

## **Inverse design of nanoporous crystalline reticular materials with deep generative models**

*Zhenpeng Yao<sup>1,2\*</sup>, Benjamín Sánchez-Lengeling<sup>1</sup>, N. Scott Bobbitt<sup>3</sup>, Benjamin J. Bucior<sup>3</sup>, Sai Govind Hari Kumar<sup>2</sup>, Sean P Collins<sup>4</sup>, Thomas Burns<sup>4</sup>, Tom K. Woo<sup>4</sup>, Omar K. Farha<sup>3,5</sup>, Randall Q. Snurr<sup>3\*</sup>, Alán Aspuru-Guzik<sup>1,2,6,7\*</sup>*

*<sup>1</sup>Department of Chemistry and Chemical Biology, Harvard University, 12 Oxford Street, Cambridge, Massachusetts 02138, United States*

*<sup>2</sup>Chemical Physics Theory Group, Department of Chemistry and Department of Computer Science, University of Toronto, Toronto, Ontario M5S 3H6, Canada*

*<sup>3</sup>Department of Chemical and Biological Engineering, Northwestern University, 2145 Sheridan Road, Evanston, Illinois 60208, United States*

*<sup>4</sup>Department of Chemistry and Biomolecular Science, University of Ottawa, 10 Marie Curie Private, Ottawa K1N 6N5, Canada*

*<sup>5</sup>Department of Chemistry, Northwestern University, 2220 Campus Drive, Evanston, Illinois 60208, United States*

*<sup>6</sup>Vector Institute for Artificial Intelligence, Toronto, Ontario M5S 1M1, Canada*

*<sup>7</sup>Canadian Institute for Advanced Research (CIFAR) Lebovic Fellow, Toronto, Ontario M5S 1M1, Canada*

*\*Correspondence: alan@aspuru.com, snurr@northwestern.edu, yaozhenpeng@gmail.com*

## Abstract

Reticular frameworks are crystalline porous materials that form *via* the self-assembly of molecular building blocks (*i.e.*, nodes and linkers) in different topologies. Many of them have high internal surface areas and other desirable properties for gas storage, separation, and other applications. The notable variety of the possible building blocks and the diverse ways they can be assembled endow reticular frameworks with a near-infinite combinatorial design space, making reticular chemistry both promising and challenging for prospective materials design. Here, we propose an automated nanoporous materials discovery platform powered by a supramolecular variational autoencoder (SmVAE) for the generative design of reticular materials with desired functions. We demonstrate the automated design process with a class of metal-organic framework (MOF) structures and the goal of separating CO<sub>2</sub> from natural gas or flue gas. Our model exhibits high fidelity in capturing structural features and reconstructing MOF structures. We show that the autoencoder has a promising optimization capability when jointly trained with multiple top adsorbent candidates identified for superior gas separation. MOFs discovered here are strongly competitive against some of the best-performing MOFs/zeolites ever reported. This platform lays the groundwork for the design of reticular frameworks for desired applications.

## 1. Introduction

Reticular frameworks (which include metal-organic frameworks and covalent-organic frameworks) are crystalline porous materials, many of which feature high internal surface area and high stability. They are formed *via* the self-assembly of molecular building blocks (*i.e.*, nodes and linkers) in different topologies. The notable variety of the possible building blocks and the diverse ways they can be assembled endow reticular frameworks with exceptional geometrical and chemical tunability.<sup>1</sup> Since the first metal-organic framework (MOF),<sup>2</sup> thousands of reticular frameworks have been made towards various applications with remarkable advances achieved in fields such as gas storage,<sup>3-5</sup> molecular separation,<sup>6-9</sup> catalysis,<sup>10,11</sup> sensing,<sup>12</sup> electrochemical energy storage,<sup>13-15</sup> and drug delivery.<sup>16</sup> Aiming at a particular application, novel reticular frameworks can be designed in a trial-and-test manner through selecting plausible building blocks that assemble in a desired topology.<sup>17</sup> Given the vastness of chemical space for small molecules<sup>18,19</sup> that can potentially be used as linkers, reticular frameworks exhibit a near-infinite combinatorial design space. The boundless design space significantly expands the scope of useful materials for prospective applications, yet its enormity also complicates its systematic exploration. Therefore, the search for new materials becomes a constrained global optimization problem in the high dimension space.

One powerful approach developed to assist the discovery of reticular frameworks is high-throughput (HT) computational<sup>20-22</sup> and experimental<sup>23-26</sup> screening. HT screening proceeds *via* generating/synthesizing and evaluating all the frameworks (building block combinations) from a selected library. The HT computational methodology has enabled the examination of a design space on the order of 10<sup>3</sup>~10<sup>5</sup>. One main drawback of this approach is the low coverage and

restriction of the search space according to the combinatorics of the building blocks. In addition to HT screening, heuristic optimization approaches include genetic algorithms (GA) and evolutionary strategies (ES). Given a score metric and a set of candidates, these methods can transform/evolve/mutate the candidates based on their scored performance, eventually leading to higher scoring structures. This approach allows the search of larger spaces and has been successful at identifying top-performing MOFs in recent studies.<sup>21,27–29</sup> The downside of this approach is that it requires specifying prior rules on how to transform the frameworks, which then creates a preceding constraint of the types of frameworks that can be explored.

Another promising approach for optimizing frameworks lies with machine learning (ML) algorithms that are able to learn from data and improve their performance automatically through experience. Among them, predictive algorithms (*i.e.*, discriminative models), those that given a datapoint  $x$  aim to predict a property  $y$ , have been used to aid or even replace physical simulations under certain circumstances. Discriminative models have been widely applied to accelerate the HT screening process of reticular frameworks for properties such as storage,<sup>30,31</sup> mechanical stability,<sup>32</sup> synthesizability,<sup>33</sup> etc.<sup>34</sup> Another class of algorithms that do not necessarily deal with predicting a property  $y$  but modeling the data itself are generative models. For example, a Bernoulli probability distribution can be used as a model to generate a coin flip. With more complicated data distributions, we can use deep generative models such as variational autoencoders<sup>35</sup> and generative adversarial networks.<sup>36</sup> In these cases, the mapping between probability distributions and data is learned *via* a deep neural network; and this map can be further enhanced with additional information (physical properties) to condition or bias the generative process. By conditioning the generative process on a property of interest, the models can be employed to generate preferentially molecules with a given property. This property-to-structure approach is called inverse design.<sup>37</sup> Generative models can be used as a key component to realize the automated “closed-loop” design of materials towards targeted performances. These models have been successfully applied to a variety of molecular<sup>38–41</sup> and material design applications.<sup>42–44</sup> In the context of reticular frameworks, we can design and generate a framework by sampling a random vector and mapping it back to a learned data distribution. Other challenges relevant to closing the loop are the planned synthesis of a reticular framework given a set of materials and the potential automatic robotic realization of this procedure.

The primary goal in the reticular framework design presented in this work is the guided optimization of crystal structures according to a targeted functionality. In a simplified manner, a reticular framework could be seen as a large collection of regularly bonded particles (atoms) in 3D space. Optimization with this representation corresponds to optimizing the number of particles, the identities of these particles, and their positions. The realization of this optimization is then quite challenging due to the high and variable dimensionality, large particle number, along with a mix of discrete and continuous variables. Therefore, finding an efficient representation becomes the essential step for machine learning based reticular framework optimization. One way to attack the problem is reduction approximation through exploiting symmetries and hierarchical structures of the systems. An ideal representation would encode the degrees of freedom, physical symmetries,

and constraints of a system and be amenable to gradient-based optimization techniques. The representation should also be decodable such that a framework can be re-constructed or decoded back. With proper representation, deep generative models show great promise for reticular framework optimization because of their potential capability to map these frameworks into a continuous vector representation. Variational autoencoders (VAEs), in particular, which can learn an invertible mapping, encode a material to a vector (*i.e.*, its latent vector), and decode it back to a framework, are a compelling solution. Optimization of materials can be ultimately made in the latent vector space within the VAE framework, which then lays the ground for the design of reticular frameworks with desired properties.

In this work, we build an automated nanoporous materials discovery platform for the property-orientated generative design of reticular frameworks, empowered by a supramolecular variational autoencoder. We develop a semantically constrained graph-based canonical code for the efficient representation of reticular frameworks (RFcode). With MOF structures from the computation-ready, experimental (CoRE 2019-ASR) MOF database<sup>45</sup> as inputs and clean energy applications (*i.e.*, CO<sub>2</sub>/N<sub>2</sub>, CO<sub>2</sub>/CH<sub>4</sub> separations) as the exemplified targets, we demonstrate the automated design process using a discovery platform for novel MOF structures with remarkably improved performance. By examining the latent space of our model, we illustrate that our representation captures structural features while also organized around properties. We demonstrate its capabilities for automatic targeted generation by proposing top candidates for gas separation adsorbent materials. We believe MOFs discovered here are strongly competitive against some of the best-performing MOFs/zeolites ever reported in the literature. We make our trained models, results and code available as open-source to aid reproducibility and adoption to broader applications (*e.g.*, covalent-organic frameworks, metal-organic polyhedra, hydrogen-bonded organic framework, coordinational polymers).

## Reticular Framework Representation and Identification

All crystalline materials can be seen as a collection of particles with different identities arranged periodically in 3D space. Given the identities and positions of the atoms, in principle, any property can be computed for the framework from the Schrödinger equation (SE). However, in practice, this may be difficult due to computational complexity and cost, which lead to tradeoffs generally made in the form of approximations. Another approach is to estimate material properties using models such as linear models or neural networks that learn transformations on their input representations. Ideally, the representation would contain the same symmetries the SE presents: translational, rotational, and permutational invariance with respect to its atomic identities.<sup>37</sup> Meanwhile, representations and models are coupled such that different types of inputs will lead to distinct choices of preferred models (*e.g.*, images and convolutional networks). Materials representation currently is an open research problem, while for non-periodic chemical systems (*e.g.*, molecules), several representations have been proven successful, such as fingerprints,<sup>46</sup> SMILES,<sup>38</sup> SELFIES,<sup>47</sup> and graphs.<sup>40,41</sup> Defining a representation for periodic crystalline materials is more challenging because of the necessity to deal with the extra-dimensional connections at the

border of unit cells. Particularly for reticular frameworks, their generally larger cell sizes ( $10^2\sim 10^4$  atoms<sup>48</sup> *versus* common crystalline materials with  $10^1\sim 10^2$  atoms<sup>49</sup>) bring further difficulties in representing them efficiently. Methods like the smooth overlap of atomic positions,<sup>50</sup> Voronoi tessellation,<sup>51,52</sup> diffraction images,<sup>53</sup> and multi-perspective fingerprints<sup>54</sup> have been suggested for crystalline materials classification, property prediction, and so on. Some of the most promising representations under development are graph-based<sup>55,56</sup> algorithms, where atoms are encoded as vertices and atom-pairs (*i.e.*, bonds) as edges. They can be effective without encoding positional coordinates explicitly. However, applying this representation to typical reticular frameworks results in graphs with  $10^2\sim 10^4$  vertices and 3-5 edges per vertex,<sup>48</sup> leading to a space with billions of potential configurations. Barely any effective optimizations can be done in a space of this size using the graph models, and thus reductions are called for. Tiling, net, and graph theories<sup>57-59</sup> can be used to aid the reduction by replacing atom-based vertices with motif-based vertices and bond-based edges with polyatomic-branch-based edges that connect these motifs.

Inspired by these reduction theories, we construct our representation of the reticular frameworks (*i.e.*, RFcode) using their unique, decomposed nets as a tuple: edges|vertices|topologies. Edges are molecular fragments with two connection points, vertices are multi-connected metal or organic nodes, and topologies define how these components are connected to form a specific reticular framework. Within the RFcode, we consider the edges as semantically constrained graphs,<sup>47</sup> while vertices and topologies are categorical variables from known frameworks considering their relatively limited variety. In addition, we consider metal and organic vertices separately in RFcode. The advantages of RFcode are: 1) efficiency, since there is no redundant information, edges and vertices are only described once in RFcode. 2) uniqueness, each representation encodes a unique framework. 3) invertibility, since all components can be readily translated back and forth. Moreover, for each component of the RFcode, generative models have been effectively developed, and therefore a model that takes the full RFcode is realizable. To illustrate this method, we use MOF-117<sup>60</sup> as an example, and its representation is shown in **Fig. 1A**.

The RFcode representations of reticular frameworks can be determined automatically using a previously developed identification algorithm supplemented with framework deconstruction<sup>61</sup> and reconstruction tools.<sup>62,63</sup> As a demonstration of our method, we decomposed all MOF structures from the CoRE MOF 2019-ASR database into their building blocks and identified all their RFcodes. Meanwhile, collections of edges, vertices (metal and organic), and topologies were also built. To have a sense of the chemical variety of all linkers in the CoRE database, we conducted a fragmentation analysis of them using molBLOCKS,<sup>64</sup> and the derivations of different linkers are illustrated using a scaffold tree plot shown in **Fig 1B**. Note that here and in the traditional MOF terminology, an organic “linker” may be a single edge (connecting two metal vertices) or may contain an organic vertex and several edges.

## Reticular Framework (*i.e.*, MOF) Library Generation

While there are no established rules on the sufficient size of training datasets for deep generative models, empirically these models start to be useful when input datasets are on the order of  $10^6$ . With the correct architecture, at this scale, the model can begin to generate new data that is likely to come from the empirical data distribution. Considering that there are only around 14000 MOFs in the CoRE MOF database, a training data augmentation is necessary. Starting with all the MOF edges obtained from CoRE MOF identification (372 edges), we did random functionalizations (**Fig. S1 and Tab. S1**) with selected common functional groups of known MOF structures (**Tab. S1**). An augmented edges dataset of  $\sim 300,000$  was generated. Vertex and topology datasets are constructed during the identification of the CoRE database as mentioned in the previous section by selecting vertices and topologies that are compatible with the current reticular framework reconstructor.<sup>62,63</sup> Therefore, all these datasets are subject to further expansions in the future with improvements of the reconstructor. We then used this augmented edge dataset with the vertex dataset (metal: 14, **Fig. S2**; organic: 47, **Fig. S3**) and topology dataset (152, **Fig. S4**), resulting in an augmented dataset with around 2 million MOF structures. An underlying assumption in our dataset is that the current vertex and topology pools represent plausible and realizable structures for reticular frameworks. Our search space does not include new vertices and topologies.

Besides generating new structures, we are interested in making our model aware of properties of interests. Doing so with deep neural networks requires having a large dataset of reticular frameworks (RFcodes) as well as properties, preferably experimental. However, such a dataset is currently lacking; therefore, we resorted to computational simulations on around 40k randomly-selected MOF structures. The randomness allows coverage of multiple types of frameworks, and the quantity is to keep the computational cost at a reasonable level. We considered properties as follows: 4 textural properties (pore-limiting diameter, largest cavity diameter, density, accessible gravimetric surface area), 3 properties related to natural gas separation ( $\text{CO}_2$  uptake,  $\text{CH}_4$  uptake,  $\text{CO}_2/\text{CH}_4$  selectivity, all at 5 bar and 300 K for a 10/90 mole fraction mixture of  $\text{CO}_2/\text{CH}_4$ ), and 3 properties related to flue gas separation ( $\text{CO}_2$  uptake,  $\text{N}_2$  uptake,  $\text{CO}_2/\text{N}_2$  selectivity, all at 1 bar and 313 K for a 15/85 mole fraction mixture of  $\text{CO}_2/\text{N}_2$ ). Textural properties were calculated geometrically, and gas uptake properties were calculated using grand canonical Monte Carlo simulations. Gas separation selectivities, which are entirely dependent on the uptake values of the mixed gas phases from the mixed-gas simulations, were then derived numerically. We use pore-blocking to prevent insertions into cavities that are inaccessible to the adsorbate molecules due to narrow windows. Therefore, some reticular frameworks may have extremely small or even zero uptakes of the larger radius molecules like  $\text{CH}_4$  and  $\text{N}_2$ . As a result, these frameworks are predicted to have enormous or even infinite ( $\infty$ ) selectivity of  $\text{CO}_2$  against  $\text{CH}_4$  and  $\text{N}_2$ . In reality, the observed selectivities may not be perfect (infinite) because the frameworks may not be completely rigid and large adsorbate molecules may not be totally blocked. Further details are described in **Supplementary Note 1**. The distributions of the textural properties for these 40k MOFs are shown in **Fig. S5**, and the distributions of gas uptake properties for these 40k MOFs are shown in **Fig. S6**.

## Supramolecular Variational Autoencoder

For our deep generative model, we utilize a variational autoencoder<sup>65</sup> (VAEs). A VAE is trained to process and reconstruct non-labeled data in an unsupervised manner. In its simplest form, a VAE is composed of two components: an encoder and a decoder. For a given datapoint  $\mathbf{x}$ , the encoder compresses the information to a vector  $\mathbf{z}$ , and the decoder decompresses the data into a reconstructed sample  $\tilde{\mathbf{x}}$ . To learn these transformations, neural networks are used as computational and optimizable building blocks for the encoder and decoder. The encoder and decoder are then optimized according to a loss, which is a low reconstruction error ( $\|\mathbf{x}-\tilde{\mathbf{x}}\|$ ). In order to generalize to new data points, a VAE imposes a prior over the structure of the vector space  $\mathbf{z}$ , and this lower-dimensional space, namely the latent space, is in our case normally distributed. To enforce this constraint, an additional term is introduced in the loss function, the Kullback–Leibler (KL) divergence of the variational approximation.<sup>38</sup> This term can also be interpreted as a regularization term. It measures how our latent space resembles a normal Gaussian distribution. A cyclical annealing scheduler, which has been proven to be effective in boosting the training performance and mitigating KL vanishing,<sup>66</sup> was also adopted.

Considering that our reticular framework representation, namely the RFcode, is a multiple component input, we build our Supramolecular Variational Autoencoder (SmVAE) with several corresponding components that are in charge of encoding and decoding each part of the RFcode. When properly trained, this model allows us to map the frameworks with discrete representations (RFcodes) into continuous vectors ( $\mathbf{z}$ ) and then back. Since the latent space is a vector space, continuous optimization and search algorithms will be used to find local minima or maxima. By decoding, we can sample and reconstruct new frameworks. To posit information relating structure to physical properties in our latent space, SmVAE has a property prediction component and is jointly trained for property prediction and framework generation. Since the size of our property dataset is much smaller than our structural dataset, we train this component in a semi-supervised fashion. In the joint training, SmVAE was fed with 40k MOFs with the property data (textural and gas uptake properties) and another  $\sim 2$  million MOFs without property data. Predictive network parameters are only optimized when labeled data is observed during training.<sup>67</sup> When the model is correctly trained, we can identify principal axes that align with increasing and decreasing values of physical properties. This feature improves the optimization capabilities of our model. Gas separation selectivities are then derived using the corresponding uptake values of the gas phases. Taking all the components into account, we propose a multi-component loss function  $L_{total}$  as follows:

$$\begin{aligned} L_{total} &= L_{edge} + L_{vertex} + L_{topo} + L_{property} + L_{KL} \\ &= L_{RFcodeRecon.} + L_{Semi-superProp} + L_{VAEConstraint} \end{aligned} \tag{1}$$

After the realization of property prediction, we ultimately add one property-guided optimization component to the SmVAE for automated reticular framework inverse design. A Gaussian process model is built and trained with labeled frameworks from the jointly-trained latent space from the SmVAE to predict the targeted properties. The entire structure of our SmVAE with all components is illustrated in **Fig. 2**. Gaussian process models are known to be effective in prediction with even a limited amount of training data.<sup>68</sup> The detailed SmVAE architecture and hyperparameter tuning process are described in **Supplementary Note 2**.

### **Demonstration of SmVAE on MOF design and optimization**

To evaluate the fidelity of the trained SmVAE and the capability of its latent space to capture MOF structure information, we estimate the kernel density of each dimension in the latent space (288 in total). As shown in **Fig. S7**, all data distributions in different dimensions are normal, indicating the effectiveness of the variational regularizer as implemented in SmVAE. Furthermore, we use MOF-117<sup>60</sup> as an example by feeding its RFcode to the encoder to obtain its latent representation and sampling its neighbouring latent points at various distances. We check the decoding results of the original representation and neighbouring points. We are able to get the original MOF-117 back at the original point, and decoded MOF structures at the sampled neighbouring points demonstrate more and more variations with increasing distance as shown in **Fig. 3C**. The autoencoder also provides us a critical opportunity to explore the geometrical correlation between different MOF structures. We encode two well known yet topologically distinct MOF structures (*i.e.*, cubic, **ftw** NU-1104<sup>69</sup> and hexagonal, **csq** NU-1000<sup>70</sup>) and perform an interpolation between their latent points in space (**Fig. 3D**). The intermediate frameworks along the interpolation path are then decoded, which demonstrate a clear geometrical evolution from the cubic framework to the hexagonal framework.

Discovering systems with improved properties is the essential goal of materials design. We examine the mapping of property values to the latent representation in the jointly trained SmVAE latent space using PCA (**Fig. 3A and 3B**), and we find that the distribution of frameworks shows an explicit gradient, with high performance MOFs located in one domain and low performance MOFs in other domains. For comparison, another SmVAE was trained with ~2 million MOFs without any property as a control group. The resulting latent representation distribution shows no noticeable pattern with respect to property values (**Fig. S8**), confirming the ability of SmVAE to organize the latent space according to property values. Performance metrics like prior and posterior scores for sampling and constructing valid MOFs, as well as the mean absolute error (MAE) on predicting MOF properties, are computed and shown in **Tab. S2**. Our SmVAE demonstrates superb accuracy in designing MOFs and predicting their properties.

Ultimately, we optimize MOF structures in the latent space of the jointly trained autoencoder. We build a Gaussian Process (GP) model, which has been proven to be lightweight and effective for smooth function prediction,<sup>68</sup> to learn the property landscape of the latent representations. A GP model is then trained to predict the target property of the latent vector of a given MOF RFcode. We then choose CO<sub>2</sub> uptake in the natural gas separation (CO<sub>2</sub>/CH<sub>4</sub>) as the



target and demonstrate two optimization processes: 1) isorecticular MOF design, where the topology is constrained and 2) globally optimized MOF design (**Fig. 4A**), with maximized property frameworks identified and intermediate structures interpolated. In the isorecticular design process, we pick the MOF NU-1104<sup>69</sup> (CO<sub>2</sub> uptake of 0.65 mol/kg) as the starting point and optimized the framework with constrained **ftw** topology. Going through a series of intermediate linkers (**Fig. 4B**), we are able to optimize the targeted CO<sub>2</sub> uptake to 4.33 mol/kg. In the global optimization process without topology constraint, we begin with MOF-5<sup>71</sup> (CO<sub>2</sub> uptake of 2.80 mol/kg<sup>72</sup>) and search for MOFs with optimized uptake (**Fig. 4C**). At the end, we discover a **spn** MOF with a remarkably high CO<sub>2</sub> uptake of 7.55 mol/kg for natural gas separation.

### Top MOF candidates proposed for gas separation

Aiming at CO<sub>2</sub> loading in natural gas (5 bar, 300 K, 10/90 CO<sub>2</sub>/CH<sub>4</sub>) and flue gas separation (1 bar, 313 K, 15/85 CO<sub>2</sub>/N<sub>2</sub>), we repeat the globally optimized design process and select the top candidates for further validations. When we rank all the generatively designed MOFs (GMOFs), we consider their gas separation properties as well as the MOF synthesizability to make suggestions for further experimental measurements. To estimate the latter, we calculate the synthetic complexity score (SCScore)<sup>73</sup> of the organic linkers used in the GMOFs. Complete linkers of all MOF candidates are assembled using the appropriate edge and organic vertex, as shown in **Fig. S9**. The top candidates with superior performance are shown in **Tab. 1A** sorted according to decreasing SCScore. We are able to identify multiple MOFs with enhanced gas separation properties, including GMOF-9, which exhibits the highest 7.55 mol/kg CO<sub>2</sub> uptake and reasonably large selectivity of 16.0 for CO<sub>2</sub>/CH<sub>4</sub> separation. All candidate MOF structures are stable through relaxation, and corresponding properties predicted have been reconfirmed with grand canonical Monte Carlo (GCMC) simulations (**Fig. S10**).

By examining the corresponding porosities, we identify two types of promising MOFs with distinct gas separation mechanisms:

- 1) **Size exclusion frameworks (GMOFs-1,2,3,4)** with small pore-limiting diameter (PLD: 3.40~3.71 Å) that fall between CO<sub>2</sub> (3.3 Å) and CH<sub>4</sub> (3.8 Å) or N<sub>2</sub> (3.64 Å), therefore effectively permitting CO<sub>2</sub> to diffuse into the MOF while excluding CH<sub>4</sub> or N<sub>2</sub>. These MOFs have very high or infinite selectivity for CO<sub>2</sub> because they absorb little or no CH<sub>4</sub> or N<sub>2</sub>. For CO<sub>2</sub>/CH<sub>4</sub> separation, they exhibit remarkable CO<sub>2</sub> uptakes (4.64, 4.33, 4.22, 3.97 mol/kg). They are also strong CO<sub>2</sub>/N<sub>2</sub> separation candidates with high CO<sub>2</sub> uptakes (3.06, 2.09, 2.47, 2.51 mol/kg) and high selectivities.
- 2) **Thermodynamic separation frameworks (GMOFs-5,6,7,8,9)**, which show large pores (LCD: 62.61~81.08 Å, PLD: 55.99~68.92 Å), compared to the size of the targeted molecules. These MOFs all have high gravimetric surface areas (> 5000 m<sup>2</sup>/g), offering many binding sites for CO<sub>2</sub>, which results in high capacity. They exhibit strong selective CO<sub>2</sub> adsorption as a result of the stronger van der Waals interactions between CO<sub>2</sub> and the frameworks versus CH<sub>4</sub> and N<sub>2</sub>. For CO<sub>2</sub>/CH<sub>4</sub> separation, we observe notably high CO<sub>2</sub> uptakes (4.80, 4.51, 4.34, 7.21, 7.55 mol/kg) at reasonably high CO<sub>2</sub>/CH<sub>4</sub> selectivities (10.0~17.5). They are also competent flue gas separation

materials with reasonable CO<sub>2</sub> uptakes (1.29, 1.26, 1.45, 2.61, 2.80 mol/kg) and good CO<sub>2</sub>/N<sub>2</sub> selectivities (11.7~27.1).

Performance comparison on gas separations between MOFs is practically difficult since the measurements are often conducted at different experimental conditions (*e.g.*, temperature, pressure, and gas phase composition). However, we believe MOFs discovered here are strongly competitive against some of the best-performing MOFs/zeolites ever reported in the literature (**Tab. 1B**). For natural gas separation, the notable MgMOF-74 and zeolite 13X show CO<sub>2</sub> capacities comparable with (8.0 mol/kg)<sup>74</sup> or even lower than (4.4 mol/kg)<sup>74</sup> our top candidates (*i.e.*, GMOF-8: 7.21 mol/kg, GMOF-9: 7.55 mol/kg) at similar conditions (~5bar, ~300K) while SIFSIX-2-Cu-i, SIFSIX-3-Zn, and UTSA-16 exhibit capacities of 4.16,<sup>9</sup> 2.46,<sup>9</sup> 4.25<sup>5</sup> mol/kg at lower pressure conditions (1~2bar, ~300K). Their selectivities against CH<sub>4</sub> (29.8~231)<sup>5,9,74</sup> are all lower than our top selectivity candidates (*i.e.*, GMOFs-1,2,3,4: 3058 ~ ∞ with zero CH<sub>4</sub> uptake). For flue gas separation at similar conditions as this study (1 bar, 313 K, 15/85 CO<sub>2</sub>/N<sub>2</sub>), our top candidate GMOF-1 exhibits CO<sub>2</sub> uptake of 3.06 mol/kg with extremely high selectivity (∞ with zero N<sub>2</sub> uptake), which is only lower than the capacity of MgMOF-74: 4.43 mol/kg with a selectivity of 175 (0.9 bar, 313 K, 15/75 CO<sub>2</sub>/N<sub>2</sub>),<sup>75</sup> while higher than SIFSIX-2-Cu-i (1.59 mol/kg at 140 selectivity),<sup>9</sup> SIFSIX-3-Zn (2.27 mol/kg at 1818 selectivity),<sup>9</sup> UTSA-16 (2.37 mol/kg at 314.7 selectivity),<sup>5</sup> and 13X (3.0 mol/kg at 20 selectivity).<sup>76</sup> Furthermore, our top candidates show potentially strong chemical and hydrothermal stabilities, with the exclusive usage of well-known stable metal nodes like Zr<sub>6</sub>O<sub>8</sub>-/Zr<sub>6</sub>O<sub>20</sub>-, Cr<sub>3</sub>O<sub>4</sub>-, and ZnN<sub>4</sub>-.<sup>77</sup> This is particularly important for carbon capture applications in a harsh flue gas environment.<sup>78,79</sup>

## Conclusions

We develop an automated nanoporous materials discovery platform using a supramolecular variational autoencoder for the generation of reticular frameworks with optimized properties. We have demonstrated the automated design process with MOF structures starting from the computation-ready, experimental (CoRE) MOF database<sup>45</sup> and generating new proposed structures with improved properties for CO<sub>2</sub> separations. Our model exhibits high fidelity in capturing structural features and reconstructing MOF structures. The autoencoder shows great prediction and optimization capability when jointly trained with multiple top candidates identified for superior gas separation and confirmed *via* atomistic Monte Carlo simulations. We use this platform to design novel MOFs with improved capacity and good selectivity for CO<sub>2</sub>/N<sub>2</sub> and CO<sub>2</sub>/CH<sub>4</sub> separations, which are important clean energy-relevant applications. The top performing MOF has a CO<sub>2</sub> capacity of 7.55 mol/kg and a selectivity over CH<sub>4</sub> of 16. This platform can be applied to a broad range of materials (*e.g.*, covalent-organic frameworks, metal-organic polyhedra, hydrogen-bonded organic framework, coordinational polymers) and lays the groundwork for the design of reticular frameworks for a variety of applications.

## Reference

1. Yaghi, O. M. *et al.* Reticular synthesis and the design of new materials. *Nature* **423**, 705–714 (2003).
2. Li, H., Eddaoudi, M., Groy, T. L. & Yaghi, O. M. Establishing microporosity in open metal-organic frameworks: Gas sorption isotherms for Zn(BDC) (BDC = 1,4-benzenedicarboxylate). *Journal of the American Chemical Society* **120**, 8571–8572 (1998).
3. Farha, O. K. *et al.* De novo synthesis of a metal-organic framework material featuring ultrahigh surface area and gas storage capacities. *Nat. Chem.* **2**, 944–948 (2010).
4. Mason, J. A. *et al.* Methane storage in flexible metal-organic frameworks with intrinsic thermal management. *Nature* **527**, 357–361 (2015).
5. Xiang, S. *et al.* Microporous metal-organic framework with potential for carbon dioxide capture at ambient conditions. *Nat. Commun.* **3**, (2012).
6. Rodenas, T. *et al.* Metal-organic framework nanosheets in polymer composite materials for gas separation. *Nat. Mater.* **14**, 48–55 (2015).
7. Yang, S. *et al.* Supramolecular binding and separation of hydrocarbons within a functionalized porous metal-organic framework. *Nat. Chem.* **7**, 121–129 (2015).
8. Chen, K.-J. *et al.* Synergistic sorbent separation for one-step ethylene purification from a four-component mixture. *Science (80-. )*. **366**, 241–246 (2019).
9. Nugent, P. *et al.* Porous materials with optimal adsorption thermodynamics and kinetics for CO<sub>2</sub> separation. *Nature* **495**, 80–84 (2013).
10. Diercks, C. S., Liu, Y., Cordova, K. E. & Yaghi, O. M. The role of reticular chemistry in the design of CO<sub>2</sub> reduction catalysts. *Nat. Mater.* **17**, 301–307 (2018).
11. Mondloch, J. E. *et al.* Destruction of chemical warfare agents using metal-organic frameworks. *Nat. Mater.* **14**, 512–516 (2015).
12. Hu, Z., Deibert, B. J. & Li, J. Luminescent metal-organic frameworks for chemical sensing and explosive detection. *Chemical Society Reviews* **43**, 5815–5840 (2014).
13. Sheberla, D. *et al.* Conductive MOF electrodes for stable supercapacitors with high areal capacitance. *Nat. Mater.* **16**, 220–224 (2017).
14. Shimizu, G. K. H., Taylor, J. M. & Kim, S. Proton conduction with metal-organic frameworks. *Science* **341**, 354–355 (2013).
15. Baumann, A. E., Burns, D. A., Liu, B. & Thoi, V. S. Metal-organic framework functionalization and design strategies for advanced electrochemical energy storage devices. *Communications Chemistry* **2**, (2019).
16. Tan, L. L. *et al.* Stimuli-responsive metal-organic frameworks gated by pillar[5]arene supramolecular switches. *Chem. Sci.* **6**, 1640–1644 (2015).
17. Li, M., Li, D., O’Keeffe, M. & Yaghi, O. M. Topological analysis of metal-organic frameworks with polytopic linkers and/or multiple building units and the minimal transitivity principle. *Chemical Reviews* **114**, 1343–1370 (2014).
18. Kirkpatrick, P. & Ellis, C. Chemical space. *Nature* **432**, 823 (2004).
19. Reymond, J. L. The Chemical Space Project. *Acc. Chem. Res.* **48**, 722–730 (2015).
20. Wilmer, C. E. *et al.* Large-scale screening of hypothetical metal-organic frameworks. *Nat. Chem.* **4**, 83–89 (2012).
21. Collins, S. P., Daff, T. D., Piotrkowski, S. S. & Woo, T. K. Materials design by evolutionary optimization of functional groups in metal-organic frameworks. *Sci. Adv.* **2**,

- (2016).
22. Boyd, P. G., Lee, Y. & Smit, B. Computational development of the nanoporous materials genome. *Nat. Rev. Mater.* **2**, 17037 (2017).
  23. Yazaydin, A. Ö. *et al.* Screening of metal-organic frameworks for carbon dioxide capture from flue gas using a combined experimental and modeling approach. *J. Am. Chem. Soc.* **131**, 18198–18199 (2009).
  24. Sumida, K. *et al.* Hydrogen storage and carbon dioxide capture in an iron-based sodalite-type metal-organic framework (Fe-BTT) discovered via high-throughput methods. *Chem. Sci.* **1**, 184–191 (2010).
  25. Sonnauer, A. *et al.* Giant pores in a chromium 2,6-Naphthalenedicarboxylate Open-Framework structure with MIL-101 topology. *Angew. Chemie - Int. Ed.* **48**, 3791–3794 (2009).
  26. Boyd, P. G. *et al.* Data-driven design of metal–organic frameworks for wet flue gas CO<sub>2</sub> capture. *Nature* **576**, 253–256 (2019).
  27. Chung, Y. G. *et al.* In silico discovery of metal-organic frameworks for precombustion CO<sub>2</sub> capture using a genetic algorithm. *Sci. Adv.* **2**, e1600909 (2016).
  28. Gustafson, J. A. & Wilmer, C. E. Intelligent Selection of Metal–Organic Framework Arrays for Methane Sensing via Genetic Algorithms. *ACS Sensors* **4**, 1586–1593 (2019).
  29. Bao, Y. *et al.* In Silico Discovery of High Deliverable Capacity Metal–Organic Frameworks. *J. Phys. Chem. C* **119**, 186–195 (2015).
  30. Fernandez, M., Boyd, P. G., Daff, T. D., Aghaji, M. Z. & Woo, T. K. Rapid and Accurate Machine Learning Recognition of High Performing Metal Organic Frameworks for CO<sub>2</sub> Capture. *J. Phys. Chem. Lett.* **5**, 3056–3060 (2014).
  31. Fanourgakis, G. S., Gkagkas, K., Tylianakis, E. & Froudakis, G. E. A Universal Machine Learning Algorithm for Large-Scale Screening of Materials. *J. Am. Chem. Soc.* **142**, 3814–3822 (2020).
  32. Moghadam, P. Z. *et al.* Structure-Mechanical Stability Relations of Metal-Organic Frameworks via Machine Learning. *Matter* **1**, 219–234 (2019).
  33. Raccuglia, P. *et al.* Machine-learning-assisted materials discovery using failed experiments. *Nature* **533**, 73–76 (2016).
  34. Shi, Z. *et al.* Machine-learning-assisted high-throughput computational screening of high performance metal–organic frameworks. *Mol. Syst. Des. Eng.* (2020). doi:10.1039/d0me00005a
  35. Kingma, D. P. & Welling, M. Auto-encoding variational bayes. in *2nd International Conference on Learning Representations, ICLR 2014 - Conference Track Proceedings* (International Conference on Learning Representations, ICLR, 2014).
  36. Goodfellow, I. J. *et al.* Generative Adversarial Networks. (2014).
  37. Sanchez-Lengeling, B. & Aspuru-Guzik, A. Inverse molecular design using machine learning: Generative models for matter engineering. *Science (80-. )*. **361**, 360–365 (2018).
  38. Gómez-Bombarelli, R. *et al.* Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules. *ACS Cent. Sci.* **4**, 268–276 (2018).
  39. You, J., Liu, B., Ying, R., Pande, V. & Leskovec, J. Graph Convolutional Policy Network for Goal-Directed Molecular Graph Generation. *Adv. Neural Inf. Process. Syst.* **31** (eds. Bengio, S. *al.*) 6410–6421 (2018).
  40. Liu, Q., Allamanis, M., Brockschmidt, M. & Gaunt, A. L. Constrained graph variational autoencoders for molecule design. in *Advances in Neural Information Processing Systems*

- 2018-December**, 7795–7804 (Neural information processing systems foundation, 2018).
41. Jin, W., Barzilay, R. & Jaakkola, T. Junction Tree Variational Autoencoder for Molecular Graph Generation. *35th Int. Conf. Mach. Learn. ICML 2018* **5**, 3632–3648 (2018).
  42. Noh, J. *et al.* Inverse Design of Solid-State Materials via a Continuous Representation. *Matter* **1**, 1370–1384 (2019).
  43. Li, X. *et al.* A Deep Adversarial Learning Methodology for Designing Microstructural Material Systems. in (ASME International, 2018). doi:10.1115/detc2018-85633
  44. Noura, A., Sokolovska, N. & Crivello, J.-C. CrystalGAN: Learning to Discover Crystallographic Structures with Generative Adversarial Networks. (2018).
  45. Chung, Y. G. *et al.* Advances, Updates, and Analytics for the Computation-Ready, Experimental Metal–Organic Framework Database: CoRE MOF 2019. *J. Chem. Eng. Data* acs.jced.9b00835 (2019). doi:10.1021/acs.jced.9b00835
  46. Duvenaud, D. *et al.* Convolutional networks on graphs for learning molecular fingerprints. in *Advances in Neural Information Processing Systems 2015-January*, 2224–2232 (Neural information processing systems foundation, 2015).
  47. Krenn, M., Häse, F., Nigam, A., Friederich, P. & Aspuru-Guzik, A. SELFIES: a robust representation of semantically constrained graphs with an example application in chemistry. (2019).
  48. Li, P. *et al.* Bottom-up construction of a superstructure in a porous uranium-organic crystal. *Science (80-. )*. **356**, 624–627 (2017).
  49. Jain, A. *et al.* Commentary: The Materials Project: A materials genome approach to accelerating materials innovation. *APL Mater.* **1**, 011002 (2013).
  50. Bartók, A. P., Kondor, R. & Csányi, G. On representing chemical environments. *Phys. Rev. B - Condens. Matter Mater. Phys.* **87**, (2013).
  51. Isayev, O. *et al.* Universal fragment descriptors for predicting properties of inorganic crystals. *Nature Communications* **8**, (2017).
  52. Ward, L. *et al.* Including crystal structure attributes in machine learning models of formation energies via Voronoi tessellations. *Phys. Rev. B* **96**, (2017).
  53. Ziletti, A., Kumar, D., Scheffler, M. & Ghiringhelli, L. M. Insightful classification of crystal structures using deep learning. *Nat. Commun.* **9**, (2018).
  54. Ryan, K., Lengyel, J. & Shatruk, M. Crystal Structure Prediction via Deep Learning. *J. Am. Chem. Soc.* **140**, 10158–10168 (2018).
  55. Xie, T. & Grossman, J. C. Crystal Graph Convolutional Neural Networks for an Accurate and Interpretable Prediction of Material Properties. *Phys. Rev. Lett.* **120**, (2018).
  56. Park, C. W. & Wolverton, C. Developing an improved Crystal Graph Convolutional Neural Network framework for accelerated materials discovery. (2019).
  57. Eon, J. G. Topological features in crystal structures: A quotient graph assisted analysis of underlying nets and their embeddings. *Acta Crystallogr. Sect. A Found. Adv.* **72**, 268–293 (2016).
  58. Delgado-Friedrichs, O., Hyde, S. T., O’Keeffe, M. & Yaghi, O. M. Crystal structures as periodic graphs: the topological genome and graph databases. *Structural Chemistry* **28**, 39–44 (2017).
  59. O’Keeffe, M. & Yaghi, O. M. Deconstructing the crystal structures of metal-organic frameworks and related materials into their underlying nets. *Chemical Reviews* **112**, 675–702 (2012).
  60. Furukawa, H., Kim, J., Ockwig, N. W., O’Keeffe, M. & Yaghi, O. M. Control of vertex

- geometry, structure dimensionality, functionality, and pore metrics in the reticular synthesis of crystalline metal-organic frameworks and polyhedra. *J. Am. Chem. Soc.* **130**, 11650–11661 (2008).
61. Bucior, B. J. *et al.* Identification Schemes for Metal–Organic Frameworks To Enable Rapid Search and Cheminformatics Analysis. *Cryst. Growth Des.* **19**, 6682–6697 (2019).
  62. Colón, Y. J., Gómez-Gualdrón, D. A. & Snurr, R. Q. Topologically Guided, Automated Construction of Metal–Organic Frameworks and Their Evaluation for Energy-Related Applications. *Cryst. Growth Des.* **17**, 5801–5810 (2017).
  63. Anderson, R. & Gómez-Gualdrón, D. A. Increasing topological diversity during computational “synthesis” of porous crystals: how and why. *CrystEngComm* **21**, 1653–1665 (2019).
  64. Ghersi, D. & Singh, M. molBLOCKS: decomposing small molecule sets and uncovering enriched fragments. *Bioinformatics* **30**, 2081–2083 (2014).
  65. Kingma, D. P. & Welling, M. An Introduction to Variational Autoencoders. *Found. Trends® Mach. Learn.* **12**, 307–392 (2019).
  66. Fu, H. *et al.* Cyclical Annealing Schedule: A Simple Approach to Mitigating KL Vanishing. (2019).
  67. Kingma, D. P., Rezende, D. J., Mohamed, S. & Welling, M. Semi-Supervised Learning with Deep Generative Models. *Adv. Neural Inf. Process. Syst.* **4**, 3581–3589 (2014).
  68. Rasmussen, C. E. & Williams, C. K. I. *Gaussian Processes for Machine Learning*. (The MIT Press, 2005). doi:10.7551/mitpress/3206.001.0001
  69. Deria, P. *et al.* Ultraporous, water stable, and breathing zirconium-based metal-organic frameworks with ftw topology. *J. Am. Chem. Soc.* **137**, 13183–13190 (2015).
  70. Mondloch, J. E. *et al.* Vapor-phase metalation by atomic layer deposition in a metal-organic framework. *J. Am. Chem. Soc.* **135**, 10294–10297 (2013).
  71. Li, H., Eddaoudi, M., O’Keeffe, M. & Yaghi, O. M. Design and synthesis of an exceptionally stable and highly porous metal- organic framework. *Nature* **402**, 276–279 (1999).
  72. Gu, Z. Y., Jiang, J. Q. & Yan, X. P. Fabrication of isorecticular metal-organic framework coated capillary columns for high-resolution gas chromatographic separation of persistent organic pollutants. *Anal. Chem.* **83**, 5093–5100 (2011).
  73. Coley, C. W., Rogers, L., Green, W. H. & Jensen, K. F. SCScore: Synthetic Complexity Learned from a Reaction Corpus. *J. Chem. Inf. Model.* **58**, 252–261 (2018).
  74. Herm, Z. R., Krishna, R. & Long, J. R. CO<sub>2</sub>/CH<sub>4</sub>, CH<sub>4</sub>/H<sub>2</sub> and CO<sub>2</sub>/CH<sub>4</sub>/H<sub>2</sub> separations at high pressures using Mg<sub>2</sub>(dobdc). *Microporous Mesoporous Mater.* **151**, 481–487 (2012).
  75. Mason, J. A., Sumida, K., Herm, Z. R., Krishna, R. & Long, J. R. Evaluating metal-organic frameworks for post-combustion carbon dioxide capture via temperature swing adsorption. *Energy Environ. Sci.* **4**, 3030–3040 (2011).
  76. Cavenati, S., Grande, C. A. & Rodrigues, A. E. Adsorption Equilibrium of Methane, Carbon Dioxide, and Nitrogen on Zeolite 13X at High Pressures. *J. Chem. Eng. Data* **49**, 1095–1101 (2004).
  77. Howarth, A. J. *et al.* Chemical, thermal and mechanical stabilities of metal-organic frameworks. *Nature Reviews Materials* **1**, 1–15 (2016).
  78. Mangano, E., Kahr, J., Wright, P. A. & Brandani, S. Accelerated degradation of MOFs under flue gas conditions. *Faraday Discussions* **192**, 181–195 (2016).

79. Rieth, A. J., Wright, A. M. & Dincă, M. Kinetic stability of metal–organic frameworks for corrosive and coordinating gas capture. *Nature Reviews Materials* **4**, 708–725 (2019).

### **Acknowledgments**

Z.Y., N.S.B., B.B., S.G.H.K., O.K.F, R.S.Q. and A.A.-G. were supported as part of the Nanoporous Materials Genome Center by the U.S. Department of Energy, Office of Science, Office of Basic Energy Sciences under award number DE-FG02-17ER16362. Funding for T.B., S.P.C. and T.K.W. were provided by NSERC. Computations were made on the supercomputer “beluga” from École de technologie supérieure, managed by Calcul Québec and Compute Canada. The operation of this supercomputer is funded by the Canada Foundation for Innovation (CFI), the ministère de l'Économie, de la science et de l'innovation du Québec (MESI) and the Fonds de recherche du Québec - Nature et technologies (FRQ-NT). This research was supported in part through the computational resources and staff contributions provided for the Quest high performance computing facility at Northwestern University which is jointly supported by the Office of the Provost, the Office for Research, and Northwestern University Information Technology.

### **Author contributions:**

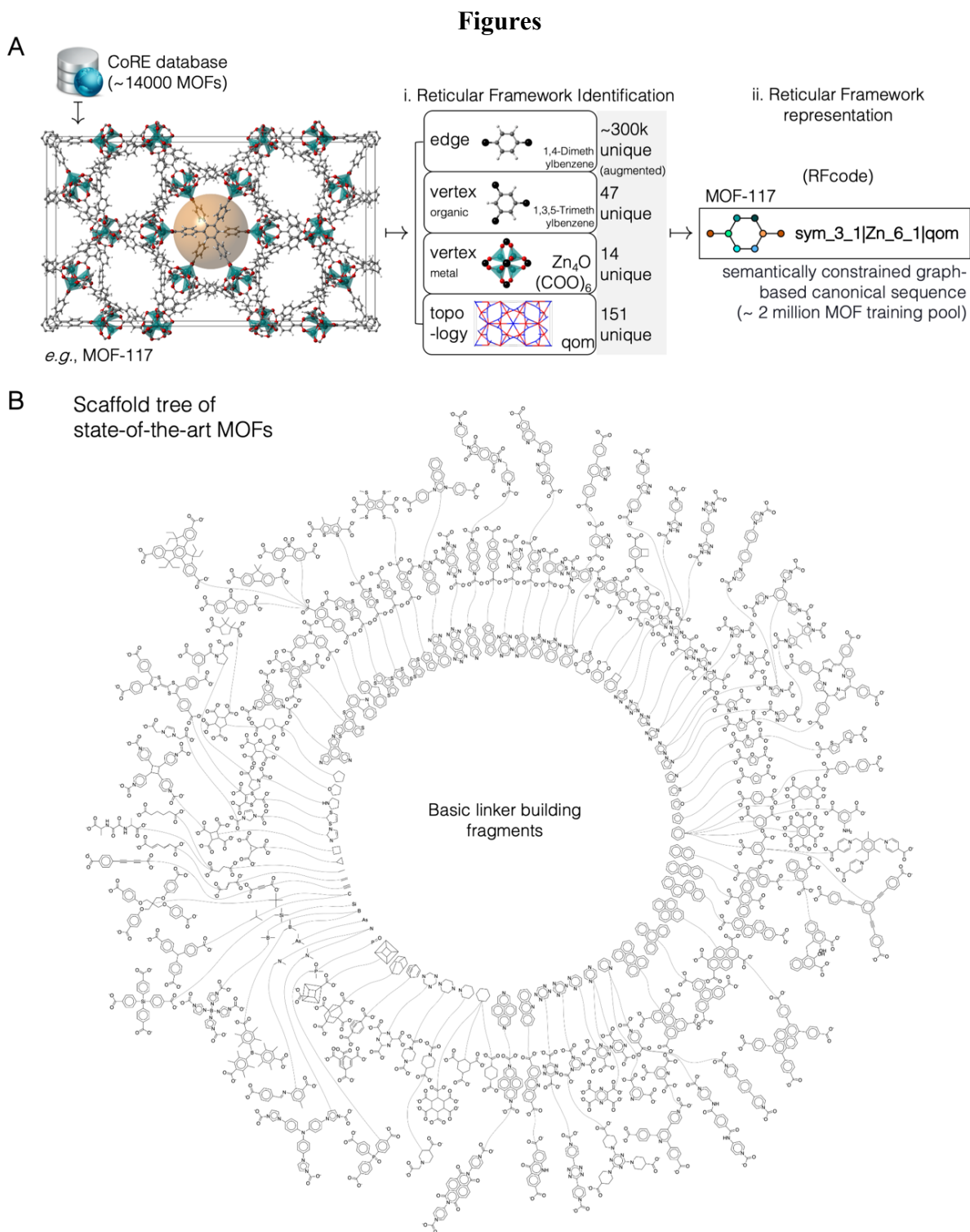
Z.Y. conceived the overall project. Z.Y. B.J.B. and R.Q.S. designed the reticular framework representation approach. N.S.B. and R.Q.S. conducted the MOF property determination calculations. Z.Y. and B.S-L. developed the deep learning variational autoencoder. S.P.C., T.B., and T.K.W. did the charge calculations for the framework charges for property simulations. A.A.-G. led the project and provided the overall directions. All authors participated in preparing the manuscript.

### **Additional information**

Supplementary information is available in the online version of the paper. Reprints and permissions information is available online at [www.nature.com/reprints](http://www.nature.com/reprints). Correspondence and requests for materials should be addressed to A.A.-G., R.Q.S., and Z.Y..

### **Competing financial interests**

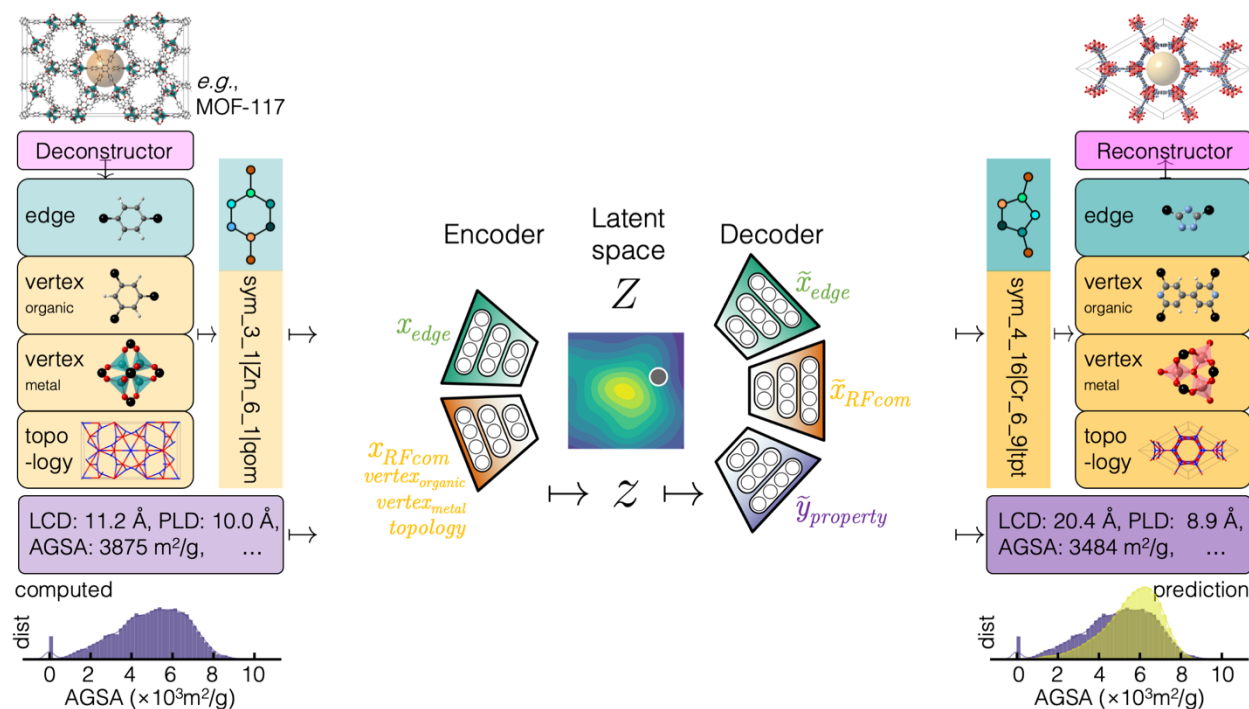
O.K.F. and R.Q.S. have a financial interest in NuMat Technologies, a startup company that is seeking to commercialize MOFs.



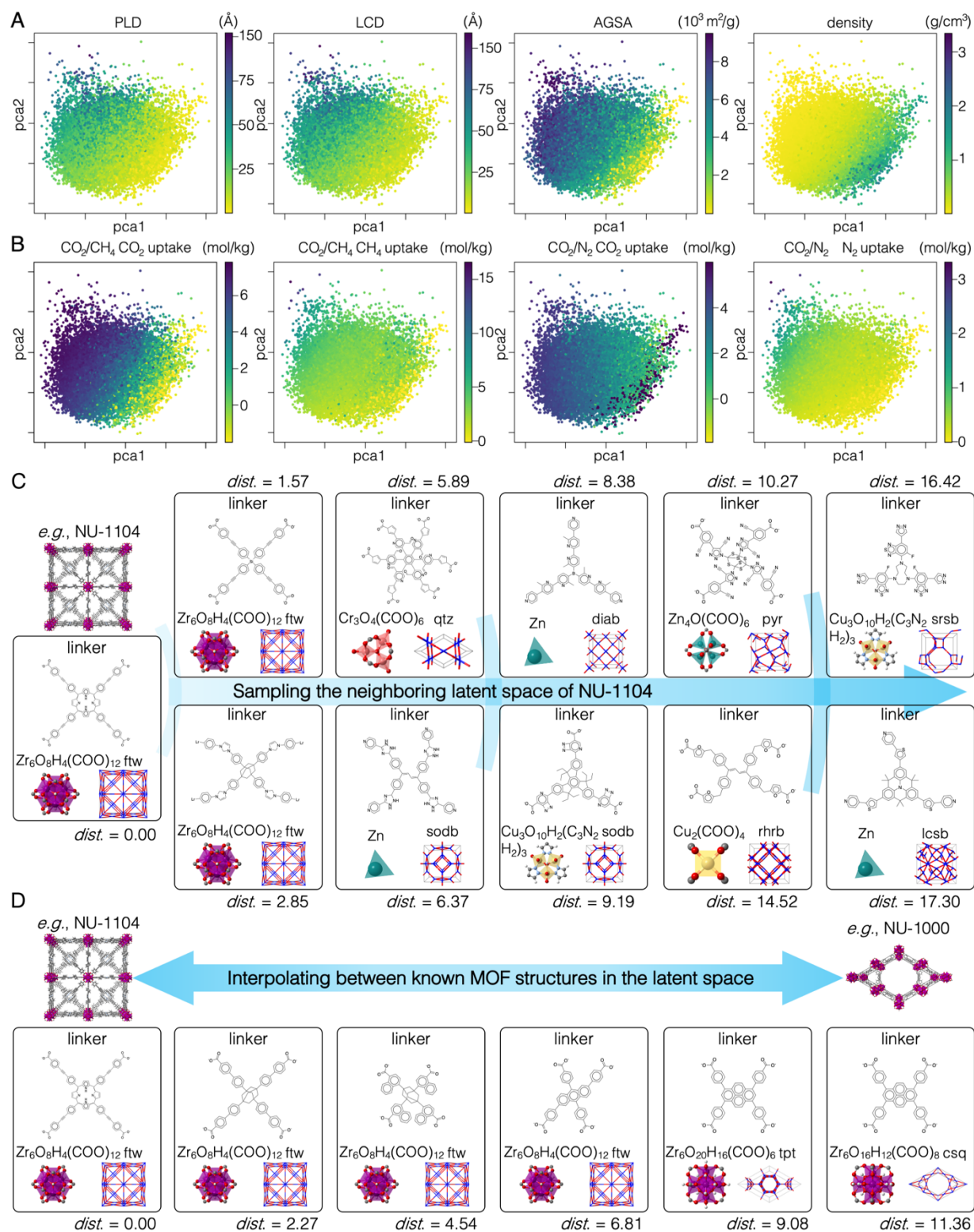
**Fig. 1 | Reticular framework identification and representation, exemplified with MOF structures from the CoRE database.**<sup>45</sup> (A) Reticular frameworks (*e.g.*, MOF-117)<sup>60</sup> are: (i) decomposed to their building blocks (edges, organic/metal vertices) and topology using a previously developed identification method,<sup>61</sup> which are then recognized and labeled; (ii) the labels are further converted to semantically constrained graph-based canonical sequences, namely the RFcode. (B) A fragmentation analysis was



conducted on the linkers of all MOF structures from the CoRE MOF database. Here we illustrate all the basic building fragments of state-of-the-art MOFs with high occurrence rate (the inner circle) and linkers derived from them (the second and third circles) in a scaffold tree plot.

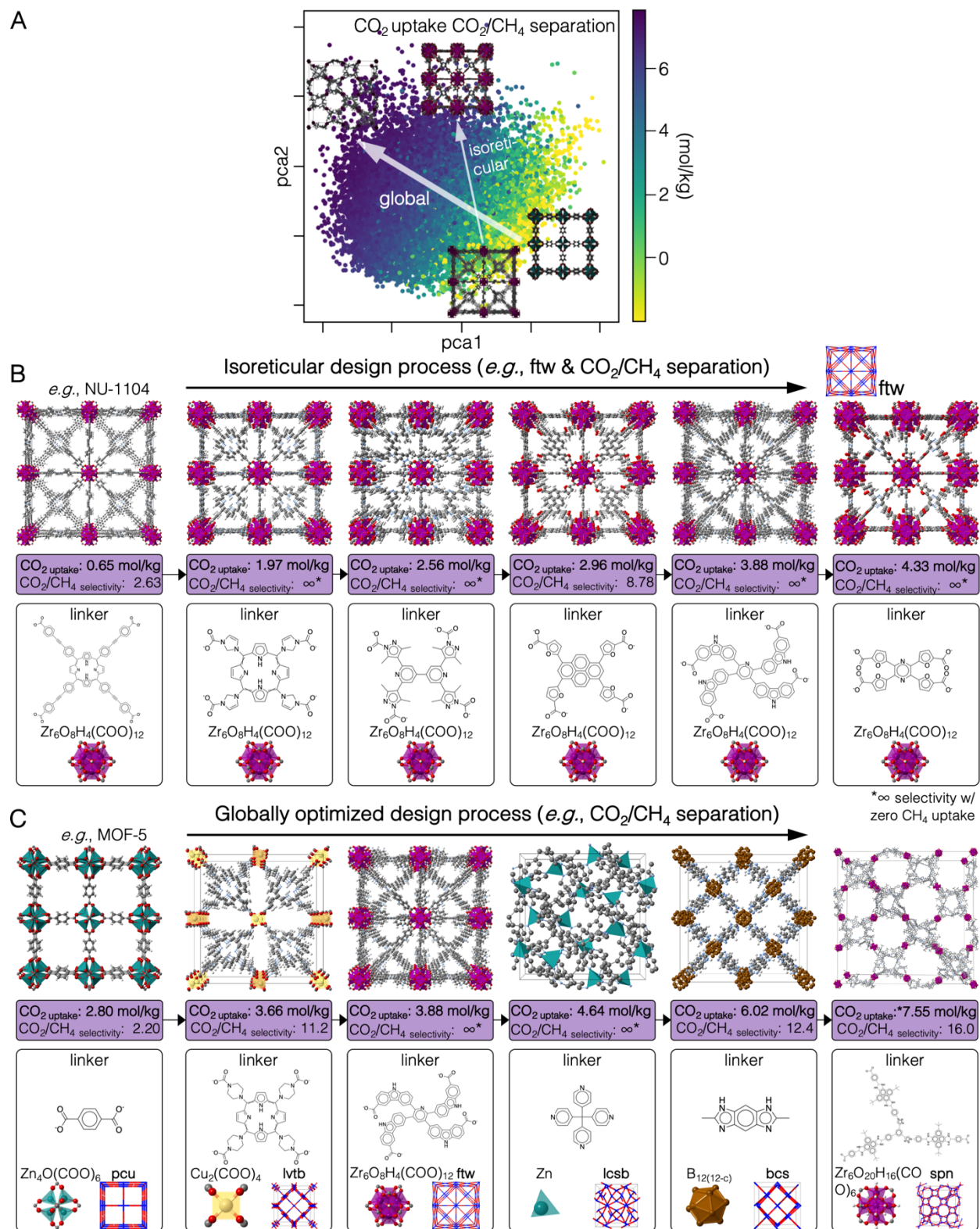


**Fig. 2 | Schematic diagram of the automated reticular framework discovery platform empowered by the supramolecular variational autoencoder (SmVAE).** SmVAE is a multi-component variational autoencoder with modules that are in charge of encoding and decoding each part of the RFCODE ( $\mathbf{x}_{edge} \rightarrow \tilde{\mathbf{x}}_{edge}$ ,  $\mathbf{x}_{RFcom} \rightarrow \tilde{\mathbf{x}}_{RFcom}$ ). Reticular frameworks are mapped with discrete RFCodes, transferred into continuous vectors ( $\mathbf{Z}$ ), and then transferred back. To have the latent space organized around properties of interest, we add an extra component to the model that uses labeled data ( $\mathcal{Y}$ ). This process is realized with the additional model that learns to predict properties ( $\tilde{\mathcal{Y}}_{property}$ ) from the latent space.



**Fig. 3 | Illustration of the latent space of the jointly trained SmVAE using PCA analysis conditioned by MOF properties and exemplified sampling of the latent space.** Latent space of the SmVAE after joint training exhibits notable gradients by A) textural and B) gas uptake property values. C) Starting with NU-1104,<sup>69</sup> we sample its neighbouring latent points at various distances and check the decoding results.

D) We interpolate the latent points of two known distinct MOF structures (*e.g.*, NU-1104 and NU-1000<sup>70</sup>) and identify the intermediate structures. A clear structure evolution is observed from two geometrically different frameworks.



**Fig. 4 | Reticular framework design and optimization using SmVAE with natural gas separation (CO<sub>2</sub> uptake) as the exemplified target. Two optimized design processes of isorecticular and global are demonstrated in B and C with the optimization paths in the latent space schematically shown in A. Through**

the isoreticular design process (starting with MOF NU-1104<sup>69</sup>) constrained to the **ftw** topology, a MOF with notable CO<sub>2</sub> uptake of 4.33 mol/kg and infinite selectivity (zero CH<sub>4</sub> uptake) is discovered (5 bar, 313 K, 10/90 CO<sub>2</sub>/CH<sub>4</sub>). The global optimization design process without topology constraint (starting with MOF-5<sup>71</sup>) leads to a MOF with the remarkably high CO<sub>2</sub> uptake of 7.55 mol/kg and high selectivity of 16.0 (5 bar, 313 K, 10/90 CO<sub>2</sub>/CH<sub>4</sub>).

**Tab. 1 | A. Generatively designed top MOF candidates targeted at gas separations (natural gas: CO<sub>2</sub>/CH<sub>4</sub> and flue gas: CO<sub>2</sub>/N<sub>2</sub>) sorted with decreasing synthesizability SCScore. B. Gas separation performance of well-known MOFs and zeolites.**

**A**

	GMOF-1	GMOF-2	GMOF-3	GMOF-4	GMOF-5	GMOF-6	GMOF-7	GMOF-8	GMOF-9
Topology	lcsb	ftw	ftw	tpt	spn	spn	spn	spn	spn
Linker									
(SCScore)	1.6	3.2	3.5	3.8	4.8	4.9	4.9	4.9	5.0
Metal node									
	Zn	Zr <sub>6</sub> O <sub>8</sub> H <sub>4</sub> (COO) <sub>12</sub>	Zr <sub>6</sub> O <sub>8</sub> H <sub>4</sub> (COO) <sub>12</sub>	Cr <sub>3</sub> O <sub>4</sub> (COO) <sub>6</sub>	Zr <sub>6</sub> O <sub>20</sub> H <sub>16</sub> (COO) <sub>6</sub>	Zr <sub>6</sub> O <sub>20</sub> H <sub>16</sub> (COO) <sub>6</sub>	Zr <sub>6</sub> O <sub>20</sub> H <sub>16</sub> (COO) <sub>6</sub>	Zr <sub>6</sub> O <sub>20</sub> H <sub>16</sub> (COO) <sub>6</sub>	Zr <sub>6</sub> O <sub>20</sub> H <sub>16</sub> (COO) <sub>6</sub>
CO <sub>2</sub> capacity (mol/kg)	CO <sub>2</sub> /CH <sub>4</sub> separation (10:90, 5bar, 300K)								
Selectivity	4.64	4.33	4.22	3.97	4.80	4.51	4.34	7.21	7.55
CO <sub>2</sub> capacity (mol/kg)	CO <sub>2</sub> /N <sub>2</sub> separation (15:85, 1bar, 313K)								
Selectivity	∞	∞	∞	3058.0	10.8	10.0	11.4	17.5	16.0
LCD (Å)	3.06	2.09	2.47	2.51	1.29	1.26	1.45	2.61	2.80
PLD (Å)	∞	∞	48.5	3157.2	12.3	11.7	16.3	27.1	24.6
AGSA (m <sup>2</sup> /g)	5.79	8.73	9.49	5.56	71.16	70.86	62.61	74.06	81.08
	3.59	3.46	3.71	3.40	58.63	60.83	55.99	68.92	61.29
	1337	1184	1470	429	5261	5423	5076	5025	5233

**B**

	SIFSIX-2-Cu-i	SIFSIX-3-Zn	MgMOF-74	UTSA-16	13X
CO <sub>2</sub> capacity (mol/kg)	CO <sub>2</sub> /CH <sub>4</sub> separation				
Selectivity	(50:50, 1bar, 298K) <sup>9</sup>	(50:50, 1bar, 298K) <sup>9</sup>	(50:50, 5bar, 313K) <sup>74</sup>	(50:50, 2bar, 296K) <sup>5</sup>	(50:50, 5bar, 313K) <sup>74</sup>
	4.16	2.46	8.0	4.25	4.4
	33	231	105.1	29.8	36
CO <sub>2</sub> capacity (mol/kg)	CO <sub>2</sub> /N <sub>2</sub> separation				
Selectivity	(10:90, 1bar, 298K) <sup>9</sup>	(10:90, 1bar, 298K) <sup>9</sup>	(15:75, 0.9bar, 313K) <sup>75</sup>	(15:85, 1bar, 296K) <sup>5</sup>	(16:84, 1.1bar, 288K) <sup>76</sup>
	1.59	2.27	4.43	2.37	3.0
	140	1818	175	314.7	20