

1

2

PMTransformer: Universal Transfer Learning

3

and Cross-material Few-shot Learning

4

in Porous Materials

5

*Hyunsoo Park[⊥], Yeonghun Kang[⊥], and Jihan Kim**

6

Department of Chemical and Biomolecular Engineering, Korea Advanced Institute of Science

7

and Technology (KAIST), 291, Daehak-ro, Yuseong-gu, Daejeon 34141, Republic of Korea

8

⊥ These authors contributed equally to this work

9

10 **ABSTRACT**

11 Porous materials have emerged as a promising solution for a wide range of energy and
12 environmental applications. However, the asymmetric development in the field of MOFs has led
13 to data imbalance when it comes to MOFs versus other porous materials such as COFs, PPNs, and
14 zeolites. To address this issue, we introduce PMTransformer (Porous Material Transformer), a
15 multi-modal pre-trained Transformer model pre-trained on a vast dataset of 1.9 million
16 hypothetical porous materials, including metal-organic frameworks (MOFs), covalent-organic
17 frameworks (COFs), porous polymer networks (PPNs), and zeolites. PMTransformer showcases
18 remarkable transfer learning capabilities, resulting in state-of-the-art performance in predicting
19 various porous material properties. To address the challenge of asymmetric data aggregation, we
20 propose cross-material few-shot learning, which leverages the synergistic effect among different
21 porous material classes to enhance fine-tuning performance with a limited number of examples.
22 As a proof of concept, we demonstrate its effectiveness in predicting bandgap values of COFs
23 using the available MOF data in the training set. Moreover, we established cross-material
24 relationships in porous materials by predicting unseen properties of other classes of porous
25 materials. Our approach presents a new pathway for understanding the underlying relationships
26 between various classes of porous materials, paving the way toward a more comprehensive
27 understanding and design of porous materials.

28

29 **Introduction**

30 Porous materials possess void spaces that can be exploited for many different applications.^{1,2}
31 Depending on the specific nature of the constituent blocks, they can be further categorized into
32 subclasses of materials including metal-organic frameworks (MOFs)³, covalent organic
33 frameworks (COFs)^{4,5}, porous polymer networks (PPNs)⁵, and zeolites⁶. Since these materials are
34 composed of diverse combinations of molecular building blocks, the nearly infinite chemical
35 design space presents an excellent opportunity to design these materials for a wide range of
36 applications, including gas storage and separation⁷, catalysis⁸, and drug delivery⁹. And due to the
37 increasing number of experimental and computational structures, recently there have been several
38 works devoted to using a data-science approach to discover and design new porous materials using
39 various different methods.^{10,11}

40 In recent years, machine learning (ML) models have shown promising results in constructing
41 structure-property relationships for porous materials. For instance, Shi et al.¹² have demonstrated
42 the effectiveness of using two-dimensional (2D) energy histogram features, which include
43 structure-gas interaction energies and energy grid gradients at grid points, as descriptors to
44 accurately predict the gas uptake of MOFs. Also, a 3D convolutional neural network (CNN) with
45 3D voxel, a volume element in 3D space that is analogous to a pixel in 2D space, has been
46 developed as a descriptor for accurate prediction of gas uptake in zeolites.¹³ For predicting
47 electronic properties such as band gap, graph neural networks (GNNs) such as Crystal Graph
48 Convolutional Neural Networks (CGCNN)¹⁴ and MatErials Graph Network (MEGNET)¹⁵ have
49 shown high performance. Also, various descriptors have been developed including geometric,
50 chemical, topological features, revised autocorrelations (RAC)¹⁶ and smooth overlap of atomic
51 positions (SOAP)¹⁷. Recently, MOFTransformer¹⁸, a multi-modal pre-training Transformer, has

52 been introduced to achieve universal transfer learning in MOFs, showcasing its exceptional ability
53 to transfer learning across various MOF properties.

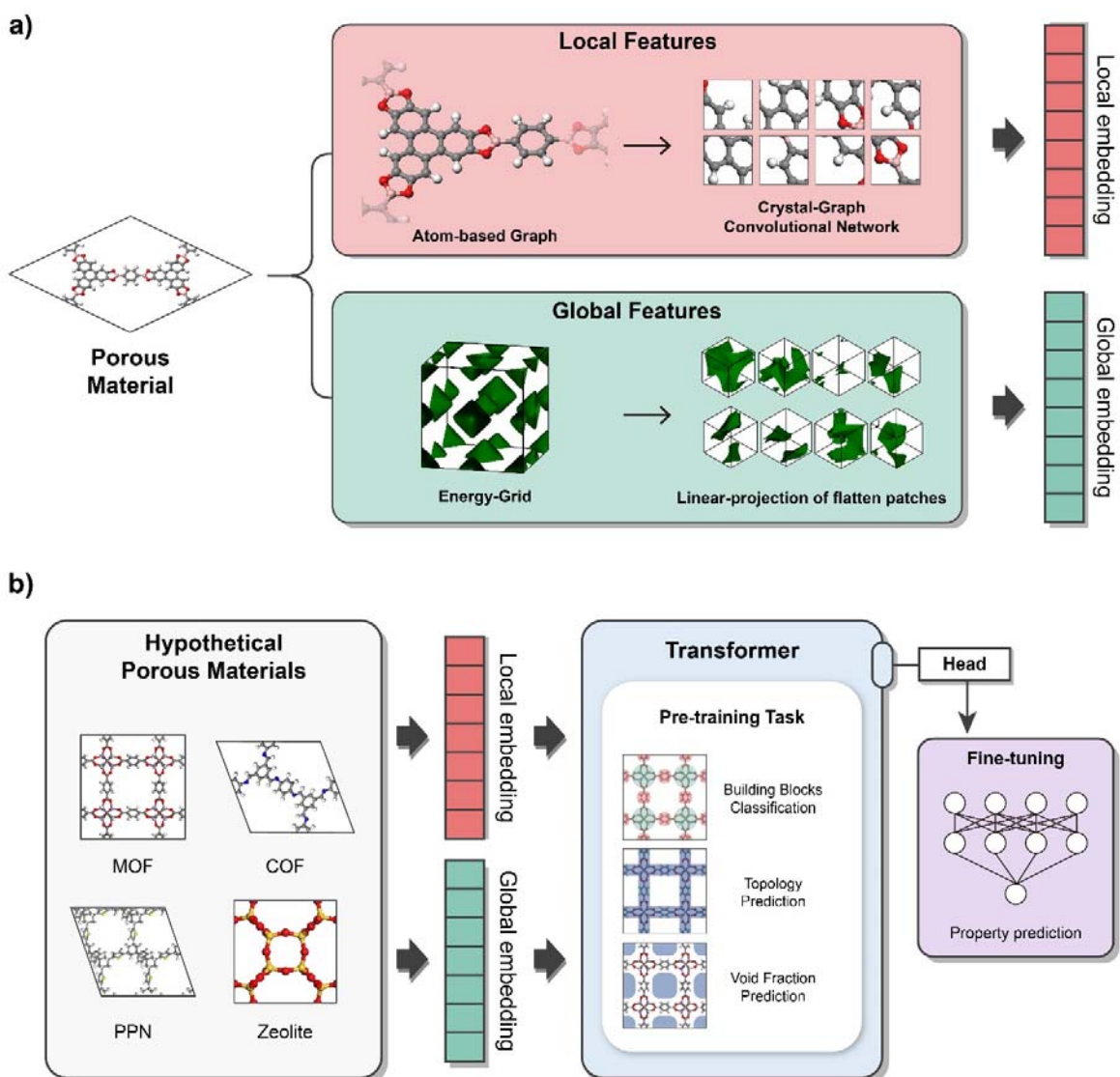
54 Despite the potential of machine learning models for predicting material properties in porous
55 materials, their usefulness remains limited by the availability of data. And while MOFs have been
56 extensively explored due to the large number of experimentally reported structures (over
57 100,000)¹⁹, other porous materials have much smaller number of experimentally reported data.
58 The CoRE COF²⁰ and Curated COF²¹ databases include around 600 experimentally reported COFs,
59 and fewer than 100 PPNs have been synthesized.²² COFs and PPNs are formed by covalent bonds
60 and strong C-C bonds, respectively, which make them harder to synthesize into crystalline
61 materials due to the lack of reversible reactions.²³ Additionally, zeolites, composed of Si and O
62 atoms have only a bit over two hundred known topologies.²⁴ This lack of available data for other
63 porous materials poses a significant challenge for developing accurate machine-learning models
64 across all porous materials and perhaps is one of the reasons on why the machine learning works
65 on porous materials thus far has been skewed towards MOFs.

66 To overcome the challenge of asymmetry data aggregation, it is our opinion that leveraging data
67 from other porous materials represents a promising solution when a specific material class lacks
68 sufficient data (both in terms of number of materials and properties) for model training. For
69 instance, the restricted data availability of only hundreds of COF structures may pose substantial
70 obstacles when it comes to developing machine learning models to predict the properties of COFs.
71 By incorporating data from abundant source materials such as MOFs, the accuracy of the model
72 predictions can be improved through exploiting the potential synergistic effect between the two
73 material classes. To the best of our knowledge, this type of cross-material transfer learning has yet
74 to be explored in any other materials. Indeed, one can envision that such an approach could

75 enhance the accuracy of machine learning model predictions in overcoming data scarcity
76 challenges through the potential synergistic effect between materials from distinct classes.

77 In this work, we introduce the Porous Material Transformer (PMTransformer), which is a multi-
78 modal Transformer architecture based on the MOFTransformer and is pre-trained with 1.9 million
79 hypothetical porous materials, including MOFs, COFs, PPNs, and zeolites. The model showcases
80 excellent transfer learning capability across various properties of porous materials, thereby
81 achieving state-of-the-art performance in predicting multiple different properties. To address the
82 challenge of asymmetry data aggregation in porous materials, we propose cross-material few-shot
83 learning to improve predictions of materials lacking available data for their properties by
84 exploiting the uniform characteristics in porous materials. Moreover, we obtain cross-relationships
85 in porous materials by predicting unseen properties of other classes of porous materials. Our
86 approach provides a novel perspective for understanding the underlying and uniform relationships
87 between various classes of porous materials, allowing for the prediction of previously unexplored
88 properties across these material classes.

89



90

91 **Figure 1.** (a) Data representations for porous materials incorporating both local features and global
 92 features used with atom-based graphs and energy grids, respectively. (b) Overall schematics of
 93 PMTransformer. The model was pre-trained with 1.9 million hypothetical porous materials with
 94 three pre-training tasks to capture local and global features in a pre-training stage. In a fine-tuning
 95 stage, the PMTransformer is fine-tuned to predict properties of porous materials where its initial
 96 weights are initialized with the pre-training weights.

97

98 **Results**

99 **Data representations of MOFs for PMTransformer**

100 Figure 1(a) shows a representative porous materials input data representations for two disparate
101 features (i.e., local features and global features), which serve as inputs of PMTransformer. The
102 local features involve atomistic information related to chemistry of building blocks and specific
103 bonds. The output features of crystal graph convolutional neural networks (CGCNN) were adopted
104 to describe the local features given that they enable capturing atoms' neighbor information such
105 as atom types, distances between neighbor atoms. On the other hand, the global features represent
106 crystalline features including topological and geometric descriptors such as pore volume, surface
107 area, which are captured by the 3D energy grids. The grids are created by calculating interaction
108 energy between a structure and a gas molecules (or gas probe) at each grid point, and can be treated
109 as 3D images, thereby leading to understand the global features. Similar to the Vision Transformer,
110 energy grids are divided by 6 x 6 X 6 patches and flattened by a linear projection. Finally, the local
111 and global embedding are fed into the Transformer encoder of PMTransformer.

112 **Pre-training of PMTransformer**

113 Figure 1(b) illustrates the overall schematic of PMTransformer indicating pre-training and fine-
114 tuning approach to achieve universal transfer learning in porous materials. The pre-training enables
115 our model to learn how to represent the input data in a way that captures its essential features,
116 which can then be used to improve the performance of the model on fine-tuning tasks. The pre-
117 training tasks are designed to enable the model to understand the essential features of porous
118 materials, resulting in superior performance in transfer learning. Previous studies have
119 demonstrated the effectiveness of pre-training tasks designed for MOFs in the MOFTransformer
120 model.¹⁸ The pre-training with topology prediction, void fraction prediction, and metal cluster &

121 organic linker classification significantly improve transfer learning in MOFs as these tasks
122 facilitate capturing both local and global features of MOFs, which is critical for accurate property
123 prediction.

124 Building on the pre-training tasks of MOFTransformer, we extended the pre-training tasks to
125 include COFs, PPNs, and zeolites. The pre-training tasks include topology prediction and void
126 fraction prediction for capturing global features of porous materials, and building block
127 classification for capturing local features. Building block classification involves classifying the (1)
128 metal cluster and organic linkers for MOFs, (2) center and linker for COFs and PPNs, and (3) Si
129 and O atoms for zeolites. The accuracies of the pre-training tasks in PMTransformer are
130 comparable to those of MOFTransformer, with topology prediction and building block
131 classification achieving accuracies of 0.98 and 0.99, respectively, and void fraction prediction
132 having a mean absolute error of 0.01.

133 **Construction of Porous Material Database**

134 Large and diverse pre-training datasets help the Transformer model learn and comprehend the
135 underlying relationships in pre-training datasets, resulting in improving transfer learning capability
136 in fine-tuning stages. When pre-training the MOFTransformer, one million hypothetical MOFs
137 (hMOFs) were created using the PORMAKE python library¹⁰, with the molecular building blocks
138 and topologies derived from the CoRE MOF,²⁵ ToBaCCo²⁶, and RSCR²⁷ database. In this work,
139 we expanded the pre-training dataset for porous materials to include COFs, PPNs, and zeolites,
140 thereby making it larger and more diverse, as illustrated in Figure 2. Notably, creating pre-training
141 datasets from scratch also facilitates the annotation of topology and building block information for
142 pre-training tasks.

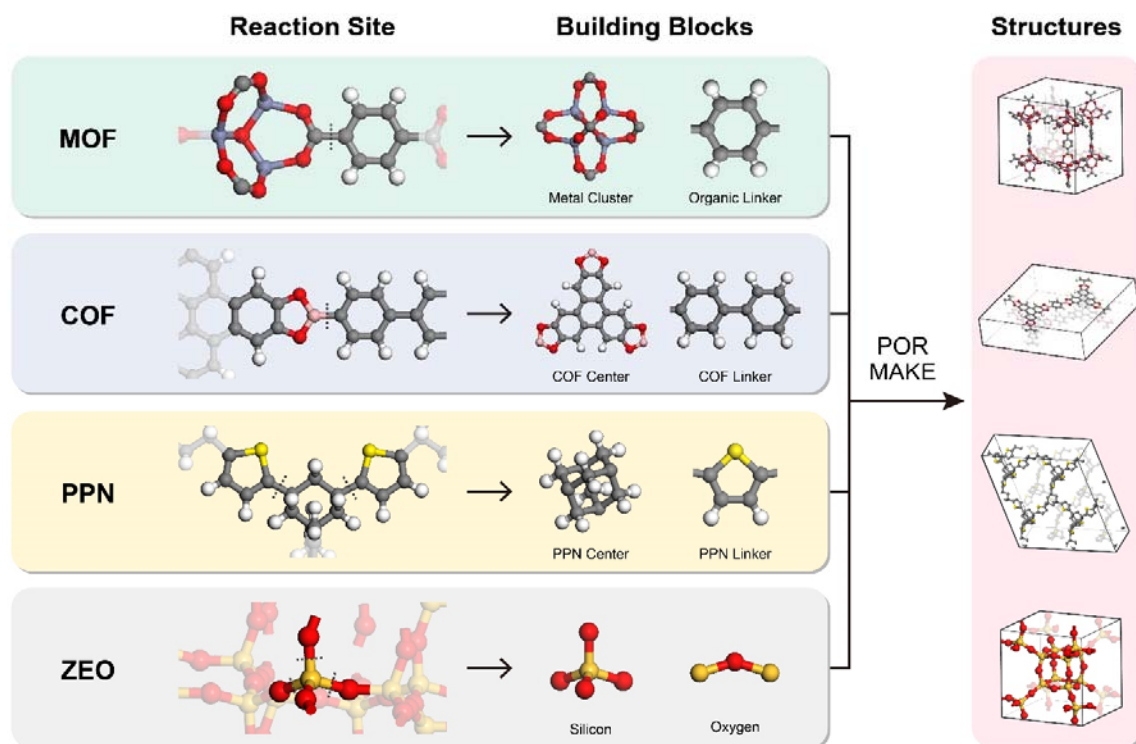
143 COFs are constructed from organic building blocks with different topologies, linked with
144 covalent bonds. The organic building blocks are relatively rigid backbones that endow the COFs
145 with crystallinity, making them distinct from organic polymers with low crystallinity. The COFs
146 can be synthesized using reactions of boron, triazine, and imine condensation.² Various databases
147 of synthesized COFs, including the CURATED COF and CoRE COF databases, as well as
148 hypothetical databases, have been established. For example, Lan et al.²⁸ developed the Genomic
149 COF database, which contains 471,990 COFs constructed from 130 genetic structure units (GSUs)
150 consisting of 58 centers, 64 linkers, and 8 functional groups with 24 topologies, using reactive
151 sites and quasi-reactive assembly algorithms (QReaxAA). For the pre-training dataset, we
152 constructed a hypothetical COF (hCOF) database using the 130 GSUs and topologies registered in
153 the RCSR database. As shown in Figure 2, we generated 519,606 COFs by PORMAKE, of which
154 only 747 topologies met the constraints with a root mean squared deviation (RMSD) of atomic
155 positions between the building blocks and target node position to measure the strain energy less
156 than 0.3. Notably, the large hCOF dataset, containing numerous COF structures with diverse
157 topologies, enables the PMTransformer to achieve a superior understanding of COFs during pre-
158 training stages.

159 PPNs constitute a class of porous polymers assembled from tunable building blocks through
160 polymerization reactions, such as homocoupling of tetrahedral monomers. These reactions are
161 typically irreversible, leading to PPNs with exceptional thermal and chemical stability, but
162 amorphous materials. The amorphous nature presents a significant challenge for computational
163 modeling. To address this issue, Martin et al.²⁹ developed the in-silico PPN database, which
164 utilized a crystalline modeling approach that successfully reproduced the gas adsorption behavior
165 of PPNs. Building on this work, we constructed a diverse hypothetical PPN (hPPN) database for

166 pre-training by PORMAKE, utilizing the same building blocks from the in-silico PPN database,
167 but with more diverse topologies requiring nodes with four connections for tetrahedral monomers.
168 The building blocks consist of Si, Ge, C, adamantane as centers and 4952 linkers. They result in
169 277,250 hPPNs including the interpenetrated structures.

170 Zeolites are a type of crystalline aluminosilicate material composed of silicon, aluminum, and
171 oxygen atoms arranged in tetrahedral structures. Compared to other porous materials like MOFs,
172 COFs, and PPNs, zeolites have a smaller chemical space due to their immutable building blocks.
173 The IZA database³⁰ currently lists around 250 known zeolite topologies, while the PCOD
174 database³¹ was developed using Monte Carlo algorithms and contains many predicted zeolite
175 structures. To prepare for pre-training, we constructed 278 zeolite structures with topologies
176 featuring four connection points using a top-down approach with the RCSR database by
177 PORMAKE. We generated 34,750 zeolites by augmenting these structures through a translational
178 motion in five parts for each cell direction. To supplement the dataset, we randomly selected
179 65,250 zeolites from the PCOD database, resulting in 100,000 zeolite structures in the pre-training
180 data. The ToposPro software³² was used to obtain topology information for the structures, but there
181 were still unknown topologies. As such, we labeled these unknown topologies as “unknown
182 topology” during the pre-training stage. All of the generated structures were geometrically optimized
183 using the LAMMPS package³³ with the UFF force field³⁴.

184



185

186 **Figure 2.** Construction of diverse and large pre-training dataset for porous materials, including
 187 COFs, PPNs, and zeolites, utilized in pre-training the PMTransformer. Hypothetical structures
 188 were generated using the PORMAKE Python library, resulting in 1 million hMOFs, 519,606
 189 hCOFs, 277,250 hPPNs, and 100,000 zeolites.

190

191 **Fine-tuning results**

192 To evaluate the performance of PMTrasformer, we compared it with the scratch model (i.e., the
193 default PMTransformer model without any pre-training), the MOFTransformer, which was pre-
194 trained with only MOFs, and several other baseline models, including energy histogram³⁵,
195 descriptor-based machine learning (ML) model³⁶, and crystal graph convolutional neural network
196 (CGCNN), using mean absolute errors (MAEs) on different properties of MOFs, COFs, PPNs, and
197 zeolites. The evaluated properties included gas uptake, diffusivity, Henry coefficient, heat of
198 adsorption, stability, and bandgap, as summarized in Table 1, 2.

199 With regards to the baseline models, the energy histogram model employed the Least Absolute
200 Shrinkage and Selection Operator (LASSO) regression³⁷, which involved taking an energy
201 histogram that had been converted from energy grids by energy bins. The descriptor-based model
202 utilized 5 geometrical properties (i.e. largest cavity diameter, pore-limiting diameter, gravimetric
203 accessible surface area, volumetric accessible surface area, and volume fraction) as well as 12
204 chemical properties (i.e. metal type present, number of specified element atoms in unit cell), and
205 6 additional chemical properties (i.e. total degree of unsaturation, metallic percentage, oxygen to
206 metal ratio, electronegative to total ratio, weighted electronegativity per atom, and nitrogen to
207 oxygen ratio) as inputs. All of these descriptors were used as input to a random forest model. On
208 the other hand, the CGCNN uses atom-based graph representation as inputs, and consists of five
209 convolution layers, one hidden layer after pooling, 64 hidden atom features in convolution layers,
210 and 128 hidden 7 features after pooling.

211 For the prediction of the MOF properties, the scratch model demonstrated superior performance
212 compared to other baseline models (i.e. energy histogram, descriptor-based ML model, CGCNN)
213 across all properties, as shown in Table 1. It indicates that the data representation of our model

214 facilitates capturing the underlying feature of MOFs, leading to high performance in predicting
215 various MOF properties. Also, the fine-tuned PMTransformer achieved lower MAE values in all
216 of the MOF properties except for O₂ uptake and N₂ diffusivity compared to the MOFTransformer.
217 This observation indicates that including other porous materials, such as COFs, PPNs, and zeolites,
218 in the pre-training dataset of MOFTransformer leads to higher performance in predicting MOF
219 properties, indicating synergetic effect due to similarity across all porous materials.
220 For properties of COFs, PPNs, and zeolites, PMTransformer exhibited the lowest MAEs across all
221 properties except for CH₄ uptake at 65 bar in COFs, in which the MOFTransformer had the lowest
222 MAE, as shown in Table 2. Our findings suggest that pre-training with a large set of diverse porous
223 materials, as opposed to pre-training with MOFs alone, plays an important role in improving
224 performance in predicting various properties of porous materials.
225

Material	Property	Number of Dataset	Energy histogram	Descriptor-based ML	CGCNN	Scratch	MOF Transformer	PM Transformer	Reference
MOF	H ₂ Uptake (100 bar)	20,000	9.183	9.456	32.864	7.018	6.377	5.963	18
MOF	H ₂ diffusivity (dilute)	20,000	0.644	0.398	0.6600	0.391	0.367	0.366	18
MOF	Band-gap	20.373	0.913	0.590	0.290	0.271	0.224	0.216	38
MOF	N ₂ uptake (1 bar)	5,286	0.178	0.115	0.108	0.102	0.071	0.069	36
MOF	O ₂ uptake (1 bar)	5,286	0.162	0.076	0.083	0.071	0.051	0.053	36
MOF	N ₂ diffusivity (1 bar)	5,286	7.82e-5	5.22e-5	7.19e-5	5.82e-05	4.52e-05	4.53e-05	36
MOF	O ₂ diffusivity (1 bar)	5,286	7.14e-5	4.59e-5	6.56e-5	5.00e-05	4.04e-05	3.99e-05	36
MOF	CO ₂ Henry coefficient	8,183	0.737	0.468	0.426	0.362	0.295	0.288	39
MOF	Thermal stability	3,098	68.74	49.27	52.38	52.557	45.875	45.766	40

226 **Table 1.** Comparison of mean absolute error (MAE) values for various baseline models, scratch,
227 MOFTransformer, and PMTransformer on different properties of MOFs. The bold values indicate
228 the lowest MAE value for each property.

229

Material	Property	Number of Dataset	Energy histogram	Descriptor-based ML	CGCNN	Scratch	MOF Transformer	PM Transformer	Reference
COF	CH ₄ uptake (65bar)	39,304	5.588	4.630	15.31	2.883	2.268	2.126	41
COF	CH ₄ uptake (5.8bar)	39,304	3.444	1.853	5.620	1.255	0.999	1.009	41
COF	CO ₂ heat of adsorption	39,304	2.101	1.341	1.846	1.058	0.874	0.842	42
COF	CO ₂ log KH	39,304	0.242	0.169	0.238	0.134	0.108	0.103	42
PPN	CH ₄ uptake (65bar)	17, 870	6.260	4.233	9.731	3.748	3.187	2.995	29
PPN	CH ₄ uptake (1bar)	17, 870	1.356	0.563	1.525	0.602	0.493	0.461	29
Zeolite	CH ₄ KH (unitless)	99,204	8.032	6.268	6.334	4.286	4.103	3.998	43
Zeolite	CH ₄ Heat of adsorption	99,204	1.612	1.033	1.603	0.670	0.647	0.639	43

230 **Table 2.** Comparison of mean absolute error (MAE) values for various baseline models, scratch,
231 MOFTransformer, and PMTransformer on different properties of COFs, PPNs, and zeolites. The
232 bold values indicate the lowest MAE value for each property.

233

234 **Discussion**

235 **Cross-material few-shot learning: Prediction of COF Bandgap**

236 Few-shot learning is a promising approach for addressing the challenges posed by limited data
237 availability (typically less than 500) in ML models.⁴⁴ In particular, fine-tuning the pre-trained
238 models in vision or language model with only few examples can lead to high performance on
239 unseen tasks. In this work, we applied few-shot learning to the PMTransformer. To address the
240 issue of asymmetry data aggregation in porous materials, we propose a cross-material few-shot
241 learning approach. This approach exploits the synergistic effects from high similarity between
242 different classes of porous materials to improve performance. Specifically, we utilize the relatively
243 abundant number of data for the metal-organic frameworks (MOFs) to train the PMTransformer
244 to predict the properties of other types of porous materials.

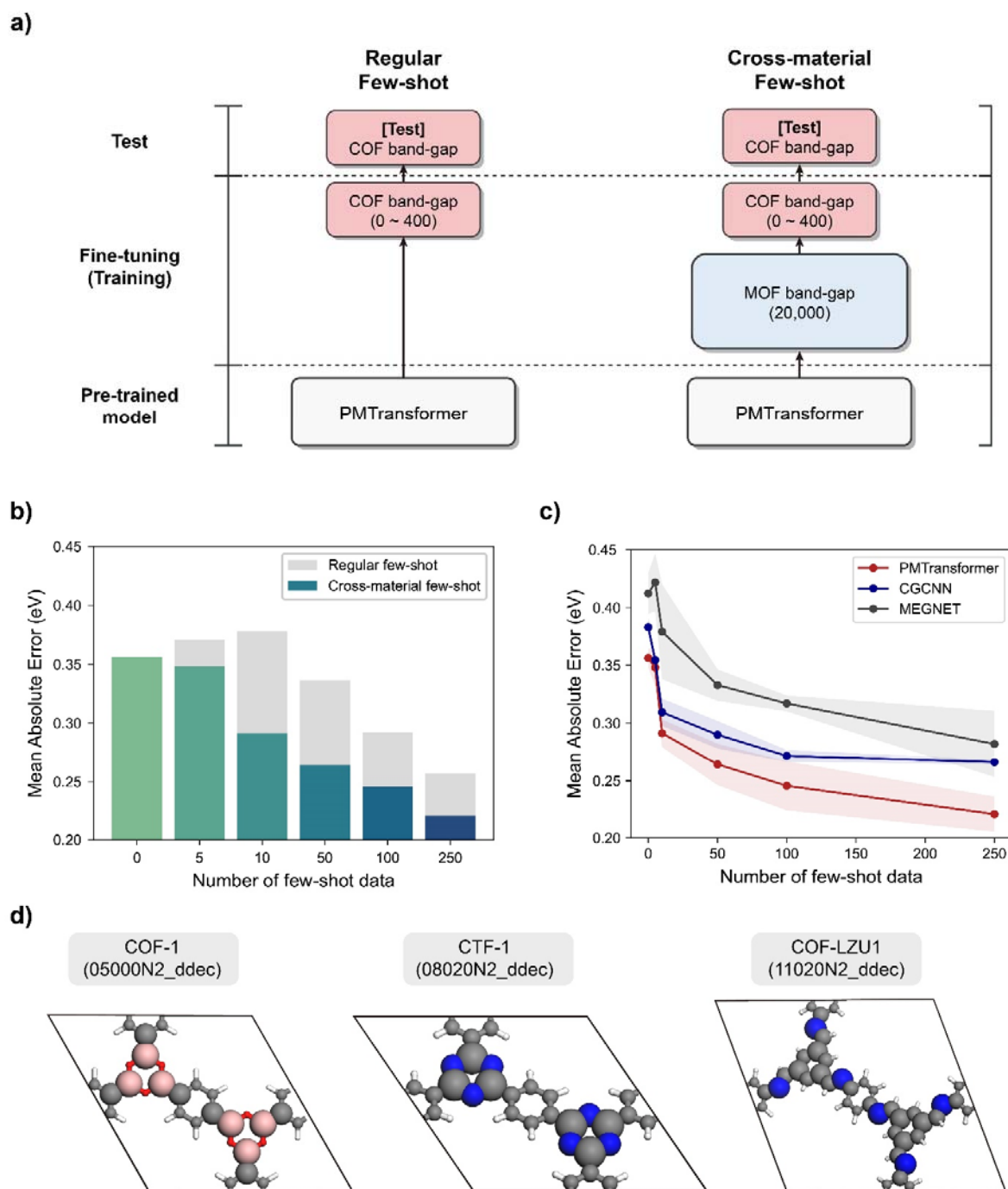
245 Figure 3(a) illustrates the case study application of cross-material few-shot learning to predict the
246 band gap values of the COFs calculated by DFT, where only 400 COF band gap data⁴⁵ in the
247 Curated COF database are available. The PMTransformer was fine-tuned to predict the COF band
248 gap values by initializing the weights of the model with the weights obtained from the fine-tuned
249 PMTransformer trained on 20,000 MOF bandgaps from the QMOF database. This approach differs
250 from the regular few-shot learning, which involves fine-tuning the PM Transformer with only 400
251 COF bandgaps. The COF bandgap data was split into 250, 50, and 100 for training, validation, and
252 test. The performance of the few-shot learning and the proposed cross-material few-shot learning
253 methods was compared in terms of mean absolute error (MAE) as the number of training examples
254 ranged from 0 to 250, as shown in Figure 3(b). The results were averaged over five trials on the
255 test set. Notably, the cross-material few-shot learning outperformed the regular few-shot learning
256 method. For instance, when the number of training examples was 250, the cross-material few-shot

257 learning achieved an r^2 score of 0.48, whereas the few-shot learning method achieved an r^2 score
258 of only 0.30. These results demonstrate the effectiveness of the proposed cross-material few-shot
259 learning method, which exploits the high similarity among porous materials to achieve a synergetic
260 effect, particularly in cases with limited available data. To further investigate the effect of the
261 number of source materials (MOFs), an ablation study was conducted by varying the number of
262 MOFs used for training from 0 to 20,000 when the number of COF training data was fixed at 250,
263 as shown in Supplementary Figure S6. The results indicate that the performance of the cross-
264 material few-shot learning converged when the number of source material for MOF was at 10,000.
265 Furthermore, we evaluated the cross-material few-shot learning performance of PMTransformer
266 when compared to other ML baseline models such as CGCNN and MEGNET which exhibited
267 high performance in predicting the band gaps in MOFs, as shown in Figure 3(c). The
268 PMTransformer exhibits superior performance compared to other baseline models. This can be
269 attributed to its pre-training, which enabled the PMTransformer to capture general patterns and
270 relationships in porous materials and adapt to new tasks with limited examples. Moreover, it can
271 be observed that the regular few-shot and cross-material few-shot learning in CGCNN and
272 MEGNET do not exhibit a significant improvement in performance compared to the
273 PMTransformer, as demonstrated in Supplementary Figures S7 and S8.

274 In general, the Transformer architectures⁴⁶ are capable of generating attention scores through their
275 attention layers, which reflect the degree of attention the model pays to input features for a given
276 task. These attention scores can be utilized as a tool for feature importance analysis. Figure 3(d)
277 presents the attention scores of representative COF structures, such as COF-1⁴⁷, CTF-1⁴⁸, and
278 COF-LZU1⁴⁹. In these scores, the larger size of atoms represents higher attention scores, which in
279 turn can be considered as more influential factors in determining band gap. It is important to note

280 that these structures are composed of benzene rings as linkers and are distinct from their
281 corresponding centers. The analysis of attention scores reveals that the centers have higher
282 attention scores than the linkers, indicating that they play a more prominent role in determining
283 the band gaps. The 2D COFs are known for their ability to extend the π -conjugation system, which
284 leads to greater emphasis being placed on centers that have more than two connection points, as
285 compared to linkers that only have two connection points. Moreover, it is noteworthy that the π -
286 conjugation ability of C-N bonds is a significant aspect to consider. The analysis of attention scores
287 for CTF-1 and COF-LZU-1 indicates that nitrogen atoms within the structures' centers exhibit
288 higher attention scores compared to other atoms. In contrast, the oxygen atoms in the B₃O₃ rings
289 of COF-1 have relatively lower attention scores among their centers, primarily due to the absence
290 of π -conjugation. This analysis demonstrates the utility of attention scores in providing insights
291 into the underlying factors that determine the band gap of COF structures, thereby facilitating the
292 development of more efficient and accurate models for porous materials.

293



294
 295 **Figure 3.** (a) Application of cross-material few-shot learning to predict COF band gaps with
 296 limited data. The PMTransformer is fine-tuned using weights from the fine-tuned model on
 297 20,000 MOF band gaps to predict COF band gaps with only 400 examples available. This
 298 approach differs from the regular few-shot learning method involving the fine-tuning with only

299 400 COFs (b) Comparison of MAEs between the regular few-shot and cross-material few-shot
300 results for prediction of band gap of COF as the number of training data (few-shot data)
301 increases from 0 to 250 for PMTransformer. (c) Comparison of MAEs for the cross-material
302 few-shot learning using PMTransformer, CGCNN, and MEGNET as the number of training data
303 (few-shot data) increases from 0 to 250. (d) The schematics for attention scores obtained from
304 the fine-tuned PMTransformer to predict COF band gaps for COF-1, CTF-1, COF-LZU1. The
305 larger atom size represents higher attention scores.
306

307 **Cross-material relationship in porous materials: H₂ Uptake**

308 Our next case study investigated the cross-material relationship in porous materials and evaluates
309 the ability of PMTransformer to predict unseen properties of other classes of porous materials. The
310 H₂ uptake at 77K and 100 bar was calculated for 5,000 MOFs, COFs, PPNs, and zeolites and
311 randomly split into 4,000 training, 500 validation, and 500 test sets.

312 In Figure 4(a), a heatmap shows the r² scores obtained from the PMTransformer fine-tuned with
313 training (or source) materials to predict H₂ uptake and tested on test (or target) materials without
314 further fine-tuning. The diagonal of the heatmap represents the r² scores when training and test
315 materials are identical. Remarkably, the PMTransformer fine-tuned with MOFs as the training
316 material achieved r² scores higher than 0.9 when predicting the H₂ uptake of COFs and PPNs in
317 the test set, which is comparable to the r² scores obtained when training and test materials are the
318 same. These results demonstrate the ability of the PMTransformer model to accurately predict the
319 H₂ uptake of COFs and PPNs when fine-tuned with MOFs as the training material. It is noteworthy
320 that the r² scores between MOFs, COFs, and PPNs exceed 0.85, indicating their synergetic effect
321 in the cross-material relationship due to their high level of similarity, except when COFs and MOFs
322 are respectively the training and test materials. Conversely, zeolites exhibit low r² scores,
323 regardless of the source materials, suggesting that zeolites have a lack of synergy with other classes
324 of porous materials.

325 This observation is supported by the t-SNE plot created by the class tokens from the fine-tuned
326 PMTransformer with MOFs, COFs, PPN, and zeolites in test set, respectively, as illustrated in
327 Figure 4(b). The plot reveals a unique clustering of zeolites, which are positioned solely within the
328 lower H₂ uptake region. It is attributed to the composition of zeolites, which consist primarily of
329 Si and O atoms, resulting in smaller pore sizes and consequently, lower H₂ uptake compared to

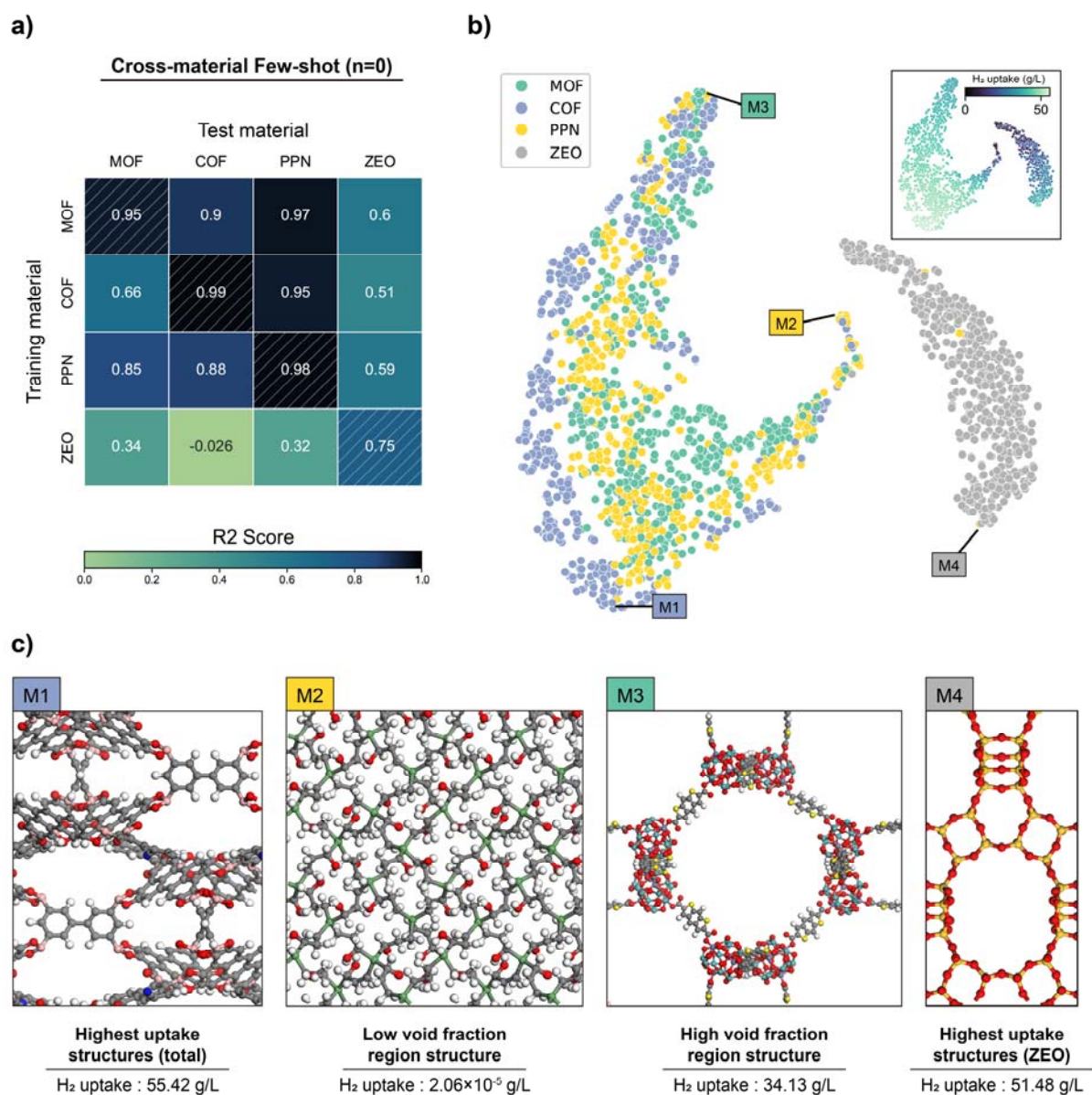
330 other porous materials that are typically composed of molecular building blocks such as metal
331 nodes, organic linkers, and polymer monomers.

332 Figure 4(c) shows four highlighted structures in the t-SNE plot, where their building blocks and
333 naming are shown in Supplementary Figure S9. **M1** exhibits the highest H₂ uptake value of 55.42
334 g/L, which is 2D COF with the *hyw* topology, composed of 3,4,9,10-Perylenetetra-carboxylic acid,
335 biphenyl, and 1-cyanopyrene. Interestingly, the high H₂ uptake region in the vicinity of **M1** is
336 mostly populated by COFs, which can be attributed to their void fraction. To investigate this
337 further, the t-SNE plot in Supplementary Figure S10 is colored according to their void fraction
338 values calculated using ZEO++. The COF structures located within the high H₂ uptake region have
339 void fractions ranging between 0.45 and 0.55, which seems to be the optimal range for high H₂
340 uptake performance, as shown in Supplementary Figure S11. In contrast, **M2** (PPN) and **M3** (MOF)
341 exhibit very low and very high void fraction values, respectively, as depicted in Figure SX, due to
342 their building block. **M2** consists of Ge atoms as centers and short linkers, specifically 1,3-
343 dibromo-1-propanol, while **M3** has a long organic linker, dithieno[3,2-b:2',3'-e]benzene-2,6-
344 dicarboxylic acid⁵⁰. Additionally, among the zeolites, **M4**, which is sourced from the PCOD
345 database, exhibits the highest H₂ uptake value of 51.48 g/L.

346 It should be noted that MOFs and PPNs cluster closely together, while COFs are more dispersed,
347 with most located in the highest H₂ uptake region. This behavior can be ascribed to the fact that
348 MOFs and PPNs have common topologies when constructed by PORMAKE, while COFs have a
349 greater diversity of 2D topologies. Indeed, most of entries in the CoRE COF and the CURATED
350 COF database are 2D COF, rather than 3D. Among GSUs from the genomic COF database, the
351 centers are mostly composed of building blocks from 2D COFs than 3D COFs, while building
352 blocks of MOFs and PPNs were derived from 3D structures. This is because the building blocks

353 of MOFs were obtained from the CoRE COF database, which contains only 3D MOFs in the CSD
354 database. This limitation suggests a need for a more large and diverse pre-training dataset,
355 including 2D MOFs, which would lead to superior transfer learning capability in the fine-tuning
356 stage. Other limitation is the lower accuracy of PMTransformer in predicting zeolite properties. It
357 can be attributed to asymmetry data of zeolites when compared to other porous materials as well
358 as lower diversity of zeolites in porous materials. The pre-training dataset contains only 100,000
359 zeolites, because their the building blocks (i.e., Si and O) are not tunable, resulting in small
360 chemical space. These limitations must be taken into consideration in future studies.

361



362

363 **Figure 4.** (a) Heatmap of r2 scores obtained from the PMTransformer fine-tuned with training
 364 materials to predict H₂ uptake and tested on test materials without further fine-tuning between
 365 MOFs, COFs, PPNs and zeolites. (b) A t-SNE plot of class tokens obtained from the fine-tuned
 366 PMTransformer with MOFs, COFs, PPN, and zeolites in the test set, with the additional small
 367 figure colored by H₂ uptake (c) The t-SNE plot highlights several structures based on their H₂
 368 uptake and void fraction characteristics, including **M1** with the highest H₂ uptake, **M2** with low

369 void fraction, **M3** with high void fraction, and the zeolite structure **M4** with the highest H₂
370 uptake.

371 **Conclusions**

372 In this work, we present the Porous Material Transformer (PMTransformer) model that combines
373 multi-modal features from MOFs, COFs, PPNs, and zeolites. By pre-training on 1.9 million
374 hypothetical porous materials, our model achieved state-of-the-art performance in predicting
375 various properties of porous materials via fine-tuning. Furthermore, we introduced cross-material
376 few-shot learning to address the challenge of asymmetry data aggregation in porous materials and
377 proposed a method for predicting previously unexplored properties across different material
378 classes (e.g. using MOF data to predict COF properties). Our approach provides an opportunity
379 for understanding the underlying relationships between various classes of porous materials,
380 allowing for the prediction of previously unexplored properties across these material classes, and
381 thus facilitating a more comprehensive understanding and design of porous materials.

382

383 **Methods**

384 **Training details**

385 We adopted a pre-training and fine-tuning approach similar to that used in previous work,
386 MOFTransformer. We note that in few-shot learning, the model is fine-tuned with only a few
387 samples, leveraging the pre-training weights as an initialization. The optimization process in all
388 stages, including pre-training, fine-tuning, and few-shot learning, employed the AdamW⁵¹
389 optimizer with a learning rate of 10^{-4} and weight decay of 10^{-2} . During the initial phase of the
390 optimization process, the learning rate was gradually increased for the first 5 % of the total epoch
391 and then linearly decayed to zero for the remaining epochs.

392 During the pre-training stage, the model was trained using a batch size of 1024 for a total of 100
393 epochs. For fine-tuning and few-shot learning, the model was trained using a smaller batch size of
394 32 for a total of 20 epochs. The dataset is split randomly into train, validation, and test, with a ratio
395 of 8 : 1 : 1. We adopted the standardization method for scaling the target properties

396 **Computational details for molecular simulation**

397 The H₂ uptake of 5000 MOFs, COFs, PPNs, and zeolites was calculated for cross-material
398 relationships using the RASPA package⁵². The property was used due to its relatively facile
399 calculation with a united atom model. The pseudo-Feynman-Hibbs model was used to describe the
400 H₂ behavior at low temperatures, and Lenard-Jones potentials were fitted to the Feynman-Hibbs
401 potential⁵³ at T = 77 K. The UFF force field was used for all molecules except H₂, with the Lorentz-
402 Berthelot mixing rule and a cutoff distance of 12.8 Å. To calculate H₂ uptake, GCMC simulations
403 were performed for 10k production cycles at 100 bar and 77 K, with 5k cycles used for
404 initialization.

405

406 **Conflicts of interest**

407 There are no conflicts to declare.

408 **Author Contributions**

409 H.P and Y.K contributed equally to this work. H.P and Y.K developed PMTransformer and
410 wrote the manuscript with J.K. The manuscript was written through the contributions of all authors.
411 All authors have given approval for the final version of the manuscript.

412 **Data availability**

413 The UFF-optimized CIF files of hypothetical porous materials database used as pre-training
414 dataset are available at <https://doi.org/10.6084/m9.figshare.21810147> for MOFs,
415 <https://doi.org/10.6084/m9.figshare.22699303> for COFs, PPNs and zeolites. The pre-training
416 PMTransformer model is available at <https://doi.org/10.6084/m9.figshare.22698655.v2>.

417 **Code availability**

418 The PMTransformer library is based on the MOFTransformer, which is available at at
419 <https://github.com/hspark1212/MOFTransformer>. From version 2.0.0, the default pre-training
420 model has been changed from MOFTransformer to PMTransformer. For the sake of
421 reproducibility, all results in this manuscript are obtained from a 2.0.0 version of the library, which
422 is available at <https://pypi.org/project/moftransformer/2.0.0>.

423 **Acknowledgements**

424 H. P., Y. K., and J. K. acknowledge funding from the National Research Foundation of Korea
425 (NRF) under Project Number 2021M3A7C208974513 and 2021R1A2C2003583. This work was
426 supported by the National Supercomputing Center with supercomputing resources including

427 technical support (KSC-2022-CRE-0515). The authors would like to thank Berend Smit and
428 Beatriz Mouriño who provided us with the band gap values calculated by DFT for the CURATED
429 COFs database.

430

431 **References**

- 432 1 Bennett, T. D., Coudert, F.-X., James, S. L. & Cooper, A. I. The changing state of porous
 433 materials. *Nature Materials* **20**, 1179-1187 (2021).
- 434 2 Geng, K. *et al.* Covalent organic frameworks: design, synthesis, and functions. *Chemical*
 435 *reviews* **120**, 8814-8933 (2020).
- 436 3 Zhou, H.-C., Long, J. R. & Yaghi, O. M. Vol. 112 673-674 (ACS Publications,
 437 2012).
- 438 4 Feng, X., Ding, X. & Jiang, D. Covalent organic frameworks. *Chemical Society Reviews*
 439 **41**, 6010-6022 (2012).
- 440 5 Makal, T. A., Li, J.-R., Lu, W. & Zhou, H.-C. Methane storage in advanced porous
 441 materials. *Chemical Society Reviews* **41**, 7761-7779 (2012).
- 442 6 Ozin, G. A., Kuperman, A. & Stein, A. Advanced zeolite, materials science. *Angewandte*
 443 *Chemie International Edition in English* **28**, 359-376 (1989).
- 444 7 Li, H. *et al.* Recent advances in gas storage and separation using metal–organic
 445 frameworks. *Materials Today* **21**, 108-121 (2018).
- 446 8 Lee, J. *et al.* Metal–organic framework materials as catalysts. *Chemical Society Reviews*
 447 **38**, 1450-1459 (2009).
- 448 9 Della Rocca, J., Liu, D. & Lin, W. Nanoscale metal–organic frameworks for biomedical
 449 imaging and drug delivery. *Accounts of chemical research* **44**, 957-968 (2011).
- 450 10 Lee, S. *et al.* Computational screening of trillions of metal–organic frameworks for high-
 451 performance methane storage. *ACS Applied Materials & Interfaces* **13**, 23647-23654
 452 (2021).
- 453 11 Yao, Z. *et al.* Inverse design of nanoporous crystalline reticular materials with deep
 454 generative models. *Nature Machine Intelligence* **3**, 76-86 (2021).
- 455 12 Shi, K. *et al.* Two-Dimensional Energy Histograms as Features for Machine Learning to
 456 Predict Adsorption in Diverse Nanoporous Materials. *Journal of Chemical Theory and*
 457 *Computation* (2023).
- 458 13 Cho, E. H. & Lin, L.-C. Nanoporous material recognition via 3D convolutional neural
 459 networks: Prediction of adsorption properties. *The journal of physical chemistry letters*
 460 **12**, 2279-2285 (2021).
- 461 14 Xie, T. & Grossman, J. C. Crystal graph convolutional neural networks for an accurate
 462 and interpretable prediction of material properties. *Physical review letters* **120**, 145301
 463 (2018).
- 464 15 Chen, C., Ye, W., Zuo, Y., Zheng, C. & Ong, S. P. Graph networks as a universal
 465 machine learning framework for molecules and crystals. *Chemistry of Materials* **31**,
 466 3564-3572 (2019).
- 467 16 Janet, J. P. & Kulik, H. J. Resolving transition metal chemical space: Feature selection
 468 for machine learning and structure–property relationships. *The Journal of Physical*
 469 *Chemistry A* **121**, 8939-8954 (2017).
- 470 17 Bartók, A. P., Kondor, R. & Csányi, G. On representing chemical environments. *Physical*
 471 *Review B* **87**, 184115 (2013).
- 472 18 Kang, Y., Park, H., Smit, B. & Kim, J. A multi-modal pre-training transformer for
 473 universal transfer learning in metal–organic frameworks. *Nature Machine Intelligence* **5**,
 474 309-318, doi:10.1038/s42256-023-00628-2 (2023).

- 475 19 Moghadam, P. Z. *et al.* Development of a Cambridge Structural Database subset: a
476 collection of metal–organic frameworks for past, present, and future. *Chemistry of*
477 *Materials* **29**, 2618-2625 (2017).
- 478 20 Tong, M., Lan, Y., Yang, Q. & Zhong, C. Exploring the structure-property relationships
479 of covalent organic frameworks for noble gas separations. *Chemical Engineering Science*
480 **168**, 456-464 (2017).
- 481 21 Ongari, D., Yakutovich, A. V., Talirz, L. & Smit, B. Building a consistent and
482 reproducible database for adsorption evaluation in covalent–organic frameworks. *ACS*
483 *central science* **5**, 1663-1675 (2019).
- 484 22 Lu, W. *et al.* Porous polymer networks: synthesis, porosity, and applications in gas
485 storage/separation. *Chemistry of Materials* **22**, 5964-5972 (2010).
- 486 23 Haase, F. & Lotsch, B. V. Solving the COF trilemma: towards crystalline, stable and
487 functional covalent organic frameworks. *Chemical Society Reviews* **49**, 8469-8500
488 (2020).
- 489 24 Blatov, V. A., Blatova, O. A., Daeyaert, F. & Deem, M. W. Nanoporous materials with
490 predicted zeolite topologies. *RSC advances* **10**, 17760-17767 (2020).
- 491 25 Chung, Y. G. *et al.* Advances, updates, and analytics for the computation-ready,
492 experimental metal–organic framework database: CoRE MOF 2019. *Journal of Chemical*
493 *& Engineering Data* **64**, 5985-5998 (2019).
- 494 26 Anderson, R. & Gómez-Gualdrón, D. A. Increasing topological diversity during
495 computational “synthesis” of porous crystals: how and why. *CrystEngComm* **21**, 1653-
496 1665 (2019).
- 497 27 O’Keeffe, M., Peskov, M. A., Ramsden, S. J. & Yaghi, O. M. The reticular chemistry
498 structure resource (RCSR) database of, and symbols for, crystal nets. *Accounts of*
499 *chemical research* **41**, 1782-1789 (2008).
- 500 28 Lan, Y. *et al.* Materials genomics methods for high-throughput construction of COFs and
501 targeted synthesis. *Nature Communications* **9**, 5274 (2018).
- 502 29 Martin, R. L., Simon, C. M., Smit, B. & Haranczyk, M. In silico design of porous
503 polymer networks: high-throughput screening for methane storage materials. *Journal of*
504 *the American Chemical Society* **136**, 5006-5022 (2014).
- 505 30 Baerlocher, C. Database of zeolite structures. <http://www.iza-structure.org/databases/>
506 (2008).
- 507 31 Pophale, R., Cheeseman, P. A. & Deem, M. W. A database of new zeolite-like materials.
508 *Physical Chemistry Chemical Physics* **13**, 12407-12412 (2011).
- 509 32 Blatov, V. A., Shevchenko, A. P. & Proserpio, D. M. Applied topological analysis of
510 crystal structures with the program package ToposPro. *Crystal Growth & Design* **14**,
511 3576-3586 (2014).
- 512 33 Thompson, A. P. *et al.* LAMMPS-a flexible simulation tool for particle-based materials
513 modeling at the atomic, meso, and continuum scales. *Computer Physics Communications*
514 **271**, 108171 (2022).
- 515 34 Rappé, A. K., Casewit, C. J., Colwell, K., Goddard III, W. A. & Skiff, W. M. UFF, a full
516 periodic table force field for molecular mechanics and molecular dynamics simulations.
517 *Journal of the American chemical society* **114**, 10024-10035 (1992).
- 518 35 Bucior, B. J. *et al.* Energy-based descriptors to rapidly predict hydrogen storage in metal–
519 organic frameworks. *Molecular Systems Design & Engineering* **4**, 162-174 (2019).

520 36 Orhan, I. B., Daglar, H., Keskin, S., Le, T. C. & Babarao, R. Prediction of O₂/N₂
521 Selectivity in Metal–Organic Frameworks via High-Throughput Computational
522 Screening and Machine Learning. *ACS Applied Materials & Interfaces* **14**, 736-749
523 (2021).

524 37 Tibshirani, R. Regression shrinkage and selection via the lasso. *Journal of the Royal*
525 *Statistical Society: Series B (Methodological)* **58**, 267-288 (1996).

526 38 Rosen, A. S. *et al.* Machine learning the quantum-chemical properties of metal–organic
527 frameworks for accelerated materials discovery. *Matter* **4**, 1578-1597 (2021).

528 39 Moosavi, S. M. *et al.* Understanding the diversity of the metal-organic framework
529 ecosystem. *Nature communications* **11**, 1-10 (2020).

530 40 Nandy, A., Duan, C. & Kulik, H. J. Using machine learning and data mining to leverage
531 community knowledge for the engineering of stable metal–organic frameworks. *Journal*
532 *of the American Chemical Society* **143**, 17535-17547 (2021).

533 41 Mercado, R. *et al.* In silico design of 2D and 3D covalent organic frameworks for
534 methane storage applications. *Chemistry of Materials* **30**, 5069-5086 (2018).

535 42 Deeg, K. S. *et al.* In silico discovery of covalent organic frameworks for carbon capture.
536 *ACS applied materials & interfaces* **12**, 21559-21568 (2020).

537 43 Kim, B., Lee, S. & Kim, J. Inverse design of porous materials using artificial neural
538 networks. *Science advances* **6**, eaax9324 (2020).

539 44 Wang, Y., Yao, Q., Kwok, J. T. & Ni, L. M. Generalizing from a few examples: A survey
540 on few-shot learning. *ACM computing surveys (csur)* **53**, 1-34 (2020).

541 45 Mouriño, B., Jablonka, K. M., Ortega-Guerrero, A. & Smit, B. In search of Covalent
542 Organic Framework photocatalysts: A DFT-based screening approach. (2023).

543 46 Vaswani, A. *et al.* Attention is all you need. *Advances in neural information processing*
544 *systems* **30** (2017).

545 47 Cote, A. P. *et al.* Porous, crystalline, covalent organic frameworks. *science* **310**, 1166-
546 1170 (2005).

547 48 Kuhn, P., Antonietti, M. & Thomas, A. Porous, covalent triazine-based frameworks
548 prepared by ionothermal synthesis. *Angewandte Chemie International Edition* **47**, 3450-
549 3453 (2008).

550 49 Ding, S.-Y. *et al.* Construction of covalent organic framework for catalysis: Pd/COF-
551 LZU1 in Suzuki–Miyaura coupling reaction. *Journal of the American Chemical Society*
552 **133**, 19816-19822 (2011).

553 50 Wang, S., Xiong, S., Wang, Z. & Du, J. Rational Design of Zinc–Organic Coordination
554 Polymers Directed by N-Donor Co-ligands. *Chemistry–A European Journal* **17**, 8630-
555 8642 (2011).

556 51 Loshchilov, I. & Hutter, F. Decoupled weight decay regularization. *arXiv preprint*
557 *arXiv:1711.05101* (2017).

558 52 Dubbeldam, D., Calero, S., Ellis, D. E. & Snurr, R. Q. RASPA: molecular simulation
559 software for adsorption and diffusion in flexible nanoporous materials. *Molecular*
560 *Simulation* **42**, 81-101 (2016).

561 53 Feynman, R. P., Hibbs, A. R. & Styer, D. F. *Quantum mechanics and path integrals*.
562 (Courier Corporation, 2010).

563