 0.007 | 262 |
| MEF | 3 | 0.007 | 385 |





Fig. S10. DNN model based on MF feature.



Fig. S11. DNN model based on EF feature.



Fig. S12. DNN model based on MEF feature.

Table S2

112 Parameters of traditional algorithms.

| Algorithms | М |
|------------|---|
| DT | random_state=35, max_depth=24 |
| RF | max_depth=17, random_state=16, max_features=84, n_estimators=11 |
| GBR | learning_rate=0.396, random_state=7, n_estimators=200 |
| XGB | random_state=0, max_depth=2, n_estimators=81 |
| SVM | Kernel='rbf', C=470, epsion=1e-06, gamma=0.01 |



- 115 Fig. S13. The structures of screened pyrazolium-based ionic liquids.
- 116
- 117 Table S3
- 118 The 12 screened pyrazolium-based ionic liquids.

| Strature | Time (h) | T (°C) | P (MPa) | Amount | Predicted |
|-----------|----------|--------|---------|--------|-----------|
| Structure | | | | (mol%) | yield |
| P-1 | 7 | 90 | 0.1 | 1 | 1.00 |
| P-2 | 7 | 90 | 0.1 | 1 | 0.96 |
| P-3 | 7 | 90 | 0.1 | 1 | 0.91 |
| P-4 | 7 | 90 | 0.1 | 1 | 0.87 |
| P-5 | 7 | 90 | 0.1 | 1 | 0.87 |
| P-6 | 7 | 90 | 0.1 | 1 | 0.85 |
| P-7 | 7 | 90 | 0.1 | 1 | 0.83 |
| P-8 | 7 | 90 | 0.1 | 1 | 0.83 |
| P-9 | 7 | 90 | 0.1 | 1 | 0.83 |
| P-10 | 7 | 90 | 0.1 | 1 | 0.82 |
| P-11 | 7 | 90 | 0.1 | 1 | 0.82 |
| P-12 | 7 | 90 | 0.1 | 1 | 0.80 |



- 121 Fig. S14. The mechanism of the CO₂ cycloaddition reaction catalyzed by ionic liquids.
- 122

123 Table S4

- 124 The energy barriers, optimized reaction conditions, and corresponding yields of 19
- 125 imidazole-based ionic liquids.

| II.a | Time | Т | Р | Amount | Predicted | Energy barriers |
|------|------|------|-------|--------|-----------|---------------------------|
| ILS | (h) | (°C) | (MPa) | (mol%) | yield | (kcal mol ⁻¹) |
| M-1 | 4 | 80 | 0.1 | 1.2 | 0.96 | 16.0 |
| M-2 | 4 | 80 | 0.1 | 1.0 | 0.98 | 17.2 |
| M-3 | 4 | 70 | 0.1 | 1.0 | 0.96 | 18.8 |
| M-4 | 4 | 80 | 0.1 | 1.0 | 0.98 | 19.8 |
| M-5 | 4 | 60 | 0.1 | 1.0 | 0.98 | 20.2 |
| M-6 | 4 | 80 | 0.1 | 1.0 | 0.96 | 20.3 |
| M-7 | 4 | 50 | 0.1 | 1.2 | 0.99 | 20.4 |
| M-8 | 4 | 80 | 0.1 | 1.0 | 0.95 | 20.4 |
| M-9 | 4 | 50 | 0.1 | 1.2 | 0.97 | 20.7 |
| M-10 | 4 | 90 | 0.1 | 0.8 | 0.94 | 20.8 |
| M-11 | 4 | 70 | 0.1 | 1.0 | 0.95 | 21.1 |
| M-12 | 4 | 90 | 0.1 | 1.0 | 0.95 | 21.8 |
| M-13 | 4 | 90 | 0.1 | 0.8 | 0.95 | 22.4 |
| M-14 | 4 | 60 | 0.1 | 0.8 | 0.97 | 22.6 |
| M-15 | 4 | 90 | 0.1 | 1.0 | 0.94 | 23.1 |
| M-16 | 4 | 80 | 0.1 | 1.2 | 0.95 | 23.3 |

| M-17 | 4 | 90 | 0.1 | 0.8 | 0.95 | 29.4 |
|------|---|----|-----|-----|------|------|
| M-18 | 4 | 80 | 0.1 | 1.0 | 0.95 | 34.8 |
| M-19 | 4 | 70 | 0.1 | 1.2 | 0.97 | 40.5 |

127



128 Fig. S15. The structures of 5 imidazole-based ionic liquids with energy barriers more than

129 23.0 kcal mol⁻¹.