

GHP-MOFassemble: Diffusion modeling, high throughput screening, and molecular dynamics for rational discovery of novel metal-organic frameworks for carbon capture at scale

Hyun Park^{1,2,3,†}, Xiaoli Yan^{1,4,†}, Ruijie Zhu^{1,5,†}, E. A. Huerta^{1,6,7},
Santanu Chaudhuri^{1,4}, Donny Cooper⁸, Ian Foster^{1,6}, Emad
Tajkhorshid^{2,3,9}

¹ Data Science and Learning Division, Argonne National Laboratory, Lemont, Illinois 60439, USA

² Theoretical and Computational Biophysics Group, Beckman Institute for Advanced Science and Technology, University of Illinois at Urbana-Champaign, Urbana, Illinois 61801, USA

³ Department of Biochemistry, University of Illinois at Urbana-Champaign, Urbana, Illinois 61801, USA

⁴ Multiscale Materials and Manufacturing Lab, University of Illinois Chicago, Chicago, Illinois 60607, USA

⁵ Department of Materials Science and Engineering, Northwestern University, Evanston, Illinois 60208, USA

⁶ Department of Computer Science, University of Chicago, Chicago, Illinois 60637, USA

⁷ Department of Physics, University of Illinois at Urbana-Champaign, Urbana, Illinois 61801, USA

⁸ Computational Science and Engineering, Data Science and AI Department, TotalEnergies EP Research & Technology USA, LLC, Houston, Texas 77002 USA

⁹ Center for Biophysics and Quantitative Biology, University of Illinois at Urbana-Champaign, Urbana, Illinois 61801, USA

[†] These authors contributed equally to this work.

E-mail: eli.hu@anl.gov

16 June 2023

Abstract. We introduce **GHP-MOFassemble**, a **Generative** artificial intelligence (AI), **High Performance** framework to accelerate the rational design of metal-organic frameworks (MOFs) with high CO₂ capacity and synthesizable linkers. Our framework combines a diffusion model, a class of generative AI, to generate novel linkers that are assembled with one of three pre-selected nodes into MOFs in a primitive cubic (pcu) topology. The CO₂ capacities of these AI-generated MOFs are predicted using a modified version of the crystal graph convolutional neural network model. We then use the LAMMPS code to perform molecular dynamics simulations to relax the AI-generated MOF structures, and identify those that converge to stable structures, and maintain their porous properties throughout the simulations. Among 120,000 pcu MOF candidates generated by the **GHP-MOFassemble** framework, with three distinct

metal nodes (Cu paddlewheel, Zn paddlewheel, Zn tetramer), a total of 102 structures completed molecular dynamics simulations at 1 bar with predicted CO₂ capacity higher than 2 mmol/g at 0.1 bar, which corresponds to the top 5% of hMOFs in the hypothetical MOF (hMOF) dataset in the MOFX-DB database. Among these candidates, 18 have change in density lower than 1% during molecular dynamics simulations, indicating their stability. We also found that the top five GHP-MOFassemble’s MOF structures have CO₂ capacities higher than 96.9% of hMOF structures. This new approach combines generative AI, graph modeling, large-scale molecular dynamics simulations, and extreme scale computing to open up new pathways for the accelerated discovery of novel MOF structures at scale.

Keywords: diffusion model, pre-trained regression model, metal-organic framework, carbon capture

1. Introduction

Metal-organic frameworks (MOFs) have garnered much research interest in recent years due to their diverse industrial applications, including gas adsorption and storage [1], catalysis [2], and drug delivery [3]. As nanocrystalline porous materials, MOFs are modular in nature [4], with a specific MOF defined by its three types of constituent *building blocks* (inorganic nodes, organic nodes, and organic linkers) and *topology* (the relative positions and orientations of building blocks). MOFs with different properties can be produced by varying these building blocks and their spatial arrangements. Nodes and linkers differ in terms of their numbers of connection points: inorganic nodes are typically metal oxide complexes whereas organic nodes are molecules, both with three or more connection points; organic linkers are molecules with only two connection points.

MOFs have been shown to exhibit superior chemical and physical CO₂ adsorption properties. They can be recycled in appropriate operating environments for varying numbers of times before undergoing significant structural degradation. However, their industrial applications have not yet reached full potential due to stability issues, such as poor long-term recyclability [5], and high moisture sensitivity [6]. For instance, the presence of moisture in adsorption gas impairs a MOF’s CO₂ capture performance, which may be attributed to MOFs having stronger affinity toward water molecules than to CO₂ molecules [7]. Other investigations of MOF stability issues [8, 9, 10] have shown [11, 12] that the gas adsorption properties of MOFs can be enhanced by tuning the building blocks used. Yet progress has been difficult due to the enormous chemical space of building blocks that makes exhaustive search with traditional experimental methods impractical [13].

Here we propose GHP-MOFassemble, a novel high-throughput computational framework to accelerate the discovery of MOF structures with high CO₂ capacities and synthesizable linkers. To demonstrate the utility of our framework, we apply it to the rational design of MOFs with pcu topology and three types of inorganic nodes: Cu

paddlewheel, Zn paddlewheel, and Zn tetramer. The key component of our generative framework is a diffusion model, **DiffLinker** [14], which enables efficient generation of new linker molecules from a set of molecular fragments.

Previous attempts to discover useful MOFs have proceeded via high-throughput screening of existing datasets [15, 16, 17], an approach that necessarily limits the search to known MOFs. Our **GHP-MOFassemble** framework instead probes the MOF design space by employing a molecular generative diffusion model (**DiffLinker** [14]) to generate chemically diverse and unique MOF linkers, which it assembles with pre-selected metal nodes to form novel MOF structures. It then screens those structures with a pre-trained regression model, a modified version of **CGCNN** [18], to identify high-performing MOF candidates for carbon capture.

2. Related Work

Previous efforts in the search for new MOF structures with exceptional gas adsorption properties include high-throughput screening and generative modeling. High-throughput screening approaches comprise both database search and machine learning (ML)-assisted screening.

Database search methods. These methods apply filters to identify attractive candidates in large databases of MOF structures plus calculated properties obtained by using molecular simulations. For example, Li et al. [16] search the computation-ready, experimental MOF database (CoRE DB) for MOFs with high CO₂ uptake when moisture is present, and Altintas et al. [19] search both CoRE DB and the Cambridge Structural Database non-disordered MOF subset for MOFs with high CH₄/H₂ selectivity.

ML-assisted screening. Database search methods require expensive molecular simulation calculations for every target molecule. ML-assisted screening avoids this difficulty by using a regression model trained on a smaller training dataset to predict target properties for a large number of new test structures. This approach has been widely applied to gas adsorption and separation property predictions of MOFs, including CO₂/H₂ separation [20] and CH₄ adsorption [21]. ML-assisted methods commonly make use of both geometrical features (e.g., dominant pore size, void fraction [22]) and chemical features (e.g., atom type, electronegativity [23]). A compound feature, the atomic property weighted radial distribution function (AP-RDF) [24] has been shown to improve regression model performance in finding MOF structures with high CO₂ capacity [25]. This feature is constructed by using both chemistry and local geometry of atomic sites in MOF structures. As an alternative to the use of hand-engineered machine learning features, neural networks have also been applied to MOF research. In a recent study, the Atomistic Line Graph Neural Network (**ALIGNN**), trained on the **hMOF** dataset [26], was applied to search CoRE DB for high-performing MOF candidates for carbon capture [27].

Generative modeling. In the generative modeling approach, candidate molecules are

not drawn from a database but are generated *de novo* via methods that produce novel compounds that are expected to have a desired set of chemical features. One notable example is the recently proposed Supramolecular Variational Autoencoder (SmVAE) [28], which uses a semantically constrained graph-based canonical code (RFcode) to encode MOF building blocks. The variational autoencoder framework structure allows this model to interpolate between existing MOF structures and enables isorecticular optimization of MOF structures toward higher CO₂ capacity and CO₂/N₂ selectivity.

This work. Proposed generative models include not only variational autoencoders (VAEs) but also generative adversarial networks (GANs), normalizing flows, autoregressive models, and diffusion models [29]. In this work, we adopt a diffusion model named DiffLinker to generate novel MOF linkers. Diffusion models use a probability distribution and Markovian properties to generate new data via a denoising step. First, Gaussian noise is added to the input samples to yield noisy data. A neural network is then trained to learn what noise was added. The trained network is then used to reversely transform (i.e., denoise) the noisy data back into target samples that resemble molecules from the training data distribution. This method has been widely applied to drug discovery to speed up the design of new ligands and ligand-protein docking configurations [30, 31, 32, 33].

Two broad categories of molecular representation schemes have been used in diffusion model-based ligand generators: molecular graphs and 3D coordinates. In the former case (Digress/Congress [34] is an example), atoms are represented as nodes and bonds as edges. In the latter case, the 3D atomic coordinates of molecules are generated directly, as in DiffLinker [14] and E(3) equivariant diffusion model (EDM) [35]. Given the success of diffusion models in drug discovery [32, 30, 33, 36, 37], in Section 3.1 we demonstrate how to transfer the idea to the iso-reticular design of MOF structures by varying MOF linkers while fixing node and topology. The diffusion model is used specifically for MOF linker design, as we describe in Section 3.2.

3. Methods

Our GHP-MOFassemble framework has three components:

- **Decompose.** We use a molecular fragmentation algorithm to decompose the MOF linkers found in high-performing MOFs within the hMOF dataset [26]—an open source dataset that provides, for each of 137,652 hypothetical MOF structures, its MOFid, MOFkey, geometric features, and isotherm data for six adsorption gases (CO₂, N₂, CH₄, H₂, Kr, Xe) at 0.01, 0.05, 0.1, 0.5, and 2.5 bar—covering the pressure ranges of cyclic adsorption gas separation processes in industrial applications. The adsorption properties of gases provided in this dataset were calculated using Grand Canonical Monte Carlo (GCMC) calculations [26], as described in Wilmer et al. [38].
- **Generate.** We use a diffusion model to generate new MOF linkers. We then screen the generated linkers by removing linkers with S, Br and I elements (we

call this step “element filter” in later text), and evaluate their quality using five different scores to quantify their synthesizability, validity, uniqueness, and internal diversity. Linkers that passed the element filter are then assembled with one of three pre-selected nodes into new MOF structures in the pcu topology.

- **Predict.** We use a pre-trained regression model to predict the CO₂ capacities of the newly generated MOF structures.

In the following sections we describe each of these components in detail.

3.1. *Decomposing MOF linkers into molecular fragments*

The first component of the GHP-MOFassemble framework, **Decompose**, decomposes linkers from high-performing MOF structures in the hMOF dataset into their molecular fragments. As illustrated in Figure 1, this process applies the following three steps to a specified node-topology pair, which in the figure is the most frequent node-topology pair in the hMOF dataset, Zn tetramer-pcu. Note that for the pcu topology, one linker is orientated along each of the x, y, and z directions, and thus at most three linker types are possible.

- Select:* We select high-performing MOFs with the given node-topology pair from hMOF. We define here a high-performing MOF structure as one with CO₂ capacity higher than 2 mmol/g at 0.1 bar, which corresponds to the top 5% of CO₂ capacity.
- Extract:* We extract the linkers Simplified Molecular-Input Line-Entry System (SMILES) strings for each high-performing MOF identified in step (i), eliminate those for MOFs with more than three unique linker types, and assemble those that remain into a linker dataset.
- Fragment:* We fragment the linkers produced in step (ii) to obtain their chemically relevant fragment-connection atom pairs, which we assemble into a molecular fragment dataset.

The extraction of linkers in step (ii) is straightforward because the MOFid provided by hMOF for each MOF specifies the SMILES strings of its constituent inorganic and organic building blocks, including its linkers, as well as a format signature and a topology code [39]. Together, these elements uniquely define the topology and building blocks of a given MOF structure.

Step (iii) use the Matched Molecular Pairs Algorithm (MMPA) [40] as implemented in DeLinker [41] to generate molecular fragments of a given molecule by breaking chemical bonds between atom pairs. We set the minimum number of connection atoms, the minimum fragment size, and the minimum path length to 3, 5, and 2, respectively. Moreover, we only consider the case where the fragments are at least two atoms away from each other. The chemically relevant fragment-connection atoms pairs are then used to form the molecular fragment dataset.

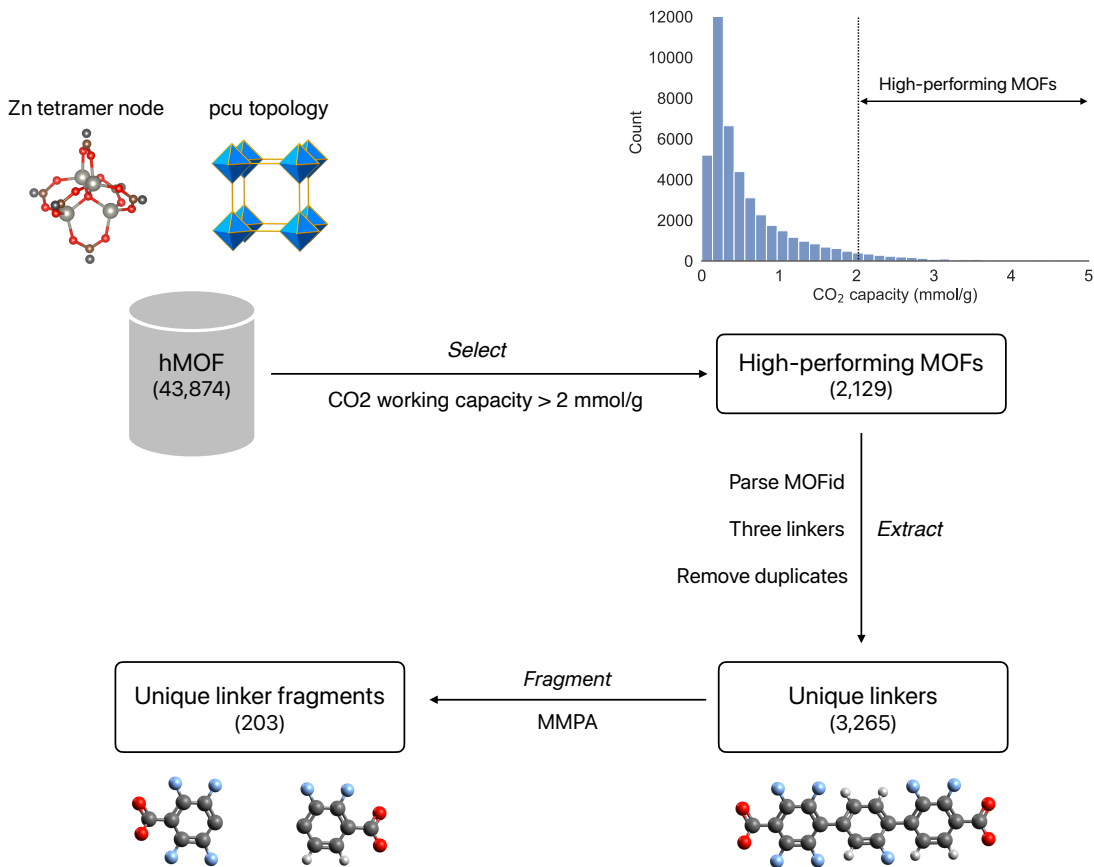


Figure 1: Schematic representation of **Decompose**, the first component of our GHP-MOFassemble framework. It consists of three steps, which we showcase for Zn tetramer-pcu MOFs. First, the *Select* step selects the Zn tetramer-pcu MOFs in hMOF that are high-performing. Then the *Extract* step extracts unique linkers from those MOFs. Finally, the *Fragment* step generates unique linker fragments. The color scheme of elements is: carbon in grey, oxygen in red, nitrogen in blue, and hydrogen in white.

3.2. Generating new MOF structures

The **Generate** component employs the pre-trained diffusion model **DiffLinker** [14] to generate new MOF linkers, and then assembles those new linkers with one of three pre-selected nodes into new MOF structures in the pcu topology. It comprises three steps: *Diffuse and Denoise*; *Screen and Evaluate*; and *Assemble*. An example of these three steps for Zn tetramer-pcu MOFs is shown in Figure 2.

3.2.1. Diffuse and Denoise. We apply the pre-trained diffusion model **DiffLinker** [14] to the molecular fragments output by the **Decompose** component to generate novel MOF linkers. This model connects the molecular fragments supplied as input (also known as context) with a specified number of sampled atoms. We vary the number of

sampled atoms from five to 10, a range chosen because we found that with fewer than five sampled atoms, straight chains or branched chains may be obtained, while with more than 10, ring structures may be present. We found that the range of five to 10 sampled atoms resulted in novel and chemically diverse MOF linkers.

DiffLinker applies a denoising process to determine the atomic species and Cartesian coordinates of sampled atoms by using a decoder neural network architecture named E(3)-Equivariant Graph Neural Network (EGNN) [42]. This graph neural network predicts zero-mean centered Cartesian coordinates noise and one-hot encoding noise of atomic species, and accounts for equivariance due to translation, rotation, and reflection of molecules, i.e., *Euclidean group 3* (E(3)). Since the pre-trained **DiffLinker** model was trained on the **GEOM** dataset [43], it enables sampling of connection atoms with high chemical and structural diversity.

DiffLinker outputs the 3D coordinates and the atomic species of heavy atoms in the linker molecules, rather than SMILES strings or 2D graphs. Since hydrogen atoms are implicit in the **DiffLinker** model, their 3D coordinates are not output by the model. To generate the spatial configurations of all-atom molecules, we employ **openbabel** to convert the **DiffLinker** generated molecules to SMILES strings, which in turn are used to generate the 3D configurations of linkers with explicit hydrogen atoms. **Openbabel** uses distance-based heuristics to determine bond connectivity in a given molecule [44], which is critical to the assignment of the number of hydrogen atoms. After the hydrogen addition step, we identify the dummy atoms which contains information about how the linkers connect with metal nodes.

The diffusion process involves two consecutive steps. The first adds Gaussian noise to the original data (i.e., \mathbf{x}), yielding noisy input (i.e., z_t at time t). Then, the denoising step applies the neural network-based noise removal operation. Mathematically, the following Markovian properties and equations are satisfied:

$$q(z_t|z_{t-1}) = \mathcal{N}(z_t; \bar{\alpha}_t z_{t-1}, \bar{\sigma}_t^2 \mathbf{I}), \quad (1)$$

$$q(z_t|\mathbf{x}) = \mathcal{N}(z_t; \alpha_t \mathbf{x}, \sigma_t^2 \mathbf{I}), \quad (2)$$

$$p(z_{t-1}|z_t) = q(z_{t-1}|\mathbf{x}, z_t), \quad (3)$$

$$q(z_{t-1}|\mathbf{x}, z_t) = \mathcal{N}(z_{t-1}; \mu_t^\theta(\mathbf{x}, z_t; \alpha_t, \sigma_t), \xi_t \mathbf{I}), \quad (4)$$

$$\mu_t^\theta(\mathbf{x}, z_t; \alpha_t, \sigma_t) = A_t z_t + B_t \mathbf{x}, \quad (5)$$

$$\hat{\mathbf{x}} = \mathbf{x} = C_t z_t - D_t \epsilon_t^\theta(z_t, t), \quad (6)$$

where $\bar{\alpha}_t = \alpha_t / (\alpha_{t-1})$, and $\bar{\sigma}_t^2 = \sigma_t^2 - \bar{\alpha}_t^2 \sigma_{t-1}^2$. \mathbf{I} is an identity matrix that computes an isotropic Gaussian upon multiplying with a constant (e.g., σ_t^2); α_t and $\bar{\alpha}_t$ are signal controls, i.e., meaningful information during training and generative steps; and σ_t and $\bar{\sigma}_t$ are noise controls, i.e., noisy information used for diversity during training and generative steps. The subscript $t = 0, 1, \dots, T$ is the time step at which a molecule is generated (a.k.a. denoised or diffused). ξ_t , A_t , B_t , C_t , and D_t are constants consisting of α_t , $\bar{\alpha}_t$, σ_t , and $\bar{\sigma}_t$.

Equations (1) and (2) represent a *diffusion* process that is conditioned on the

previous value, z_{t-1} , and initial value, \mathbf{x} , to predict a current value, z_t . The *denoising* processes are described by Equations (3) and (4), which predict a previous value conditioned on (either real or predicted) initial value and current value. Equations (5) and (6) provide a detailed evolution for μ_t^θ which represents the learned denoising mean parameterized by θ , as well as \mathbf{x} , a real initial value, approximated by $\hat{\mathbf{x}}$. The training objective is to learn a neural network, i.e., EGNN, to predict a denoising value so that random noise can be denoised to a physically and chemically valid molecule during the generative process. In practice, we predict ϵ_t^θ (learned noise value) rather than $\hat{\mathbf{x}}$ directly (i.e., Equations (5) and (6)) for better prediction [45]. Since DiffLinker generates chemically valid molecules, EGNN needs to take more regularization into account, such as E(3) and O(3) (i.e., *orthogonal group in dimension 3*: rotations and reflections). Thus, invariant features such as one-hot and equivariant features such as atomic coordinates updates need to be considered [14].

3.2.2. Screen and Evaluate. To ensure consistency of elements in hMOF linkers and generated linkers, we manually filter out linkers that contain elements not present in the hMOF dataset. Since the pre-trained DiffLinker model we used in this work was trained on the GEOM dataset, a total of nine heavy elements may be sampled, including C, N, O, F, P, S, Cl, Br, and I. Among these elements, S, Br, and I do not appear in the hMOF dataset, therefore we remove linkers that contain these three elements. For each molecular fragment, we then perform sampling 20 times. Each sampling step may yield a different set of molecules because of the random nature of the denoising process. The probabilistic nature of the diffusion model allows it to generate linkers from an extensive linker design space beyond that of the hMOF dataset.

We then use five metrics to evaluate the quality of the remaining linkers. The first two metrics are commonly used heuristic measures of synthesizability: the synthetic accessibility score (SAscore), and the synthetic complexity score (SCscore) [46].

The SAscore, as defined by Ertl and Schuffenhauer [47], is based on analysis of one million PubChem molecules, and combines fragment contributions from molecule substructures with a complexity penalty that accounts for molecular size and for structural features of molecules such as the presence of rings. The SCscore is computed by using a neural network trained on 12 million chemical reactions from the Reaxys database to estimate the number of reaction steps required to produce a molecule [46]. For both SAscore and SCscore, the higher the values are, the more difficult it is to synthesize the linker, hence less desirable.

To evaluate the capability of GHP-MOFassemble to generate a valid, novel, and diverse set of linkers, we leverage the MOSES framework [48] to compute three additional metrics: the fraction with valid SMILES strings (validity), the fraction that are unique (uniqueness), and their dissimilarity (internal diversity). We compute the fraction of linkers that are unique by comparing their canonical SMILES strings, and use their average pairwise molecular Tanimoto dissimilarity [49] as measures of diversity. We compute the diversity of the generated linker set by using two internal diversity scores

($\text{IntDiv}_1(G), \text{IntDiv}_2(G)$). We present the results of these metrics in [Section 4.2](#).

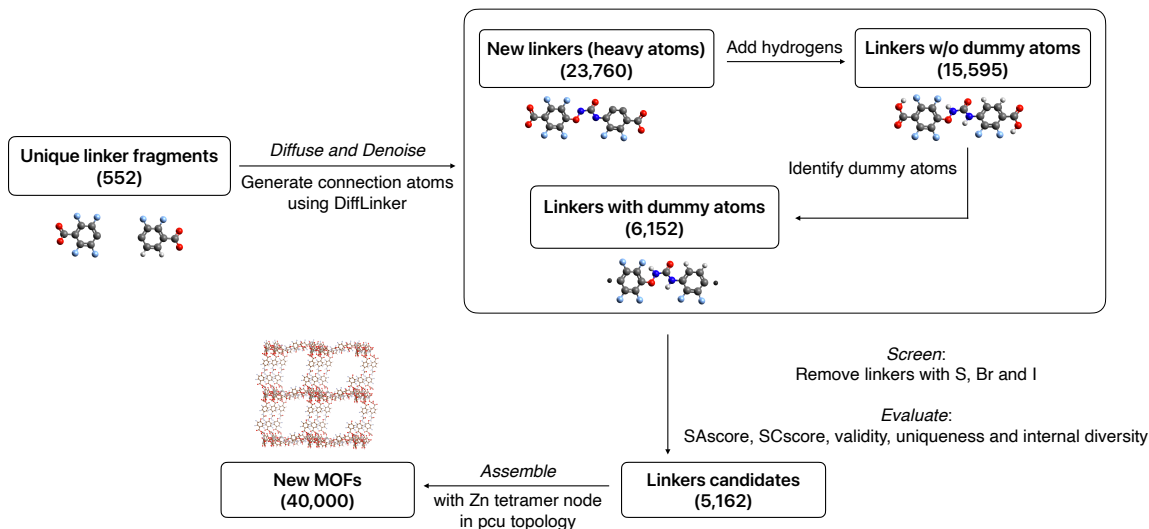


Figure 2: The **Generate** component of the GHP-MOFAssemble framework involves *Diffuse and Denoise*, *Screen*, and *Assemble* steps, as shown here for Zn tetramer-pcu MOFs. First, we generate new linkers via the *Diffuse and Denoise* step. We then add hydrogen atoms to ensure correct valency of the generated linkers. Next, we identify dummy atoms by replacing the two carbon atoms in the carboxyl groups of each generated linker with dummy atoms. Linkers with dummy atoms then undergo the *Screen* step, which removes linkers with S, Br and I. These three elements are present in the GEOM dataset but not in hMOF dataset, and thus we filter out linkers with these three elements. As a result, this step reduces the number of potential linkers to 5,162. The generated linkers’ molecular statistics are quantified using five metrics, including SCScore, SAScore, validity, uniqueness, and internal diversity. Finally, in *Assemble* step, we build 40,000 new MOFs with the Zn tetramer node in pcu topology. As before, the color code used for the atoms shown is: carbon in grey, oxygen in red, nitrogen in blue, and hydrogen in white.

3.2.3. Assemble. Once we have generated linkers that meet the five metrics mentioned above, we can assemble them by connecting the node and linker building blocks. To construct MOFs, we need to guide this assembly process. We do that by using dummy atoms, which indicate the points where the building blocks are to be connected. In practice, this is done as follows. Our parsing of MOFids generates, for MOFs with Zn tetramer nodes, linkers that carry two carboxyl groups, which by definition are part of metal nodes instead of linkers. Such wrong assignment of carboxyl groups is due to how the MOFid algorithm parses MOF structures. To reflect the correct molecular structure of linkers, the two carbon atoms in the carboxyl groups are identified as dummy atoms, and the redundant four oxygen atoms and two hydrogen atoms are removed. An illustration of dummy atom identification for linkers containing carboxyl groups is shown in the left panel of [Figure A1](#).

For MOFs with Cu paddlewheel and Zn paddlewheel node, however, another type of linker containing heterocyclic rings exists. For this type of linker, the two atoms that are nitrogen-metal bond distance away from the terminal nitrogen atoms on the heterocyclic rings are identified as dummy atoms. After the dummy atoms are correctly identified, three randomly selected linkers (duplicates allowed) are assembled with one of the three pre-selected nodes into MOFs in the pcu topology. An illustration of dummy atom identification of linkers containing heterocyclic groups is shown in the right panel of [Figure A1](#).

More than half of the **hMOF** structures are catenated MOFs, i.e., MOFs with interpenetrated lattices. By varying the level of interpenetration, it is possible to generate MOFs with different pore sizes, with a higher catenation level generally corresponding to smaller pores. We denote the four catenation levels in **hMOF**, with increasing number of interpenetrated lattices, as cat0, cat1, cat2 and cat3. To generate MOFs with high structural diversity, we applied the site translation method as implemented in **Pymatgen** [50] to generate MOFs with different catenation levels. For catenated MOF structures, we keep the spacing between interpenetrated lattices the same, so as to ensure equal pore sizes.

3.3. Applying regression model to estimate MOF CO₂ capacity

The **Predict** component of the **GHP-MOFassemble** framework uses a regression model to estimate the CO₂ capacity of newly generated pcu MOF structures. We use for this purpose a modified version of the Crystal Graph Convolutional Neural Network (CGCNN) model, developed in our previous work [18], which uses an adjacency list to format node and edge embeddings, rather than the adjacency matrix format of the original CGCNN, a change that enhances model training speed, training stability, and prediction accuracy. We trained this modified CGCNN model on all **hMOF** structures, with a target property of CO₂ capacity at 0.1 bar.

We create our ensemble regression model by training the modified CGCNN model three times and averaging the results. We applied the resulting model to the assembled pcu MOF structures produced by the **Generate** step of [Section 3.2](#) to identify those with a predicted CO₂ capacity higher than 2 mmol/g. These MOFs (see [Section 4.4.2](#)) are candidates for further investigation by, for example, by performing CO₂ capacity validation using Grand Canonical Monte Carlo (GCMC) simulations.

4. Results

By combining the generative power of diffusion models and the computing power of supercomputers, the **GHP-MOFassemble** framework enables rapid generation and assembly of chemically diverse MOF structures. Using the **hMOF** dataset as a starting point, our framework accelerates the computational design of high-performing MOF structures for carbon capture by high-throughput screening of novel MOF structures

with high CO₂ capacities. Our key findings are summarized below.

4.1. Analysis of the hMOF dataset

The three most frequent node-topology pairs in the hMOF dataset, accounting for around 74% of its MOFs, are the Cu paddlewheel-pcu, Zn paddlewheel-pcu, and Zn tetramer-pcu structures shown in the left panel of Figure 3. As summarized in Table 1, Zn tetramer-pcu is the most abundant, followed by Cu paddlewheel-pcu and then Zn paddlewheel-pcu. The cumulative distribution functions of the CO₂ capacities of these three types of MOF structure are shown in the right panel of Figure 3. Only the 78,238 MOFs from hMOF with correctly parsed MOFids and valid SMILES were used for analysis. We gather the molecular fragments numbers in Table 1, produced through the *Fragment* step described in Section 3.1. In Table 1 high-performing MOFs with three parsed linkers were selected, and their unique linkers parsed by using the MMPA (Matched Molecular Pairs Algorithm) algorithm.

In Table 1, the number of output molecular fragment conformers (last column) is much less than the number of unique linker SMILES (second to last column) for two reasons. First, around 56% of linkers did not pass the valency check, which may be due to the intrinsic limitations in the parsing of MOF linkers.

Second, around 90% of linkers that passed the valency check share similar molecular fragments, and thus many duplicated molecular fragments exist. The successfully parsed molecular fragment conformers are subsequently used for linker generation.

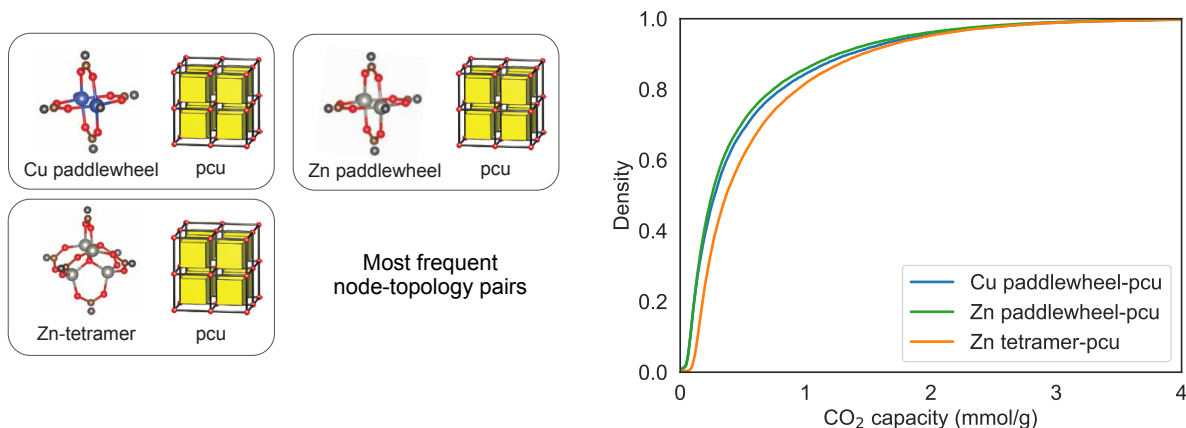


Figure 3: Depictions of the most frequent node-topology pairs in hMOF (left panel) and cumulative distribution functions of the 0.1 bar CO₂ capacities (right panel).

Figure 4 shows the empirical cumulative distribution functions of the CO₂ capacities of hMOF structures at different catenation levels. Therein, we observe that a higher percentage of catenated MOFs (cat1, cat2 and cat3) are high-performing, as compared to the uncatenated MOFs (cat0). This result confirms that catenation is an important factor when designing new MOF structures with high CO₂ capacity. In Table B1, we summarize the amount of origin shift for each of the four levels of catenation.

Table 1: Data on the hMOF with most frequent node-topology pairs. PW, TM, and HP-MOF stand for paddlewheel, tetramer, and high-performing MOF, respectively.

node-topology	total MOFs	HP-MOFs	HP-MOFs with three linkers	unique linker SMILES	unique molecular fragment conformers
Cu PW-pcu	29,714	1,458	1,016	3,330	180
Zn PW-pcu	28,529	1,314	834	3,221	162
Zn TM-pcu	43,874	2,129	1,388	3,265	198

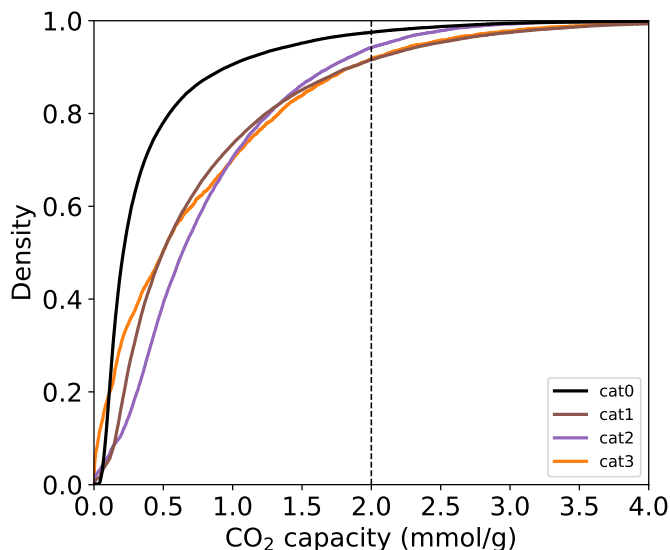
**Figure 4:** Empirical cumulative distribution functions of the 0.1 bar CO_2 capacities of hMOF structures at different catenation levels. The x axis is capped at 4 mmol/g to preserve details.

Figure 5 shows the pairwise relationships of the CO_2 capacities of hMOF structures at the five pressures. We observe a strong correlation of CO_2 capacities between 0.05 bar and 0.1 bar, with a Pearson’s correlation coefficient of 0.86. For other pairs of pressures, the CO_2 capacities are only weakly correlated, with a decreasing correlation for larger pressure differences. Moreover, the distributions of CO_2 capacities at all pressures exhibit long tails at high value ranges, which indicate that the majority of MOFs are low-performing and high-performing MOFs are uncommon.

4.2. Linker generation and evaluation

For all three MOF types, we start with 540 molecular fragments extracted from high-performing hMOF structures, and use DiffLinker to generate a large pool of 64,800 linkers with the number of sampled connection atoms ranging from 5 to 10. For each linker, sampling is performed 20 times, as described in Section 3.2. Since these linkers only contain heavy atoms, to generate all-atom linker molecules, we apply openbabel to add hydrogen atoms, which results in 56,257 linkers. Next, dummy atom identification

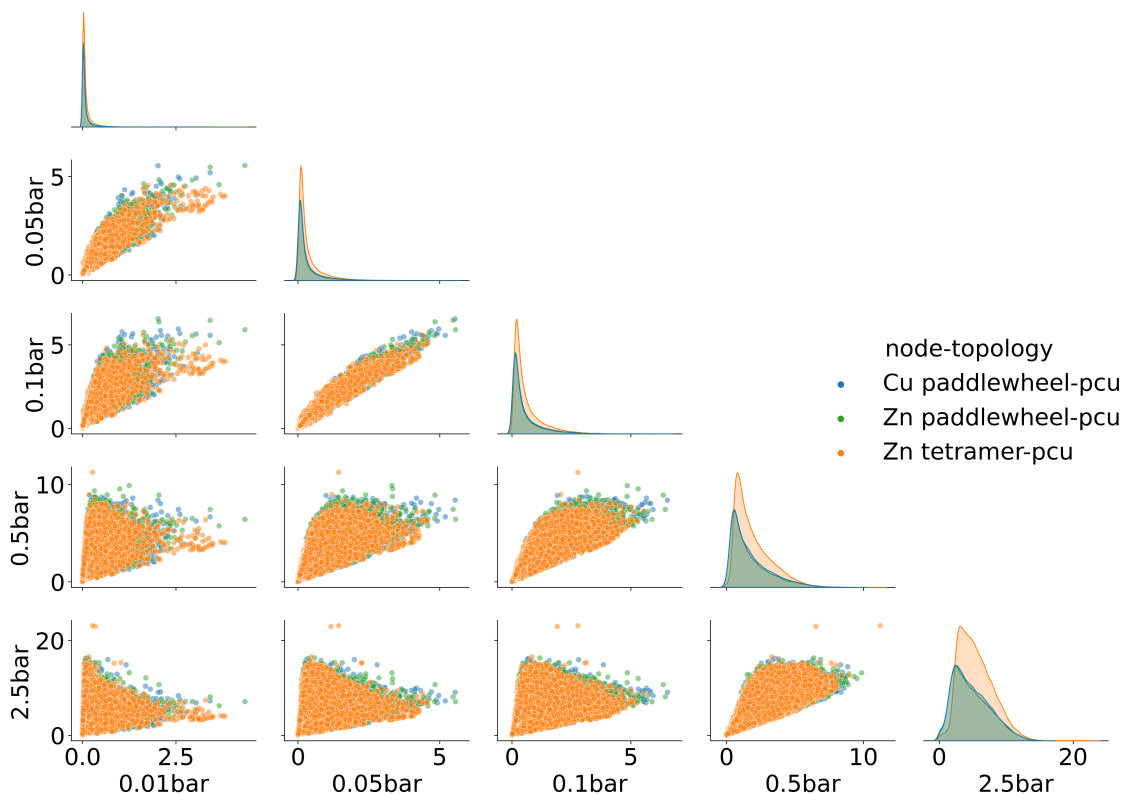


Figure 5: Pairplot of the CO₂ capacities of hMOF structures at different adsorption pressures. The x- and y-axis are adsorption pressures. Different colors indicate MOFs with different node-topology pairs.

is performed to generate information that enables assembly with metal nodes. A total of 16,162 linkers with generated dummy atoms are generated. These linkers are passed through the element filter, which removes linkers that contain S, Br and I elements, further reducing number to 12,305 linkers. The details of the number of linkers in each step are summarized in Table G1.

We show in Figure 6 the synthetic accessibility score (SAscore) and synthetic complexity score (SCscore) values, defined in Section 3.2.2, for the remaining linkers. We observe that as the number of sampled atoms increases, both distributions generally shift to the right, with the exception of the SAscore distribution with six sampled atoms, which is to the left of that with five atoms. The general trend indicates that linkers become harder to synthesize as the molecules become bigger. This result is expected because as more atoms are sampled, more complex substructures may be present. We note that no linker has zero or very large SAscore or SCscore, values that would indicate unsynthesizability.

We show in Table 2 the validity, uniqueness, and internal diversity metrics, defined in Section 3.2.2, grouped by number of sampled atoms and node-topology pairs of the corresponding MOFs. The *validity* column confirms that all generated linkers are valid. The last three columns, when reviewed from top to bottom, reveal a considerable

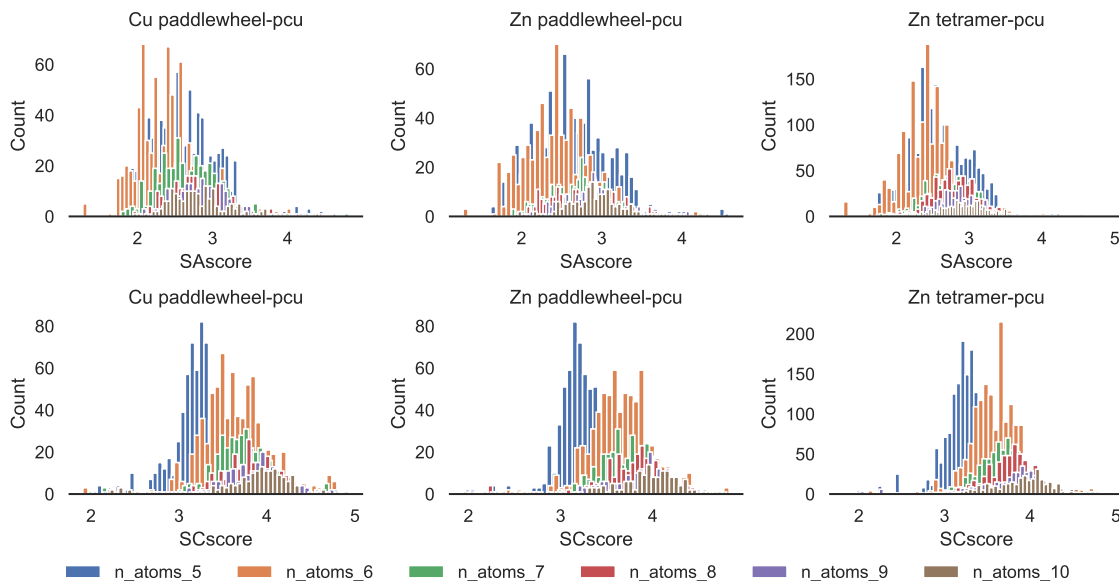


Figure 6: Distributions of synthetic accessibility score (SAscore, top row) and synthetic complexity score (SCscore, bottom row) of *DiffLinker*-generated MOF linkers with dummy atoms. The subfigure titles are node-topology pairs, and the colors listed in the legend indicate the number of sampled atoms, from five to 10 inclusive.

increase in linker uniqueness and a more modest increase in internal diversity as the number of sampled atoms increases. High uniqueness values indicate that our model is generating non-duplicate molecules, whereas high internal diversity values indicate that the generated molecules are chemically diverse. Internal diversity (which indicates how similar/dissimilar a specific linker is to the rest of the population) is computed using two internal diversity scores: IntDiv_1 and IntDiv_2 [48], also shown in Table 2—see Appendix C for details on how these scores are computed. The increase in these metrics as more atoms are sampled is expected because the degrees of freedom of atomic species and their spatial coordinates increase with molecular size.

4.3. MOF assembly

We generate new MOF structures by assembling three randomly selected *DiffLinker*-generated linkers with one of the three most frequent nodes in the *hMOF* dataset. Random sampling of linkers ensures that the selected linkers cover a large design space. For each node-linker-topology combination, we considered four levels of catenation (cat0, cat1, cat2, cat3). We generated a total of 120,000 MOFs as follows: we have four catenation levels and three node candidates. The random sampling of 10,000 linkers for each catenation level-node candidate pair generates 120,000 total MOFs of different catenation-node-linker combinations. We provide more details in Section 3.2.3 on how these new MOFs were assembled.

Table 2: Statistics of the generated linkers that correspond to different MOF types (first column), number of sampled atoms (\mathcal{N} in the second column), number of carboxyl linkers (\mathcal{C} , fourth column), and number of heterocyclic linkers (\mathcal{H} , fifth column). In total, there are 12,305 linkers with identified dummy atoms.

node-topology	\mathcal{N}	n_linker	\mathcal{C}	\mathcal{H}	valid	unique	IntDiv ₁	IntDiv ₂
Cu paddlewheel-pcu	5	1,240	774	466	1	0.708	0.789	0.769
Zn paddlewheel-pcu	5	1,184	775	409	1	0.692	0.783	0.762
Zn tetramer-pcu	5	1,666	1,666	0	1	0.734	0.739	0.721
Cu paddlewheel-pcu	6	1,117	761	356	1	0.764	0.819	0.796
Zn paddlewheel-pcu	6	992	681	311	1	0.761	0.809	0.787
Zn tetramer-pcu	6	1,532	1,532	0	1	0.766	0.749	0.732
Cu paddlewheel-pcu	7	499	383	116	1	0.958	0.836	0.820
Zn paddlewheel-pcu	7	453	329	124	1	0.929	0.829	0.812
Zn tetramer-pcu	7	709	709	0	1	0.958	0.782	0.770
Cu paddlewheel-pcu	8	330	245	85	1	0.992	0.835	0.822
Zn paddlewheel-pcu	8	338	247	91	1	0.978	0.831	0.817
Zn tetramer-pcu	8	550	550	0	1	0.983	0.791	0.780
Cu paddlewheel-pcu	9	281	191	90	1	0.99	0.837	0.824
Zn paddlewheel-pcu	9	257	175	82	1	0.993	0.834	0.821
Zn tetramer-pcu	9	389	389	0	1	0.986	0.794	0.784
Cu paddlewheel-pcu	10	235	165	70	1	0.996	0.842	0.830
Zn paddlewheel-pcu	10	217	148	69	1	0.996	0.837	0.826
Zn tetramer-pcu	10	316	316	0	1	0.995	0.798	0.789

4.4. Regression model for MOF CO₂ capacity prediction

4.4.1. Pre-training on *h*MOF structures We employ the modified version of the CGCNN model [18] as a predictive model of the CO₂ capacity of MOF structures. We trained this model on 90% of the MOF structures in *h*MOF, along with their CO₂ capacities at 0.1 bar, leaving 10% as a holdout set for testing. We trained the modified version of the CGCNN model three times to create an ensemble of models in order to mitigate biases in individual models. We summarize in Table D1 and Figure D1 the prediction errors of each model on the 10% test set and the distribution of standard deviations, respectively.

Looking next at the ensemble model, we show in the left panel of Figure 7 its predictions for the 10% test set. The ensemble model has a mean absolute error (MAE) of 0.093 mmol/g on the test set. Using CO₂ capacity of 2 mmol/g as a threshold, we repurposed our predictive model as a classifier to categorize MOFs into low and high performers (with predicted CO₂ capacities below and above the threshold, respectively)

and high performers. Using this scheme, the confusion matrix for identifying low and high performers is shown in the right panel of Figure 7.

We conclude from the confusion matrix that the pre-trained model classifies both low and high performers with high accuracy, with 98.4% (13551 out of 13765) of test samples correctly classified. As the test set is heavily imbalanced, with many more low performers than high performers, we also calculated the balanced accuracy of classification [51] (the sum of true positive rate and true negative rate divided by 2), obtaining a value of 90.7%. The majority of the 214 misclassified samples lie close to the decision boundaries, as shown in the red and purple regions of the left panel of Figure 7. Therefore, we also conclude that the ensemble model is capable of differentiating low and high performers.

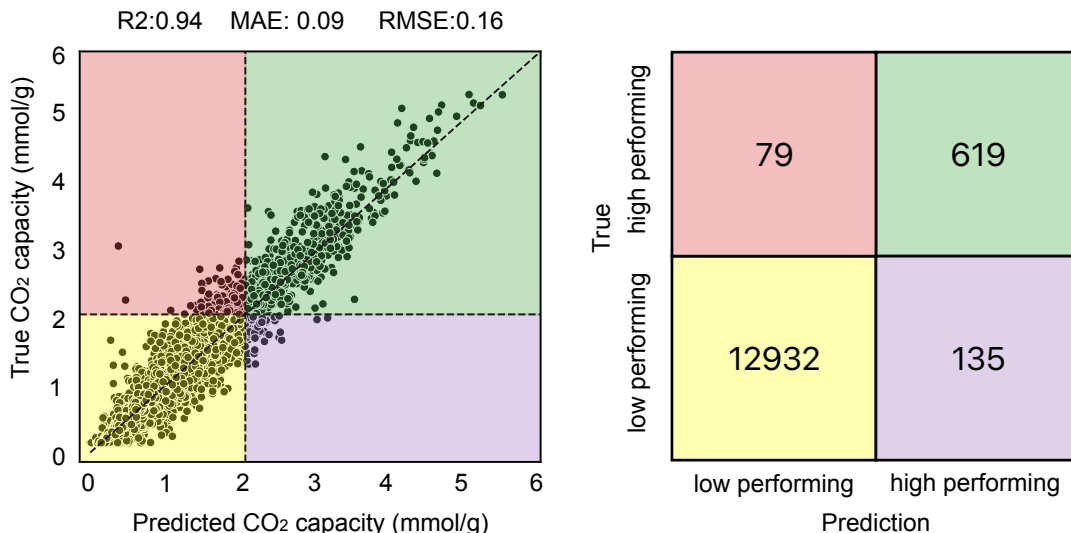


Figure 7: Left panel, predictive performance of the ensemble model on the 10% test set. Right panel, confusion matrix of the ensemble model in classifying low and high performers in the test set. The classification threshold is 2 mmol/g.

4.4.2. Predictions for newly generated MOF structures When using the ensemble model to infer the CO₂ capacities of newly generated MOF structures, we employ the average of the predictions made by the three independent models as the predicted capacity. In addition, we filter out structures for which the standard deviation of the predicted values obtained from the three models is larger than 0.2 mmol/g.

We treat the standard deviation of predictions from the three models as a measure of the uncertainty of model predictions, which we view as arising from the model’s difficulty in learning certain data points. We chose the threshold for the standard deviation filter of 0.2 mmol/g because it is sufficiently small (with 96% of data points

below the threshold) that for low CO₂ capacity predictions, the predicted values across the ensemble are similar. On the other hand, for high CO₂ capacity predictions, it also ensures that outlier values (e.g., predictions made by the three models are vastly different) are filtered out, therefore minimizing the overall error.

We show in Figure 8 the empirical cumulative distribution functions of both the hMOF (dashed lines) and predicted (solid lines) CO₂ capacities. We observe that there are more high-performing cat2 and cat3 MOFs among the generated structures than in the hMOF dataset, as indicated by lower cumulative density values at 2 mmol/g (black dashed line). On the other hand, there are fewer high-performing cat0 and cat1 MOFs, as indicated by higher cumulative density values at 2 mmol/g. This result indicates that within the generated MOF design space, there are more high-performing candidates with higher catenation levels compared to lower catenation levels.

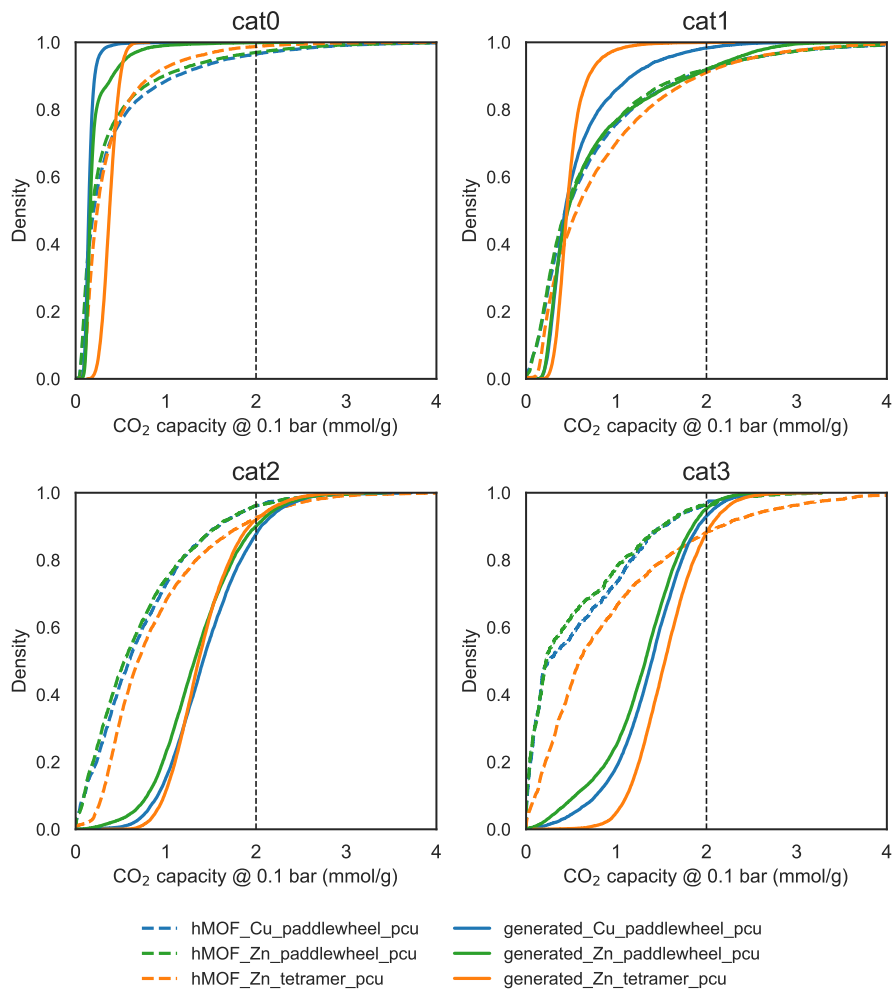


Figure 8: Comparison of empirical cumulative distribution functions of the predicted CO₂ capacities of generated and hMOF MOFs at different catenation levels.

Out of the 120,000 hypothetical MOF structures (see Section 4.3 for MOF assembly details), a total of 6,020 were predicted to be high-performing MOFs with CO₂ capacity

higher than 2 mmol/g at 0.1 bar. Categorization of the predicted high-performing MOFs by node-topology pair and catenation level (see Table 3) shows that the majority of Cu paddlewheel-pcu and Zn tetramer-pcu are cat2 and cat3, whereas the majority of high-performing Zn paddlewheel-pcu are cat1 and cat2. Table 4 shows the building blocks and catenation levels in the top five pcu MOF candidates. For linkers, the corresponding molecules (without dummy atoms) are shown for ease of visualization. Note that some substructures appear more often than others. For example, ring structures appear in most of the linkers in the top five candidates, and many of the DiffLinker-generated molecular substructures (in between the terminal fragments) also contain ring structures. The frequent occurrence of rings in linkers of high performing MOFs may be due to the presence of aromatic linkers, as reported by phenylene et al. [52].

We compared the chemical similarity of the linkers in predicted high-performing MOF structures with those in high-performing hMOF structures. As shown in Figure E1, we find that most linkers in the predicted high-performing MOF structures are vastly different from those in hMOF, indicating that the GHP-MOFassemble framework generate novel MOF structures with chemically unique linkers.

Table 3: Numbers of predicted high-performing MOFs found in the hMOF dataset.

node-topology	cat0	cat1	cat2	cat3
Cu paddlewheel-pcu	0	170	1,184	628
Zn paddlewheel-pcu	9	789	929	395
Zn tetramer-pcu	0	3	794	1,119

5. Validation of AI-generated MOFs with molecular dynamics simulations

We now examine the stability and porous properties of the 6,020 MOFs described in the previous section by using molecular dynamics (MD) simulations. We use the LAMMPS [53] code, version 28 Mar 2023 Update 1, distributed from the conda-forge channel, to perform MD simulations that relax the AI-generated MOF structures. An adapted version of cif2LAMMPS [54] is used to generate LAMMPS input files automatically. For each MOF, a $2 \times 2 \times 2$ supercell structure is generated, and an NPT simulation is performed, i.e., we use an isothermal-isobaric ensemble with a constant particle number N , a pressure p fluctuating around an equilibrium value, $\langle p \rangle = 1\text{atm}$, and a temperature T fluctuating around an equilibrium value $\langle T \rangle = 300\text{K}$. The cell lengths and angles of the triclinic simulation box are allowed to relax.

We run the NPT simulations for 200,000 steps with a timestep of 1 femtosecond. The relative changes of structure density before and after the simulation are inspected to evaluate how well the porous characteristics are maintained throughout the simulation. This process produces 102 MOFs whose structures stayed intact during the MD simulation. Next, we use a $<1\%$ change in density as the threshold to determine whether

Table 4: Details of the five pcu MOF candidates generated by GHP-MOFassemble with the highest predicted CO₂ capacities. Their stability was confirmed using molecular dynamics simulations, and their 3D structures are shown in Figure 10.

MOF ranking	1	2	3	4	5
Node	Zn paddlewheel	Zn paddlewheel	Zn paddlewheel	Zn tetramer	Zn tetramer
Linker 1					
Linker 2					
Linker 3					
Catenation	cat2	cat2	cat2	cat3	cat3
CO ₂ capacity (mmol/g)	2.51	2.49	2.48	2.44	2.43

a MOF structure is stable or not. Among the 102 MOF candidates, 18 have change in density below the 1% threshold. Figure 9 presents the top five AI-generated, stable MOF structures that passed the threshold for structure determination—see Table 4 for combinations of their building blocks and Figure 10 for 3D visualization of their crystal structures.

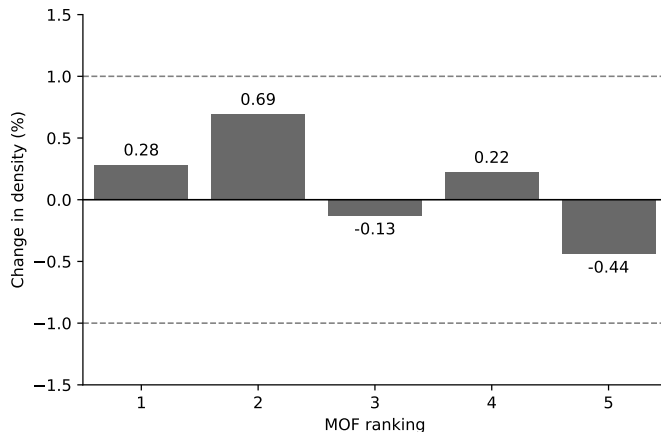


Figure 9: Top five AI-generated MOFs with structure density changed less than 1% during molecular dynamics simulations. Their CO₂ capacities are higher than 96.9% of hMOF structures in the hMOF database.

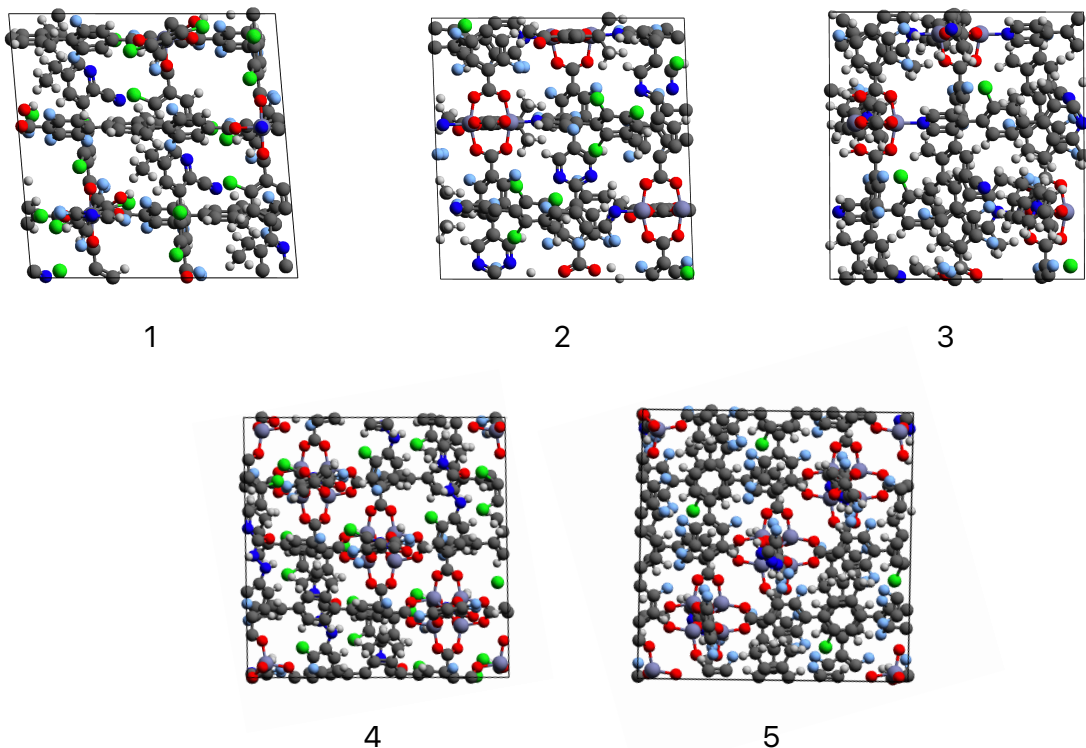


Figure 10: 3D visualization of the crystal structures of top five AI-generated, stable MOF candidates that retained their porous properties throughout molecular dynamics simulations. Atoms are presented using the color code: carbon in grey, nitrogen in dark blue, fluorine in cyan, zinc in purple, hydrogen in white, and lithium in green.

6. Conclusion

We have introduced **GHP-MOFassemble**, an AI-driven discovery framework for high-throughput generation of novel high-performing MOF structures for carbon capture. The **GHP-MOFassemble** framework employs a diffusion model, **DiffLinker**, to generate new, unique, and chemically diverse MOF linkers, which it then assembles with one of three pre-selected metal nodes into MOFs in the pcu topology.

Out of 120,000 MOF structures generated by **GHP-MOFassemble**, 6,020 were predicted to be high-performing MOFs with CO_2 capacity higher than 2 mmol/g at 0.1 bar, which corresponds to the top 5% of the **hMOF** dataset. We relaxed each of these 6,020 MOFs in a 200-picosecond NPT molecular dynamics simulation with the **UFF4MOF** force field [55, 56], and found that 102 structures did not collapse after molecular dynamics simulations, among which 18 were identified to be stable structures with less than 1% change in density.

We have deployed and extensively tested **GHP-MOFassemble** on computers at the Argonne Leadership Computing Facility and at the National Center for Supercomputing Applications, with the intent of providing scalable and computationally efficient AI tools

to accelerate the modeling and discovery of novel MOF structures. The tools introduced in this work may be readily fine-tuned and adapted to other available datasets beyond hMOF to enable accelerated design and discovery of novel MOF structures for carbon capture at industrial scale.

7. Acknowledgments

This work was supported by Laboratory Directed Research and Development (LDRD) funding from Argonne National Laboratory, provided by the Director, Office of Science, of the U.S. Department of Energy under Contract No. DE-AC02-06CH11357, and by the Braid project of the U.S. Department of Energy, Office of Science, Advanced Scientific Computing Research, under contract number DE-AC02-06CH11357. The work used resources of the Argonne Leadership Computing Facility, a DOE Office of Science User Facility supported under Contract DE-AC02-06CH11357. EAH and IF acknowledge support from National Science Foundation (NSF) award OAC-2209892. SC and XY acknowledge partial support from NSF Future of Manufacturing Research Grant 2037026. This research also used the Delta advanced computing and data resources, which is supported by the National Science Foundation (award OAC 2005572) and the State of Illinois. Delta is a joint effort of the University of Illinois at Urbana-Champaign and its National Center for Supercomputing Applications.

Code availability

The code and data are available upon reasonable request from the authors.

ORCID IDs

Hyun Park [0000-0001-5550-5610](#)
Xiaoli Yan [0000-0002-6512-2338](#)
Ruijie Zhu [0000-0001-9316-7245](#)
Eliu Huerta [0000-0002-9682-3604](#)
Santanu Chaudhuri [0000-0002-4328-2947](#)
Ian Foster [0000-0003-2129-5269](#)
Emad Tajkhorshid [0000-0001-8434-1010](#)

References

- [1] H. Li, K. Wang, Y. Sun, C. T. Lollar, J. Li, and H.-C. Zhou, “Recent advances in gas storage and separation using metal–organic frameworks,” *Materials Today*, vol. 21, no. 2, pp. 108–121, 2018.
- [2] M. Hao, M. Qiu, H. Yang, B. Hu, and X. Wang, “Recent advances on preparation and environmental applications of MOF-derived carbons in catalysis,” *Science of the Total Environment*, vol. 760, p. 143333, 2021.

- [3] H. D. Lawson, S. P. Walton, and C. Chan, "Metal-organic frameworks for drug delivery: A design perspective," *ACS Applied Materials & Interfaces*, vol. 13, no. 6, pp. 7004–7020, 2021.
- [4] M. J. Kalmutzki, N. Hanikel, and O. M. Yaghi, "Secondary building units as the turning point in the development of the reticular chemistry of MOFs," *Science Advances*, vol. 4, no. 10, p. eaat9180, 2018.
- [5] A. Chatterjee, X. Hu, and F. L.-Y. Lam, "Towards a recyclable mof catalyst for efficient production of furfural," *Catalysis Today*, vol. 314, pp. 129–136, 2018.
- [6] K. Tan, N. Nijem, Y. Gao, S. Zuluaga, J. Li, T. Thonhauser, and Y. J. Chabal, "Water interactions in metal organic frameworks," *CrystEngComm*, vol. 17, no. 2, pp. 247–260, 2015.
- [7] I. Erucar and S. Keskin, "Unlocking the effect of H₂O on CO₂ separation performance of promising MOFs using atomically detailed simulations," *Industrial & Engineering Chemistry Research*, vol. 59, no. 7, pp. 3141–3152, 2020.
- [8] Y. Zhang, Y. Zhang, X. Wang, J. Yu, and B. Ding, "Ultrahigh metal-organic framework loading and flexible nanofibrous membranes for efficient CO₂ capture with long-term, ultrastable recyclability," *ACS Applied Materials & Interfaces*, vol. 10, no. 40, pp. 34802–34810, 2018.
- [9] S. Zuluaga, E. M. Fuentes-Fernandez, K. Tan, F. Xu, J. Li, Y. J. Chabal, and T. Thonhauser, "Understanding and controlling water stability of MOF-74," *Journal of Materials Chemistry A*, vol. 4, no. 14, pp. 5176–5183, 2016.
- [10] Y. Jiao, C. R. Morelock, N. C. Burtch, W. P. Mounfield III, J. T. Hungerford, and K. S. Walton, "Tuning the kinetic water stability and adsorption interactions of Mg-MOF-74 by partial substitution with Co or Ni," *Industrial & Engineering Chemistry Research*, vol. 54, no. 49, pp. 12408–12414, 2015.
- [11] G. E. Cmarik, M. Kim, S. M. Cohen, and K. S. Walton, "Tuning the adsorption properties of UiO-66 via ligand functionalization," *Langmuir*, vol. 28, no. 44, pp. 15606–15613, 2012.
- [12] H. Huang, W. Zhang, F. Yang, B. Wang, Q. Yang, Y. Xie, C. Zhong, and J.-R. Li, "Enhancing CO₂ adsorption and separation ability of Zr (IV)-based metal-organic frameworks through ligand functionalization under the guidance of the quantitative structure-property relationship model," *Chemical Engineering Journal*, vol. 289, pp. 247–253, 2016.
- [13] S. M. Moosavi, A. Nandy, K. M. Jablonka, D. Ongari, J. P. Janet, P. G. Boyd, Y. Lee, B. Smit, and H. J. Kulik, "Understanding the diversity of the metal-organic framework ecosystem," *Nature Communications*, vol. 11, no. 1, pp. 1–10, 2020.
- [14] I. Igashov, H. Stärk, C. Vignac, V. G. Satorras, P. Frossard, M. Welling, M. Bronstein, and B. Correia, "Equivariant 3D-conditional diffusion models for molecular linker design," *arXiv preprint arXiv:2210.05274*, 2022.
- [15] S. Han, Y. Huang, T. Watanabe, Y. Dai, K. S. Walton, S. Nair, D. S. Sholl, and J. C. Meredith, "High-throughput screening of metal-organic frameworks for CO₂ separation," *ACS Combinatorial Science*, vol. 14, no. 4, pp. 263–267, 2012.
- [16] S. Li, Y. G. Chung, and R. Q. Snurr, "High-throughput screening of metal-organic frameworks for CO₂ capture in the presence of water," *Langmuir*, vol. 32, no. 40, pp. 10368–10376, 2016.
- [17] J. Rogacka, A. Seremak, A. Luna-Triguero, F. Formalik, I. Matito-Martos, L. Firlej, S. Calero, and B. Kuchta, "High-throughput screening of metal-organic frameworks for CO₂ and CH₄ separation in the presence of water," *Chemical Engineering Journal*, vol. 403, p. 126392, 2021.
- [18] H. Park, R. Zhu, E. Huerta, S. Chaudhuri, E. Tajkhorshid, and D. Cooper, "End-to-end AI framework for interpretable prediction of molecular and crystal properties," *Machine Learning: Science and Technology*, 2023.
- [19] C. Altintas, G. Avci, H. Daglar, A. N. V. Azar, I. Erucar, S. Velioglu, and S. Keskin, "An extensive comparative analysis of two MOF databases: High-throughput screening of computation-ready MOFs for CH₄ and H₂ adsorption," *Journal of Materials Chemistry A*, vol. 7, no. 16, pp. 9593–9608, 2019.
- [20] H. Dureckova, M. Krykunov, M. Z. Aghaji, and T. K. Woo, "Robust machine learning models for predicting high CO₂ working capacity and CO₂/H₂ selectivity of gas adsorption in metal

- organic frameworks for precombustion carbon capture,” *The Journal of Physical Chemistry C*, vol. 123, no. 7, pp. 4133–4139, 2019.
- [21] M. Pardakhti, E. Moharreri, D. Wanik, S. L. Suib, and R. Srivastava, “Machine learning using combined structural and chemical descriptors for prediction of methane adsorption performance of metal organic frameworks (MOFs),” *ACS Combinatorial Science*, vol. 19, no. 10, pp. 640–645, 2017.
- [22] G. S. Fanourgakis, K. Gkagkas, E. Tylianakis, and G. E. Froudakis, “A universal machine learning algorithm for large-scale screening of materials,” *Journal of the American Chemical Society*, vol. 142, no. 8, pp. 3814–3822, 2020.
- [23] C. Altintas, O. F. Altundal, S. Keskin, and R. Yildirim, “Machine learning meets with metal organic frameworks for gas storage and separation,” *Journal of Chemical Information and Modeling*, vol. 61, no. 5, pp. 2131–2146, 2021.
- [24] M. Fernandez, N. R. Trefiak, and T. K. Woo, “Atomic property weighted radial distribution functions descriptors of metal–organic frameworks for the prediction of gas uptake capacity,” *The Journal of Physical Chemistry C*, vol. 117, no. 27, pp. 14095–14105, 2013.
- [25] M. Fernandez, P. G. Boyd, T. D. Daff, M. Z. Aghaji, and T. K. Woo, “Rapid and accurate machine learning recognition of high performing metal organic frameworks for CO₂ capture,” *The Journal of Physical Chemistry Letters*, vol. 5, no. 17, pp. 3056–3060, 2014.
- [26] N. S. Bobbitt, K. Shi, B. J. Bucior, H. Chen, N. Tracy-Amoroso, Z. Li, Y. Sun, J. H. Merlin, J. I. Siepmann, D. W. Siderius, and R. Q. Snurr, “MOFX-DB: An online database of computational adsorption data for nanoporous materials,” *Journal of Chemical and Engineering Data*, 01 2023.
- [27] K. Choudhary and B. DeCost, “Atomistic Line Graph Neural Network for improved materials property predictions,” *npj Computational Materials*, vol. 7, no. 1, pp. 1–8, 2021.
- [28] Z. Yao, B. Sánchez-Lengeling, N. S. Bobbitt, B. J. Bucior, S. G. H. Kumar, S. P. Collins, T. Burns, T. K. Woo, O. K. Farha, R. Q. Snurr, and A. Aspuru-Guzik, “Inverse design of nanoporous crystalline reticular materials with deep generative models,” *Nature Machine Intelligence*, vol. 3, no. 1, pp. 76–86, 2021.
- [29] S. Bond-Taylor, A. Leach, Y. Long, and C. G. Willcocks, “Deep generative modelling: A comparative review of VAEs, GANs, normalizing flows, energy-based and autoregressive models,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 11, pp. 7327–7347, 2021.
- [30] A. Schneuing, Y. Du, C. Harris, A. Jamasb, I. Igashov, W. Du, T. Blundell, P. Lió, C. Gomes, M. Welling, M. Bronstein, and B. Correia, “Structure-based drug design with equivariant diffusion models,” *arXiv preprint arXiv:2210.13695*, 2022.
- [31] L. Huang, “A dual diffusion model enables 3D binding bioactive molecule generation and lead optimization given target pockets,” *bioRxiv*, pp. 2023–01, 2023.
- [32] G. Corso, H. Stärk, B. Jing, R. Barzilay, and T. Jaakkola, “DiffDock: Diffusion steps, twists, and turns for molecular docking,” *International Conference on Learning Representations*, 2023.
- [33] M. Xu, L. Yu, Y. Song, C. Shi, S. Ermon, and J. Tang, “GeoDiff: A geometric diffusion model for molecular conformation generation,” *arXiv preprint arXiv:2203.02923*, 2022.
- [34] C. Vignac, I. Krawczuk, A. Siraudin, B. Wang, V. Cevher, and P. Frossard, “DiGress: Discrete denoising diffusion for graph generation,” *arXiv preprint arXiv:2209.14734*, 2022.
- [35] E. Hoogeboom, V. G. Satorras, C. Vignac, and M. Welling, “Equivariant diffusion for molecule generation in 3D,” in *International Conference on Machine Learning*, pp. 8867–8887, 2022.
- [36] Z. Qiao, W. Nie, A. Vahdat, T. F. Miller III, and A. Anandkumar, “Dynamic-backbone protein-ligand structure prediction with multiscale generative diffusion models,” *arXiv preprint arXiv:2209.15171*, 2022.
- [37] M. Thomas, A. Bender, and C. de Graaf, “Integrating structure-based approaches in generative molecular design,” *Current Opinion in Structural Biology*, vol. 79, p. 102559, 2023.
- [38] C. E. Wilmer and R. Q. Snurr, “Towards rapid computational screening of metal-organic frameworks for carbon dioxide capture: Calculation of framework charges via charge

- equilibration,” *Chemical Engineering Journal*, vol. 171, no. 3, pp. 775–781, 2011.
- [39] B. J. Bucior, A. S. Rosen, M. Haranczyk, Z. Yao, M. E. Ziebel, O. K. Farha, J. T. Hupp, J. I. Siepmann, A. Aspuru-Guzik, and R. Q. Snurr, “Identification schemes for metal-organic frameworks to enable rapid search and cheminformatics analysis,” *Crystal Growth & Design*, vol. 19, no. 11, pp. 6682–6697, 2019.
- [40] J. Hussain and C. Rea, “Computationally efficient algorithm to identify matched molecular pairs (MMPs) in large data sets,” *Journal of Chemical Information and Modeling*, vol. 50, no. 3, pp. 339–348, 2010.
- [41] F. Imrie, A. R. Bradley, M. van der Schaar, and C. M. Deane, “Deep generative models for 3D linker design,” *Journal of Chemical Information and Modeling*, vol. 60, no. 4, pp. 1983–1995, 2020.
- [42] V. G. Satorras, E. Hoogeboom, and M. Welling, “E(n) equivariant graph neural networks,” in *International Conference on Machine Learning*, pp. 9323–9332, 2021.
- [43] S. Axelrod and R. Gomez-Bombarelli, “GEOM, energy-annotated molecular conformations for property prediction and molecular generation,” *Scientific Data*, vol. 9, no. 1, p. 185, 2022.
- [44] N. M. O’Boyle, M. Banck, C. A. James, C. Morley, T. Vandermeersch, and G. R. Hutchison, “Open Babel: An open chemical toolbox,” *Journal of Cheminformatics*, vol. 3, no. 1, pp. 1–14, 2011.
- [45] J. Ho, A. Jain, and P. Abbeel, “Denoising diffusion probabilistic models,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 6840–6851, 2020.
- [46] C. W. Coley, L. Rogers, W. H. Green, and K. F. Jensen, “SCScore: Synthetic complexity learned from a reaction corpus,” *Journal of Chemical Information and Modeling*, vol. 58, no. 2, pp. 252–261, 2018.
- [47] P. Ertl and A. Schuffenhauer, “Estimation of synthetic accessibility score of drug-like molecules based on molecular complexity and fragment contributions,” *Journal of Cheminformatics*, vol. 1, pp. 1–11, 2009.
- [48] D. Polykovskiy, A. Zhebrak, B. Sanchez-Lengeling, S. Golovanov, O. Tatanov, S. Belyaev, R. Kurbanov, A. Artamonov, V. Aladinskiy, M. Veselov, A. Kadurin, S. Johansson, H. Chen, S. Nikolenko, A. Aspuru-Guzik, and A. Zhavoronkov, “Molecular sets (MOSES): A benchmarking platform for molecular generation models,” *Frontiers in Pharmacology*, vol. 11, p. 565644, 2020.
- [49] G. Landrum, P. Tosco, B. Kelley, sriniker, gedec, NadineSchneider, R. Vianello, Ric, A. Dalke, B. Cole, AlexanderSavelyev, M. Swain, S. Turk, D. N, A. Vaucher, E. Kawashima, M. Wójcikowski, D. Probst, guillaume godin, D. Cosgrove, A. Pahl, JP, F. Berenger, strets123, JLVarjo, N. O’Boyle, P. Fuller, J. H. Jensen, G. Sforza, and DoliathGavid, “rdkit/rdkit: 2020_03.1 (q1 2020) release,” Mar. 2020.
- [50] S. P. Ong, W. D. Richards, A. Jain, G. Hautier, M. Kocher, S. Cholia, D. Gunter, V. L. Chevrier, K. A. Persson, and G. Ceder, “Python Materials Genomics (pymatgen): A robust, open-source Python library for materials analysis,” *Computational Materials Science*, vol. 68, pp. 314–319, 2013.
- [51] K. H. Brodersen, C. S. Ong, K. E. Stephan, and J. M. Buhmann, “The balanced accuracy and its posterior distribution,” in *20th International Conference on Pattern Recognition*, pp. 3121–3124, IEEE, 2010.
- [52] A. M. Plonka, D. Banerjee, W. R. Woerner, Z. Zhang, N. Nijem, Y. J. Chabal, J. Li, and J. B. Parise, “Mechanism of carbon dioxide adsorption in a highly selective coordination network supported by direct structural evidence,” *Angewandte Chemie International Edition*, vol. 52, no. 6, pp. 1692–1695, 2013.
- [53] A. P. Thompson, H. M. Aktulga, R. Berger, D. S. Bolintineanu, W. M. Brown, P. S. Crozier, P. J. in ’t Veld, A. Kohlmeyer, S. G. Moore, T. D. Nguyen, R. Shan, M. J. Stevens, J. Tranchida, C. Trott, and S. J. Plimpton, “Lammps - a flexible simulation tool for particle-based materials modeling at the atomic, meso, and continuum scales,” *Computer Physics Communications*,

- vol. 271, p. 108171, 2022.
- [54] R. Anderson, “cif2lammps.”
- [55] M. A. Addicoat, N. Vankova, I. F. Akter, and T. Heine, “Extension of the universal force field to metal–organic frameworks,” Journal of Chemical Theory and Computation, vol. 10, no. 2, pp. 880–891, 2014.
- [56] D. E. Coupry, M. A. Addicoat, and T. Heine, “Extension of the universal force field for metal–organic frameworks,” Journal of Chemical Theory and Computation, vol. 12, no. 10, pp. 5215–5225, 2016.

Appendix A. Dummy atom identification

Figure A1 illustrates the workflow used to identify the dummy atoms and to remove redundant atoms for linkers containing carboxyl groups (left panel) or heterocyclic rings (right panel).

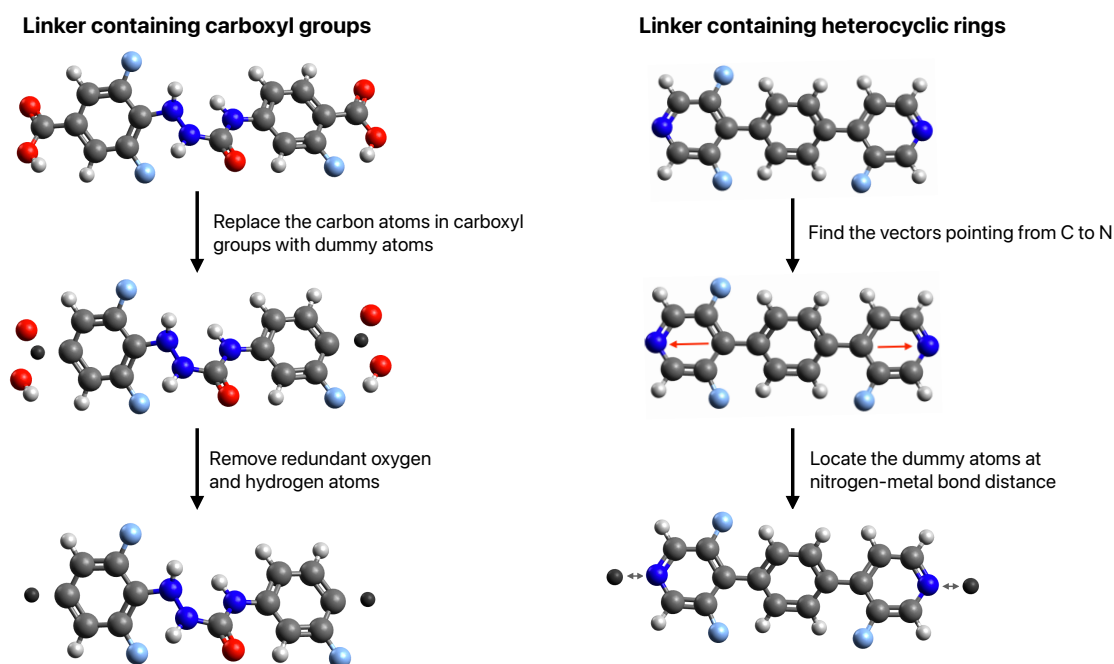


Figure A1: Left panel, identification of dummy atoms for linkers containing carboxyl groups. The dummy atoms are found by substituting the carbon atoms in the carboxyl groups. The remaining oxygen and hydrogen atoms in the carboxyl groups are removed. Right panel, identification of dummy atoms for linkers containing heterocyclic rings. The dummy atoms are found at nitrogen-metal bond distance from the terminal nitrogen atoms along the vectors pointing from the opposing carbon atoms to nitrogen atoms.

Appendix B. MOF catenation via site translation

Catenated MOFs were generated by using the site translation method, which is achieved by displacing the reference lattice along the diagonal line of the unit cell. For all four catenation levels, the amount of relative lattice displacement is given in Table B1. The numbers are fractional displacements relative to the unit cell diagonal.

Table B1: Amount of lattice displacement for catenated MOFs

catenation	lattice1	lattice2	lattice3	lattice4
cat0	(0,0,0)	/	/	/
cat1	(0,0,0)	$(\frac{1}{2}, \frac{1}{2}, \frac{1}{2})$	/	/
cat2	(0,0,0)	$(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$	$(\frac{2}{3}, \frac{2}{3}, \frac{2}{3})$	/
cat3	(0,0,0)	$(\frac{1}{4}, \frac{1}{4}, \frac{1}{4})$	$(\frac{1}{2}, \frac{1}{2}, \frac{1}{2})$	$(\frac{3}{4}, \frac{3}{4}, \frac{3}{4})$

Appendix C. Linker evaluation

Recall from Section 4.2 that we use the synthetic accessible score (SAscore) and synthetic complexity score (SCscore) metrics to confirm that all linkers are synthesizable. Table 2 provides additional data on the generated linkers. We also provide results for the validity and uniqueness metrics, and the internal diversity scores IntDiv_1 and IntDiv_2 [48]. The internal diversity scores were calculated based on the Tanimoto distances among all pairs of molecules, which are obtained by calculating the normalized Jaccard score of Morgan fingerprint bit vector between all pairs of generated linkers. This yields similarities between a specific linker and the rest of the linkers in the linker pool, which are averaged out. Therefore, we end up with a metric showing each linker’s similarity compared with the population of generated linkers. We used the MOSES [48] framework to compute these two scores using the relations:

$$\text{IntDiv}_1(G) = 1 - \frac{1}{|G|^2} \sum_{m_1, m_2 \in G} T_d(m_1, m_2), \quad (\text{C.1})$$

$$\text{IntDiv}_2(G) = 1 - \sqrt{\frac{1}{|G|^2} \sum_{m_1, m_2 \in G} T_d(m_1, m_2)^2}, \quad (\text{C.2})$$

where G is the generated set of molecules, $|G|$ is the size of molecule set, (m_1, m_2) represents a pair of molecules in the molecule set. T_d is the Tanimoto distance, which relates to the Tanimoto similarity (T_s) by:

$$T_d(m_1, m_2) = 1 - T_s(m_1, m_2) \quad (\text{C.3})$$

Appendix D. Statistics and distribution of the standard deviation of the ensemble of pre-trained models

Table D1 shows the R^2 score, mean absolute error (MAE), and root mean squared error (RMSE) of the three pre-trained models used for creating the ensemble model. Figure D1 shows the distribution of the standard deviations of ensemble model predictions. The threshold 0.2 mmol/g serves as a criterion for assessing whether the predictions made by the three pre-trained models are in agreement. The ensemble model predictions with standard deviation above the threshold were discarded.

Table D1: Statistics of the three pre-trained models

Model	R^2	MAE	RMSE
Model1	0.932	0.098	0.171
Model2	0.937	0.100	0.170
Model3	0.936	0.099	0.170

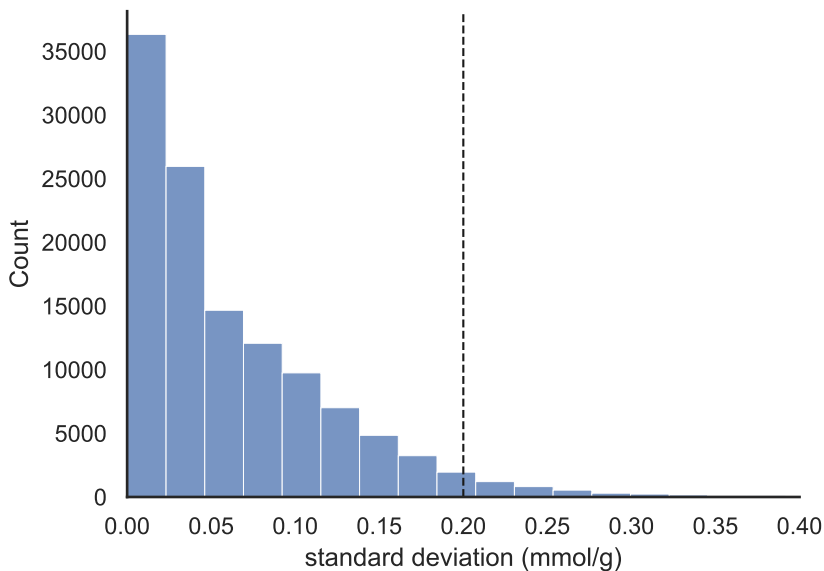


Figure D1: Distribution of the standard deviation of ensemble model predictions. The ensemble model predictions with standard deviation of less than 0.2 mmol/g consist of around 96% of the test set. This shows that our three independently trained models are in great agreement for most of the unobserved test set, showing robustness of our model. The remaining 4% with larger than 0.2 mmol/g standard deviation implies that these are the data points which may be difficult to predict due to errors in target property or extra information other than atomic species and periodic neighbor is necessary for accurate prediction.

Appendix E. Linker similarity for high-performing MOF structures

To measure the similarity between the generated linkers and hMOF linkers in high-performing MOF structures, we show in Figure E1 the distribution of maximum Tanimoto similarity between the two linker sets. For each unique linker in the predicted high-performing MOF structures generated by the GHP-MOFassemble framework, we calculate its Tanimoto similarity with all of the unique linkers in the hMOF structures. The maximum value of the Tanimoto similarities gives a quantitative measure of how different the generated linkers are from the hMOF linkers.

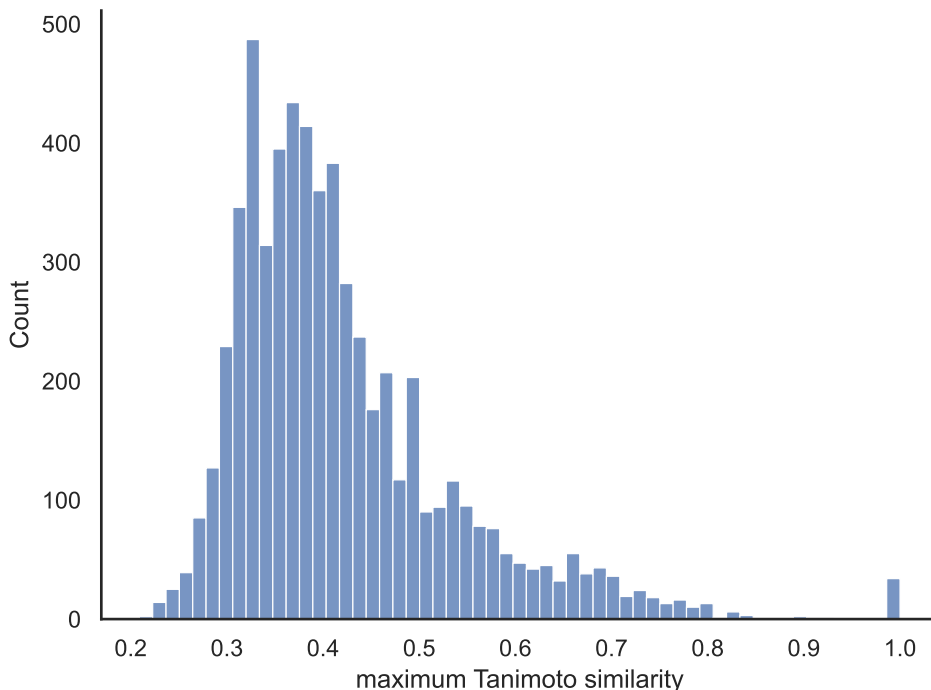


Figure E1: Distribution of maximum Tanimoto similarity between generated and hMOF linkers in high-performing MOF structures. The peak around 0.3 to 0.4 indicates that most generated linkers are just 30–40% similar to those in hMOF—showing that we are generating novel linkers not present in hMOF structures. On the other hand, the trailing heavy right tail above 0.4 Tanimoto similarity indicates that we are also able to generate linkers that are structurally similar to those present in hMOF, showing that GHP-MOFassemble enables generation of a diverse set of novel linkers.

Appendix F. Regression model training details

We independently trained three modified versions of the CGCNN model for 5,000 epochs with a batch size of 160, and with optimizer, learning rate, and weight decay of `torch_adam`, $1e-4$, and $2e-5$, respectively.

Appendix G. Linker statistics

Table G1: Statistical summary of the number of linkers after each step categorized by the corresponding MOF types. PW and TM stand for paddlewheel and tetramer, respectively. Element filter means that linkers with S, Br and I elements are removed.

Method	Total	Cu PW-pcu	Zn PW-pcu	Zn TM-pcu
DiffLinker	64,800	21,600	19,440	23,760
Hydrogen addition	56,257	18,979	17,126	20,152
Dummy atom identification	16,162	4,964	4,450	6,748
Element filter	12,305	3,702	3,441	5,162