

## PHYSICAL SCIENCES

## Machine learning enables interpretable discovery of innovative polymers for gas separation membranes

Jason Yang<sup>1†</sup>, Lei Tao<sup>2†</sup>, Jinlong He<sup>2†</sup>, Jeffrey R. McCutcheon<sup>3,4</sup>, Ying Li<sup>2,4\*</sup>

Polymer membranes perform innumerable separations with far-reaching environmental implications. Despite decades of research, design of new membrane materials remains a largely Edisonian process. To address this shortcoming, we demonstrate a generalizable, accurate machine learning (ML) implementation for the discovery of innovative polymers with ideal performance. Specifically, multitask ML models are trained on experimental data to link polymer chemistry to gas permeabilities of He, H<sub>2</sub>, O<sub>2</sub>, N<sub>2</sub>, CO<sub>2</sub>, and CH<sub>4</sub>. We interpret the ML models and extract valuable insights into the contributions of different chemical moieties to permeability and selectivity. We then screen over 9 million hypothetical polymers and identify thousands that lie well above current performance upper bounds, including hundreds of never-before-seen ultrapermeable polymer membranes with O<sub>2</sub> and CO<sub>2</sub> permeability greater than 10<sup>4</sup> and 10<sup>5</sup> Barrers, respectively. High-fidelity molecular dynamics simulations confirm the ML-predicted gas permeabilities of the promising candidates, which suggests that many can be translated to reality.

## INTRODUCTION

Polymer membranes are a flexible, processable, and inexpensive platform to provide a myriad of separations that fill critical roles in climate change mitigation (i.e., carbon capture) and resiliency (i.e., water treatment). For gas separations, polymer membranes have been widely used in the separation of mixtures in many industrial processes, including oxygen enrichment, biogas purification (1), and post-combustion carbon capture (2). In particular, carbon capture processes are garnering increased attention to reduce emissions to the environment, and membrane technologies offer known advantages such as high energy efficiency and operational simplicity due to flexibility and scalability (3). In post-combustion, pre-combustion, and oxy-combustion, CO<sub>2</sub>/N<sub>2</sub>, CO<sub>2</sub>/H<sub>2</sub>, and O<sub>2</sub>/N<sub>2</sub> separations are respectively important for environmental conservation.

During membrane-based gas separation, a gas mixture is typically driven through a membrane by pressure where separation is achieved through differences in individual gas permeabilities (4). The performance of membrane processes is determined by the membrane's permeability for a specific gas species,  $P_i$ , where  $i$  specifies the gas type. The membrane permeability is defined by Fick's law of diffusion,  $|J_i| = P_i \Delta p / l$ , where  $J_i$  is the flux of gas  $i$ , and  $\Delta p$  is the pressure drop across a membrane of thickness  $l$ . On the basis of the solution-diffusion model of gas transport in microporous membranes, permeability can alternatively be calculated as the product of diffusivity ( $D$ ) and solubility ( $S$ ):  $P_i = D_i \times S_i$ . When comparing the permeability of gas A with that of gas B, another performance measure is the membrane's selectivity between two gases,  $\alpha$ , which is defined as  $\alpha = P_A / P_B$ . An ideal membrane for a given binary gas separation would have high permeability and high selectivity. Increasing gas permeability and selectivity in these membranes would

allow for more efficient industrial processes by increasing the process throughput, reducing energy costs, and achieving a purer product (5, 6). However, there is a well-known permeability-selectivity trade-off for polymer gas separation membranes (4), which is defined by the Robeson upper bound (7). Over time, advancements in polymer designs have pushed the Robeson upper bound from 1991 values to updated 2008 values [and most recently 2015 values for O<sub>2</sub>/N<sub>2</sub> separations and 2019 values for CO<sub>2</sub>/CH<sub>4</sub> and CO<sub>2</sub>/N<sub>2</sub> separations; (8, 9)] that reflect improved membrane performance. Identifying new materials that break this upper bound has driven and continues to drive materials discovery efforts for membranes (10, 11).

In the decades of technological development in the membrane science field, design of new membrane materials has been, and remains, a largely trial-and-error process, guided by experience and intuition. Current approaches generally involve tuning chemical groups to increase affinity and solubility toward a desired gas or incorporating greater free volume to increase overall diffusivity. When assembling a new polymer, typically a desired enhancement is targeted (i.e., higher CO<sub>2</sub> affinity, higher overall permeability, and aging resistance), and a chemical group that is likely to achieve that enhancement is incorporated into the polymer chemistry (12–15). For achieving higher permeability, polymers of intrinsic microporosity (PIMs) have been extensively studied during the past two decades (16, 17). PIMs generally enhance fractional free volume via inefficient chain packing to increase permeability while simultaneously stiffening the polymer backbone and improving solubility selectivity (17, 18). Efforts to design improved chemistries for PIMs generally involve tuning the contortion group, increasing steric frustration via modifications to side chains, or further stiffening the polymer backbone (19–21). Still, many of these studies remain limited to an Edisonian approach, unable to identify or use big-picture rules of chemistry-property relationships in polymer membranes.

Further complicating matters, synthesis of new polymeric materials and subsequent testing of permeability and selectivity is a time-consuming, expensive, and incomplete process that can miss high-performance candidates. Molecular modeling approaches, such as Monte Carlo/molecular dynamics (MC/MD) simulations, can

<sup>1</sup>Division of Chemistry and Chemical Engineering, California Institute of Technology, Pasadena, CA 91125, USA. <sup>2</sup>Department of Mechanical Engineering, University of Connecticut, Storrs, CT 06269, USA. <sup>3</sup>Department of Chemical & Biomolecular Engineering, Center for Environmental Sciences and Engineering, University of Connecticut, Storrs, CT 06269, USA. <sup>4</sup>Polymer Program, Institute of Materials Science, University of Connecticut, Storrs, CT 06269, USA.

\*Corresponding author. Email: ying.li@uconn.edu

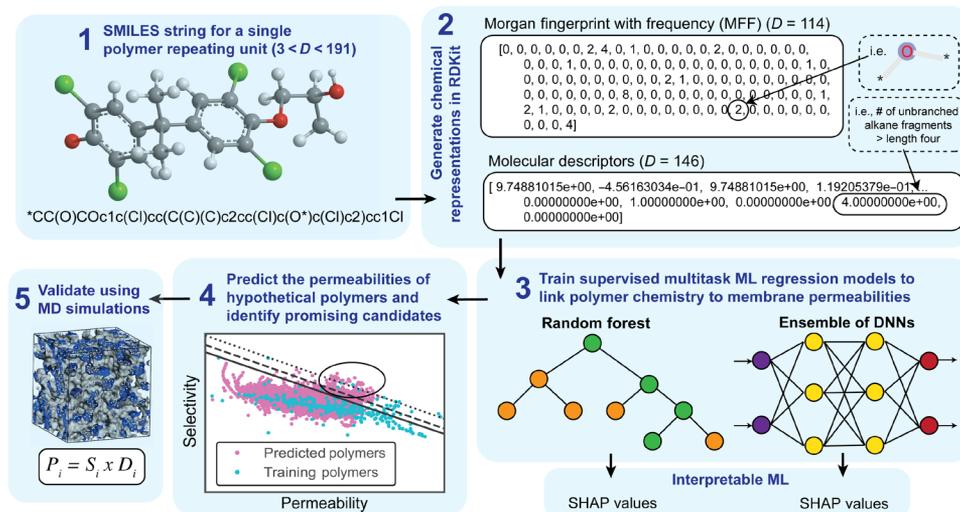
†These authors contributed equally to this work.

reasonably predict a polymer membrane's gas permeabilities without costly experiments (22–25). However, even these high-throughput molecular simulations are too computationally expensive to explore the vast chemical space of polymers on the order of  $10^6$  to  $10^{10}$ . By contrast, simplified approximations to predict gas permeability for a given membrane are low cost but inaccurate. Most simply, group contribution methods sum together the gas permeability contribution of each chemical moiety in a polymer, but they do not necessarily consider connectivity and cannot expand into new classes of polymers (26). Permeability can also be calculated via diffusivity based on the polymer's free volume and the solution-diffusion model of gas transport using various theoretical models, but these theories are incomplete (27, 28). In short, there is no efficient and accurate predictive model for gas permeability based on polymer-membrane chemistry.

Machine learning (ML) is a promising data-centric approach for prediction of gas permeabilities by learning a functional model based on polymer chemistry (29, 30). ML methods using chemical inputs have been successfully applied to accurately predicting many polymer properties including glass transition temperature (31–33), thermal conductivity (34), dielectric constants (35), organic photovoltaic properties (36, 37), and transport properties (38–40). The primary challenge for learning a generalizable ML model is training on robust and diverse data, which requires compiling multiple databases with the most recent literature values and imputing missing values (29). While Barnett *et al.* (39) have trained accurate ML models that link polymer chemistry to gas permeability, their training set notably lacks PIMs, and they only screened a limited chemical space of 11,000 existing homopolymers. Therefore, ML approaches would benefit from considering an expanded chemical space while simultaneously learning from additional training data on PIMs. Overall, ML-directed molecular design of polymer membranes still faces substantial challenges in the following aspects: (i) How can we

define an appropriate chemical space to explore the molecular design of high-performance polymer membranes? (ii) Even if ML models can be established for the gas permeability prediction of polymer membranes, how can we achieve a physical understanding of how membrane chemistry affects gas separations? (iii) Can we exceed the Robeson upper bound simultaneously for separations of different gas pairs, such as  $O_2/N_2$ ,  $CO_2/CH_4$ ,  $CO_2/N_2$ , and  $H_2/CO_2$ ?

To tackle the above challenges, we demonstrate interpretable, supervised ML models that can accurately predict the He,  $H_2$ ,  $O_2$ ,  $N_2$ ,  $CO_2$ , and  $CH_4$  permeabilities of gas separation membranes based on polymer chemistry—as part of our ML-assisted discovery workflow outlined in Fig. 1. Our training data consists of polymer chemistry and experimental gas permeabilities from two large databases, PoLyInfo and Membrane Society of Australasia (MSA). In these datasets, hundreds of homopolymers, including PIMs, are identified by their SMILES strings—a notation for chemical structures that represents a molecule as a unique string of ASCII characters. We use two representations for the polymer repeating unit, namely, chemical descriptors as generated by RDKit (listed in table S1) and the Morgan fingerprint with frequency (MFF) (41), which captures the frequency of chemical moieties (substructures) present in molecules. We then train multitask supervised ML models to establish synthesis-property relations for these polymer membranes. While various supervised ML models have been used in polymer informatics, including recurrent neural networks (RNNs), support vector machines, Gaussian processes, and others, we choose to focus our study to random forest (RF) regression and deep neural networks (DNNs), which have demonstrated outstanding performance in our recent benchmark study (33). Because of the high variance of DNNs, we perform bootstrap ensembling to further improve our predictions—achieving test  $R^2 \sim 0.90$  between predicted and actual permeability values. We also interpret our ML models by extracting feature importances using SHAP (SHapley Additive exPlanations)



**Fig. 1. Workflow for ML-assisted discovery of innovative polymer membranes with ideal gas separation performance, e.g., beyond the traditional Robeson upper bound.** (1) We begin with the SMILES string for each polymer's repeating unit and its associated gas permeabilities for model training. (2) Each polymer's relevant fingerprint substructures (MFFs) and molecular descriptors are extracted, which are used as chemical inputs for ML model training. The two representations have a dimensionality  $D$  on the order of 100. (3) Multitask RFs and ensembles of DNNs are trained to predict gas permeabilities, and physical insights can be extracted from the models using their SHAP values. (4) The models are used for high-throughput permeability prediction of hypothetical polymers with unknown permeabilities but known chemistries in a substantially expanded chemical space. (5) Last, high-fidelity MD simulations are performed to verify the membrane permeabilities/selectivities of top polymer candidates.

(42). Our analysis provides a chemical explanation for the well-known permeability-selectivity trade-off in membranes, and many of the other physical insights that we draw are consistent with established membrane design principles. Using the trained ML models, we perform high-throughput screening of over 9 million hypothetical polymers with unknown permeabilities, including many polyimides and ladder polymers that can be classified as PIMs. Thus, we identify thousands of promising polymers for gas separation membranes with desirable performance, which lie well beyond current performance upper bounds. These polymers are characterized by many of the chemical substructures revealed as beneficial from SHAP analysis. Last, we perform high-fidelity MD simulations to confirm that the ML-predicted permeabilities of top-performing polymers are very accurate. Overall, our ML-assisted workflow is a promising method to facilitate the discovery of innovative polymers for next-generation gas separation membranes to advance energy and environmental sustainability.

## RESULTS

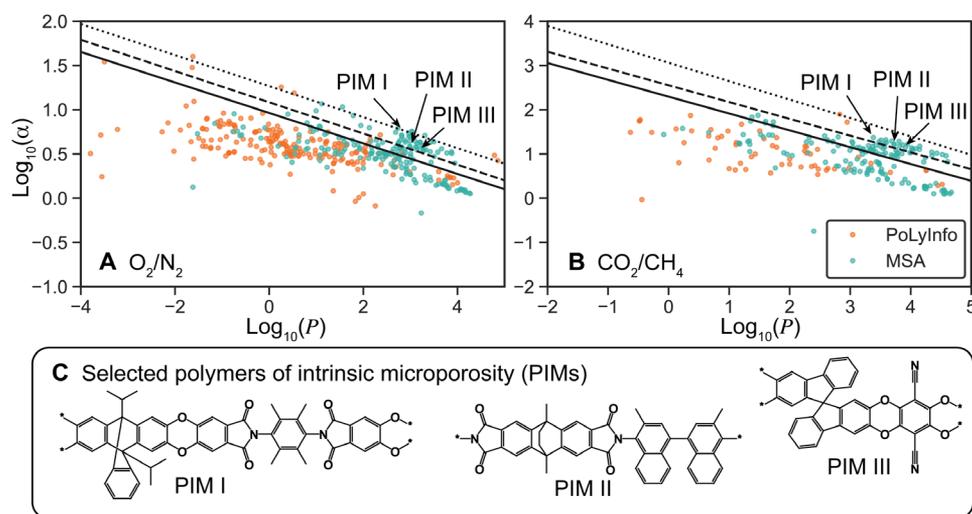
### Datasets and chemical space under exploration

Our training dataset, dataset A, consists of 778 homopolymers (353 unique polymer chemistries), linked to at least one gas permeability reported among He, H<sub>2</sub>, O<sub>2</sub>, N<sub>2</sub>, CO<sub>2</sub>, and CH<sub>4</sub>. Dataset A is manually collected from the PoLyInfo database (experimental data from before 2005) and is merged with data from the MSA database (beyond 2005). As shown in Fig. 2 (A and B), in general, the more recent MSA database contains polymers with higher permeability, e.g., CO<sub>2</sub> permeability greater than 10<sup>3</sup> Barrers, and entries that surpass the 2008 Robeson upper bound. In these plots, we identify several known PIMs, with their corresponding chemical structures shown in Fig. 2C. Many of these PIMs are ladder polymers, which have two connection points between consecutive monomers, such as PIM I and PIM II. Some of the polymers are polyimides, such as PIM III.

Table 1 provides a summary of the training and screening datasets used in this work. While dataset A is used for training, we additionally

generate three unique datasets for screening (polymer discovery): datasets B, C, and D. Dataset B, PIIM, consists of polymers learned via an RNN trained on SMILES strings of existing polymers in PoLyInfo, as constructed by Ma and Luo (43). Note that dataset B covers an overlapping chemical space with dataset A because the RNN model is also trained on the PoLyInfo database, but dataset B densely populates regions where PoLyInfo data are sparse (43). Still, dataset B mostly spans polymers similar to known polymers, which are generally not tailored for membrane separations. Thus, our motivation for constructing datasets C and D is for rationally targeted exploration of the polymer design space, based on the established interest in PIMs within the membrane design community. First, polyimides have garnered particular attention due to their superior permeability/selectivity trade-off and high chemical and physical stability, largely due to a rigid aromatic backbone (13, 44). Thus, dataset C is constructed as 8 million hypothetical polyimides formed by the polycondensation of known diamines/diisocyanates with dianhydrides from the PubChem library (45). These 8 million hypothetical polyimides substantially expand the current chemical space of around 2000 polyimides in PoLyInfo. Ladder polymers adopt an alternative approach to stiffen the polymer backbone. These unique polymers have two-bond connections between repeating units and thus restricted rotation, except at a conformation site, which is often a spirocenter (46). In our work, dataset D contains hypothetical ladder polymers generated through the binary combinations of components of existing ladder polymers (47), supplemented by an RNN model. More details about the construction of datasets C and D are provided in fig. S1.

While we train ML models using both chemical descriptors and MFFs as inputs, for simplicity, we only screen new polymers using MFFs as inputs for ML models. The feature spaces for MFFs across the datasets studied in this work are visualized using uniform manifold approximation and projection (UMAP) (48) in fig. S2. In general, our training set, dataset A, spans across the screening space of datasets B, C, and D. Thus, our ML models can learn across a wide chemical feature space of interest. While datasets B and C have



**Fig. 2. Visualization of the permeability distribution of dataset A, the training set.** (A) O<sub>2</sub>/N<sub>2</sub> and (B) CO<sub>2</sub>/CH<sub>4</sub> Robeson plots for the raw data as obtained from the PoLyInfo and MSA databases. Units of permeability are Barrers. (C) Chemical structures of three existing examples of PIMs, with their performances identified in (A) and (B). Asterisks in chemical structures indicate connection points between repeating units.

**Table 1. Summary of the datasets explored in this work.** Dataset A is the training set, which contains polymers with known chemistries and permeabilities. Datasets B, C, and D contain hypothetical polymers with known chemistries but unknown permeabilities (used for screening/discovery). Datasets B, C, and D span three different chemical spaces: PoLYInfo-like polymers, polyimides, and ladder polymers, respectively.

	No. of polymers	Permeabilities	Description	Source
Dataset A	778 (353 unique)	At least one gas known	Training set	PoLYInfo and MSA databases
Dataset B	995,799	Unknown	PI1M (43)	Hypothetical polymers generated from PoLYInfo through a recurrent neural network
Dataset C	8,205,087	Unknown	Polyimides	Hypothetical polyimides formed by known dianhydride and diamine/diisocyanate pairs from PubChem (45)
Dataset D	1,124	Unknown	Ladder polymers	Hypothetical polymers generated based on existing ladder polymers

more complete coverage due to the sheer number of samples, dataset D only includes ladder polymers, which explains why they are more confined in the feature space.

### Performance of ML models for gas permeability prediction

To quantify performance, we evaluate the accuracy and generalizability of our ML models, namely, RFs and DNN ensembles trained on chemical descriptors and MFFs. For our supervised ML models, the metric of study is the  $R^2$  correlation between the predicted and actual permeabilities on the training and test sets, as summarized in Table 2. Before training, we impute missing permeabilities using the multivariable imputation by chained equations (MICE) algorithm (49). Visualizations of the results of permeability imputation can be found in fig. S3, which show that there is not a notable difference when missing gas permeabilities are imputed via extremely randomized trees (ERTs) versus Bayesian linear regression (BLR). Thus, we focus our analysis to models trained on the permeabilities imputed via BLR for consistency. First, we find that the choice of ML model is more important than the choice of chemical features. The average test  $R^2$  across all six gases for the RF is approximately 0.74 when trained on descriptors and very similar when trained on MFFs. Similarly, the test  $R^2$  values for the DNN ensembles are around 0.90 for both descriptors and MFFs. We infer that MFFs offer slightly better performance, which has been observed in the prediction of polymer glass transition temperature (31).

By contrast, we find that the choice of ML model has a meaningful impact on performance. The RF learns a model with train  $R^2$ s of about 0.96 on descriptors and 0.90 on MFFs, which reduce to test  $R^2$ s of about 0.74 for both inputs. In particular, the RF model seems to struggle to fit the data points with very low or very high permeabilities, as demonstrated in fig. S4 by the points that have a high actual permeability but lie below the unit line. This would suggest that the RF does not prioritize fitting to the PIMs with high gas permeabilities in the training set, as PIMs make up a relatively small fraction of the training data and tend to have distinct chemistry compared to the rest of the training set.

On the other hand, the DNN ensemble learns a model with train  $R^2$ s of around 0.87 on descriptors and 0.89 on MFFs, which generalizes very well to test  $R^2$ s of approximately 0.89 for both inputs. The similarity between train and test  $R^2$ s for the DNN ensembles suggests that the model is very generalizable and learns the underlying functional relationship between chemistry and permeability. In fig. S5, we note that the DNN ensemble generally predicts permeability reasonably well, although there are some outliers at low or high

permeability. Performance quantification in table S2 further suggests that the DNN ensemble generalizes well, even for zero-shot predictions on an unseen distribution of polymers. In addition, uncertainty quantification of the DNN ensemble reveals that the ensemble of models performs better than the sum of its parts, as the average test  $R^2$  for each individual model is only  $\sim 0.70$ . Overall, there is also around 10% average normalized variance in the predicted permeabilities across the 16 DNN models, which is quite high.

### Chemical insights from interpretation of ML models

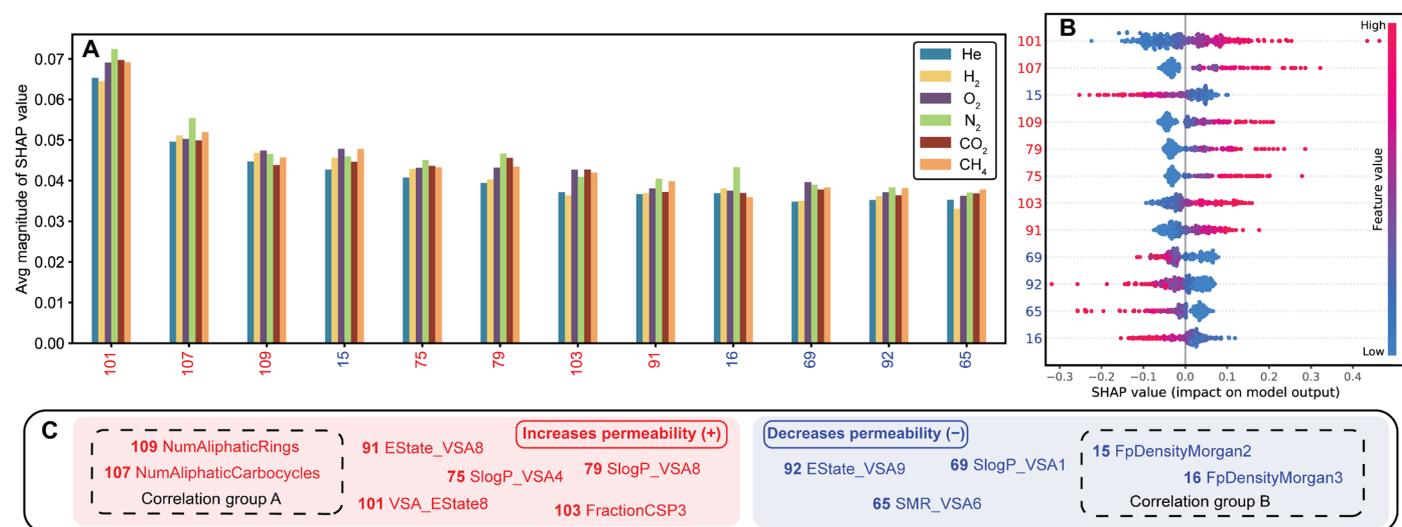
Usually, ML models are treated as black boxes, which makes it challenging to understand any physical principles learned by the models. However, we find that obtaining SHAP values from our ML models on chemical descriptors and MFFs makes our models not only accurate but also interpretable. By extracting the most important chemical features that predict gas permeability, we draw physical insights into the molecular design of polymer membranes. Here, we decide to focus our analysis on the DNN ensemble because of its better performance, but other model types can also be explained using the same method.

Figure 3 summarizes the results of SHAP analysis on the DNN ensemble trained on chemical descriptors. Figure 3A highlights the 12 most important chemical descriptors based on their average SHAP values—their relative impacts on the six gas permeabilities under study. A summary of the definitions of these descriptors is provided in table S3. The most important descriptor is VSA\_EState8, which is a hybrid electronic state and van der Waals surface area (VSA) descriptor, based on precalculated surface area values derived from a list of functional groups. While some of these descriptors do not have obvious, intuitive physical meaning, the permeability of polymer membranes is determined by the solubility and diffusivity of gas molecules (24), which is affected by the electrostatic interactions and free-volume elements, respectively. Therefore, these identified chemical descriptors should play important roles in gas permeability of polymer membranes.

In Fig. 3B, we show how the values of each of the top descriptors affect the model's  $\text{CH}_4$  permeability prediction. If higher feature values result in more positive SHAP values, then the feature has a positive effect on permeability: The feature is directly correlated with gas permeability. On the other hand, if higher feature values result in more negative SHAP values, then the feature has a negative effect on permeability: The feature is inversely correlated with gas permeability. While we only show the impacts of features on  $\text{CH}_4$  permeability output, we draw the same conclusions from the other

**Table 2. Summary of the performances of supervised ML models as scored by the  $R^2$  value between the predicted and actual permeabilities.** All ML models make multitask predictions for the six gas permeabilities of He, H<sub>2</sub>, O<sub>2</sub>, N<sub>2</sub>, CO<sub>2</sub>, and CH<sub>4</sub> and are trained on the data that are augmented using BLR imputation. The DNN ensemble models perform better than the RF models, and models trained on MFFs perform slightly better than models trained on molecular descriptors.

ML model	He		H <sub>2</sub>		O <sub>2</sub>		N <sub>2</sub>		CO <sub>2</sub>		CH <sub>4</sub>	
	Train	Test	Train	Test	Train	Test	Train	Test	Train	Test	Train	Test
RF (descriptors)	0.96	0.73	0.96	0.74	0.96	0.75	0.96	0.74	0.96	0.75	0.96	0.74
DNN ensemble (descriptors)	0.85	0.87	0.87	0.88	0.89	0.89	0.90	0.90	0.88	0.90	0.89	0.89
RF (MFFs)	0.89	0.73	0.89	0.74	0.89	0.74	0.90	0.74	0.89	0.75	0.90	0.74
DNN ensemble (MFFs)	0.88	0.91	0.88	0.90	0.90	0.92	0.90	0.91	0.89	0.90	0.89	0.88



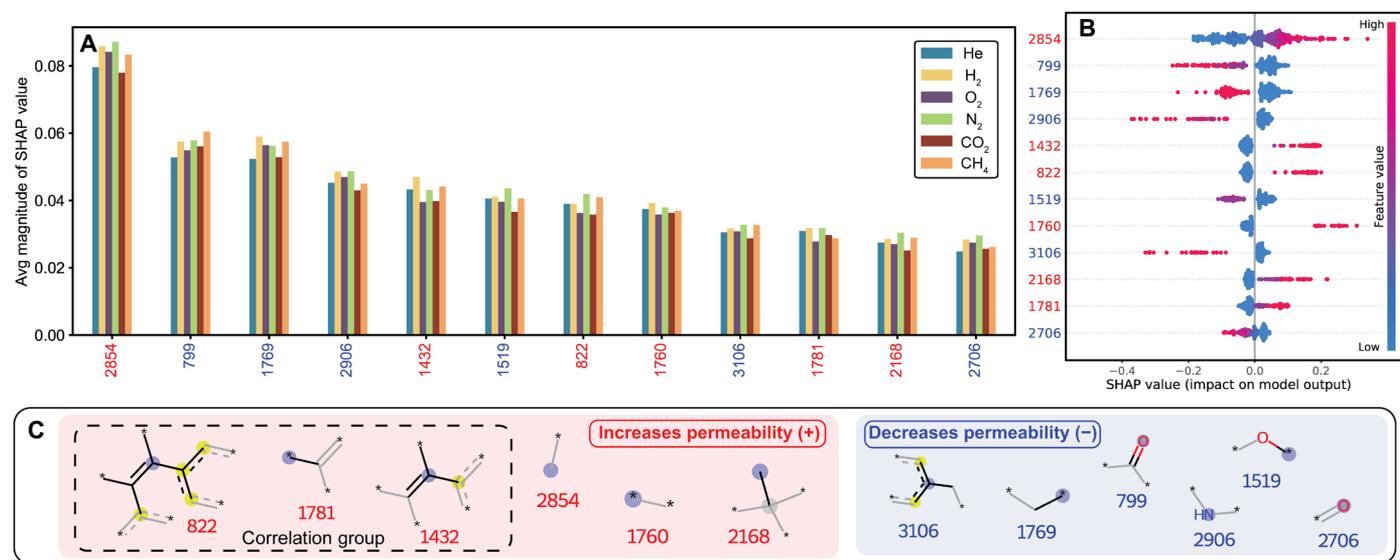
**Fig. 3. Important molecular descriptors as identified using SHAP on the DNN ensemble ML model trained on descriptors and BLR-imputed permeabilities.** (A) Average SHAP importances for the top 12 descriptors on each of the six gas permeabilities (He, H<sub>2</sub>, O<sub>2</sub>, N<sub>2</sub>, CO<sub>2</sub>, and CH<sub>4</sub>). (B) Impact of the top 12 descriptors on CH<sub>4</sub> permeability output. Each dot represents the impact of a particular sample in the training set. (C) Names of top descriptors, with highly correlated features circled. Red text signifies features that have positive effects on permeability, and blue signifies a negative effect.

gas permeability predictions, which produce almost identical SHAP impact plots (fig. S6), as all permeabilities are trained via the same multitask model.

On the basis of the correlation matrix between descriptor features in fig. S7, there are two main pairs of correlated features, which suggests that some of the top features do not have independent physical significance. Namely, descriptors 107 and 109 (Aliphatic Cycle Counts, correlation group A) are highly correlated with one another and relatively anticorrelated with features 15 and 16 (FpDensityMorgan, correlation group B). The steric space occupied by rings generally results in a lower molecular density, which explains the opposition between group A and group B. The features in group A have a positive impact on gas permeability, while the features in group B have a negative impact on gas permeability. This suggests that repeating units with more nonaromatic rings allow for larger free-volume elements and lower densities, thereby higher gas permeabilities. This supports the emerging direction of polymer research on nonplanar structures, such as *kink*, *spiro*, *cardo*, and pendant groups ( $-\text{CF}_3$ ), bulky and flexible groups ( $-\text{O}-$ ), or different

spatial linkage configurations in polyimides for enlarging their microporosities (19, 50).

In Fig. 4, we perform the same type of feature importance analysis using SHAP values, for the DNN model trained on MFFs. Here, we highlight the most important chemical substructures in the prediction of gas permeability. As shown in Fig. 4A, the most important substructure overall is 2854, the methyl group. We believe that this feature facilitates permeability because it is hydrophobic, and its shape contributes to steric frustration between polymer chains. Similarly, the quaternary carbon connected to an aliphatic ring (substructure 2168) contributes to increasing permeability, which supports our findings above. The DNN model also learns that the number of connection bonds, substructure 1781, is correlated with gas permeability, because many high-permeability PIMs are ladder polymers with four connection points per repeating unit, as opposed to two for a typical polymer. The correlation matrix between chemical substructures (fig. S9) suggests that most of the important substructure features are independent of one another. However, substructures 1781, 1432, and 822 are highly correlated and all have



**Fig. 4. Important molecular substructures as identified using SHAP on the DNN ensemble ML model trained on MFFs and BLR-imputed permeabilities. (A)** Average SHAP importances for the top 12 substructures for each of the six gas permeabilities (He, H<sub>2</sub>, O<sub>2</sub>, N<sub>2</sub>, CO<sub>2</sub>, and CH<sub>4</sub>). **(B)** Impact of top 12 substructures on CH<sub>4</sub> permeability output. Each dot represents the impact of a particular sample in the training set. **(C)** Illustration of top substructures, with correlated features circled. Red text signifies features that have positive effects on permeability, and blue signifies a negative effect. In the substructure drawings, blue highlights the central atom in the environment, yellow indicates aromatic atoms, and gray indicates aliphatic ring atoms.

a positive relation to gas permeability (Fig. 4B). Upon closer examination, we find that substructure 1781 is contained within substructure 1432, which is contained in substructure 822. Substructures 1432 and 822, two double-bonded carbons connected to an aromatic ring, define polyacetylenes, which demonstrate some of the highest permeabilities among nonporous polymers in gas separations (51). By contrast, polar groups generally have negative contributions to gas permeability, as shown in Fig. 4B. For example, double-bonded oxygens (799 and 2706), ethers (1519), and nitrogen atoms (2906) are all inversely correlated with gas permeability. Since most gas molecules are nonpolar, the presence of these polar groups generally reduces the solubility of gases, which explains the negative effect.

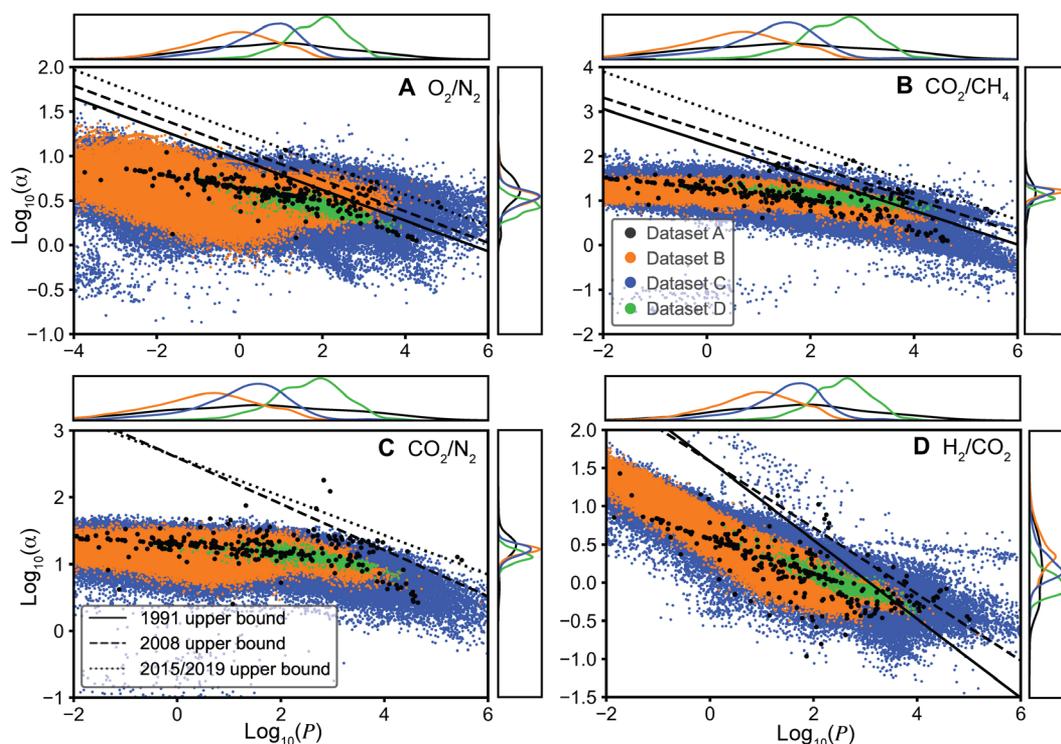
However, in Fig. 4A, our ML models show that these groups tend to have a greater negative impact (measured by SHAP value) on N<sub>2</sub> and CH<sub>4</sub> permeability, compared to O<sub>2</sub> and CO<sub>2</sub>, which explains why the presence of these groups in polyimides, ladder polymers, and poly(ethylene oxides) (15) can increase selectivity by widening the permeability difference between certain pairs of gases, which is desirable for gas separations. This supports a known heuristic in membrane design—that CO<sub>2</sub> selectivity can be increased via increased CO<sub>2</sub> solubility by incorporating oxygen atoms into polymer membranes (15, 52). Across the board for substructures, SHAP values tend to be higher for N<sub>2</sub> and CH<sub>4</sub> compared to O<sub>2</sub> and CO<sub>2</sub> (Fig. 4A), which suggests that incorporating chemistries that increase permeability (i.e., methyl groups) is likely to come at the cost of selectivity. Our ML models thus elucidate a chemical basis for the permeability/selectivity trade-off: Chemical features that increase permeability are likely to do so to a greater extent for molecules that are less permeable (N<sub>2</sub> and CH<sub>4</sub>), but chemical features that reduce permeability are also likely to affect these molecules to a greater magnitude—thereby increasing selectivity for the more permeable gas (O<sub>2</sub> and CO<sub>2</sub>). Achieving high permeability and selectivity thus becomes a balancing act. This unique understanding is unlocked from the ability of ML to learn complex patterns in data.

### Discovery of high-performance polymers and validation through MD simulations

After training our RF and DNN ensemble ML models, we use the models based on MFFs for high-throughput screening and discovery of high-performance polymers for gas separations. We choose the ML models using MFFs for simplicity due to their slightly better performance and lower system memory requirements. We calculate MFFs for millions of hypothetical polymers in datasets B, C, and D, which span a wide and relevant chemical space. These inputs are then passed through the RF and DNN ensemble models in a feed-forward manner. The predicted permeabilities for the DNN model, broken down by dataset, are plotted for O<sub>2</sub>/N<sub>2</sub>, CO<sub>2</sub>/CH<sub>4</sub>, CO<sub>2</sub>/N<sub>2</sub>, and H<sub>2</sub>/CO<sub>2</sub> separations in Fig. 5. Similarly, the predictions for the RF model are visualized in fig. S10.

Broadly, we find that the RF model predicts permeabilities in a much narrower space than the DNN ensemble, which explains its lower *R*<sup>2</sup> values on the test set and further supports the observation that the DNN ensemble is more accurate and generalizable. Predicted permeabilities for each screening dataset lie in their expected region in the permeability-selectivity space, which further supports the accuracy of our ML models. Namely, both models predict permeabilities close to existing Robeson upper bounds for polymers in dataset D, which consists entirely of ladder polymers (a subclass of PIMs). Similarly, dataset C consists of polyimides (including many PIMs), and their permeability predictions span a space that includes polymers below and above the Robeson upper bound, reflecting the dataset's diversity. However, dataset B corresponds to mostly polymers with low permeability and selectivity. We believe that this can be explained by the fact that PoLyInfo is a broad database that contains many polymers that are not suitable for gas separation applications, and dataset B is populated from existing polymers in PoLyInfo.

Most promisingly, the DNN model predicts thousands of polymers from dataset C to be above the 2008 Robeson upper bound, for O<sub>2</sub>/N<sub>2</sub>, CO<sub>2</sub>/CH<sub>4</sub>, CO<sub>2</sub>/N<sub>2</sub>, and H<sub>2</sub>/CO<sub>2</sub> separations, which are



**Fig. 5. Visualization of predicted permeabilities for hypothetical polymers in datasets B, C, and D, based on the ensemble of DNNs trained on MFFs with BLR-imputed permeabilities.** The training dataset (dataset A) is overlaid on the predicted permeabilities. The data are visualized for (A)  $O_2/N_2$ , (B)  $CO_2/CH_4$ , (C)  $CO_2/N_2$ , and (D)  $H_2/CO_2$  separations, with thousands of promising polymers lying at or above the Robeson upper bounds. Units of permeability are Barrers.

summarized in Table 3. We further find that the DNN ensemble trained on MFFs not only generalizes but also extrapolates. We find a class of hypothetical polymers in dataset C with never-before-seen ultrahigh  $CO_2$  permeability (greater than  $10^5$  Barrers) and a class of polymers with ultrahigh  $O_2$  permeability (greater than  $10^4$  Barrers)—even though our training set only contains 12 polymers with  $O_2$  permeability greater than  $10^4$  Barrers and only 2 polymers with  $CO_2$  permeability greater than  $10^5$  Barrers.

We select a handful of hypothetical polymers with high predicted permeability and selectivity across all four separations under study to validate their performance using MD simulations to calculate permeability. SMILES strings, synthetic accessibilities (53), and solubility scores for the components of these polymers are given in table S4 and fig. S11. These polymers demonstrate reasonable synthetic accessibility scores, and the polyimides are formed from known PubChem chemicals and compounds, which suggest that the syntheses of these promising candidates are feasible through the polycondensation of existing dianhydride and diamine/diisocyanate molecules. The predicted solubility parameters of these polymers in various solvents are also given in tables S5 to S12, which further suggests that many may be solution processable to form membranes.

Details of intermediate values calculated during the atomistic simulation process can be found in table S13. Note that our MD simulation protocol has been benchmarked, with excellent agreement with the literature, against experimental and simulation values for the gas permeabilities of PIM-1 (46) (table S14) and 10 other relevant polyimide/ladder polymers that are chemically similar to the polymer candidates with promising performance (fig. S12 and table S15).

Both the RF and DNN model can identify polymers with high performance in datasets C and D. The chemical structures of the selected polymers are drawn in Fig. 6A. We highlight some of the top substructures identified from SHAP analysis (Fig. 4) in these chemical drawings, which corroborates our earlier conclusions. For instance, the higher permeability polymers tend to have more methyl groups (substructure 2854) and methyl groups attached to aliphatic rings (substructure 2168) to increase steric frustration. Meanwhile, double-bonded oxygens (substructure 2706) in the polyimide backbone help to maintain selectivity for gases such as  $O_2$  and  $CO_2$ .

As shown in Fig. 6 (B to E), ML-predicted performances lie very close to their respective MD-simulated performances for separations involving  $O_2$ ,  $N_2$ ,  $CO_2$ ,  $CH_4$ , and  $H_2$ . Error ranges for permeability calculations from simulations and predictions from ML models are provided in table S16. In general, the DNN ensemble model predictions differ less from the values given by MD simulations, compared to those of the RF model. While the permeability predictions tend to have larger error and uncertainty as the DNN model extrapolates to higher permeability values, our experimentally validated MD simulations confirm the predicted performances of these top candidates. Thus, thousands of polymers in our screening datasets with predicted permeabilities above the Robeson upper bound, or ultrahigh predicted permeabilities, could translate to real polymer membranes with exceptional separation performance.

Notably, P-DNN-C3 and P-DNN-C4, hypothetical polyimides in dataset C, demonstrate  $O_2/N_2$  selectivity well beyond the 2015 upper bound of existing known polymers, as predicted by our DNN ensemble model and further validated by MD simulations. To our knowledge, these previously unidentified discoveries have the highest  $O_2/N_2$

**Table 3. Summary of the number of polymers with exceptional performance found from dataset C.** Dataset C contains about 8 million hypothetical polyimides formed by known dianhydride and diamine/diisocyanate pairs from PubChem (45).

	Gas or separation	No. of polymers
Above 2008 Robeson upper bound	O <sub>2</sub> /N <sub>2</sub>	~80,000
	CO <sub>2</sub> /CH <sub>4</sub>	~3,000
	CO <sub>2</sub> /N <sub>2</sub>	~800
	H <sub>2</sub> /CO <sub>2</sub>	~10,000
Permeability above 10 <sup>4</sup> Barrers	O <sub>2</sub>	197
Permeability above 10 <sup>5</sup> Barrers	CO <sub>2</sub>	225

selectivities for their respective O<sub>2</sub> permeabilities, found to date. These hypothetical polyimides can each be formed through the polycondensation of a previously synthesized diisocyanate with a dianhydride, as shown in table S4.

Because of the multitask nature of the DNN ensemble ML model, many polymers are predicted to perform well across several metrics. For example, P-DNN-C3 surpasses current upper bounds for O<sub>2</sub>/N<sub>2</sub> and H<sub>2</sub>/CO<sub>2</sub> separations, while P-DNN-C4 demonstrates exceptional performance for O<sub>2</sub>/N<sub>2</sub> and CO<sub>2</sub>/N<sub>2</sub> separations. However, these two polymers have poor CO<sub>2</sub>/CH<sub>4</sub> selectivity. P-DNN-C1 performs near or above the 2008 Robeson upper bounds for O<sub>2</sub>/N<sub>2</sub>, CO<sub>2</sub>/CH<sub>4</sub>, and CO<sub>2</sub>/N<sub>2</sub> separations. In another vein, P-DNN-C2 has both ultrahigh O<sub>2</sub> permeability and CO<sub>2</sub> permeability while maintaining high selectivity. Similarly, P-DNN-C1 and P-DNN-C2 can each be formed through the polycondensation of a known diamine and a dianhydride, as given in table S4.

To further investigate the superior permeability of the selected top polymer candidates, we generate their realistic structural models and analyze their pore structures *in silico* in fig. S13. In comparison with PIM-1, our top candidates have more voids, enhanced microporosity and larger pore radii. The pore size distribution of the top candidate polymers is wider and shifted to the right, further suggesting enhanced microporosity and permeability. Table S17 validates that our computational method generalizes accurately to other polymers with smaller free-volume elements, such as flexible fluorinated polyimides and semicrystalline polymers.

## DISCUSSION

In this work, we demonstrate an accurate and cost-effective ML implementation that can effectively explore the ever-expanding design space for polymeric gas separation membrane materials, by learning their synthesis-property relationships. First, our study reveals that fixed chemical descriptors and MFFs are both excellent representations for predicting gas permeabilities of polymer membranes. Corroborating our benchmark study on polymer glass transition temperature (33), we conclude that the choice of chemical representation generally plays a limited role in each ML model's performance, as long as sufficient chemical substructures are captured. Additional features, such as microstructure, could be considered in future ML models, given the importance of microstructural characteristics such as free-volume elements in the solution-diffusion transport theory of membranes (54). Incorporation of these characteristics as

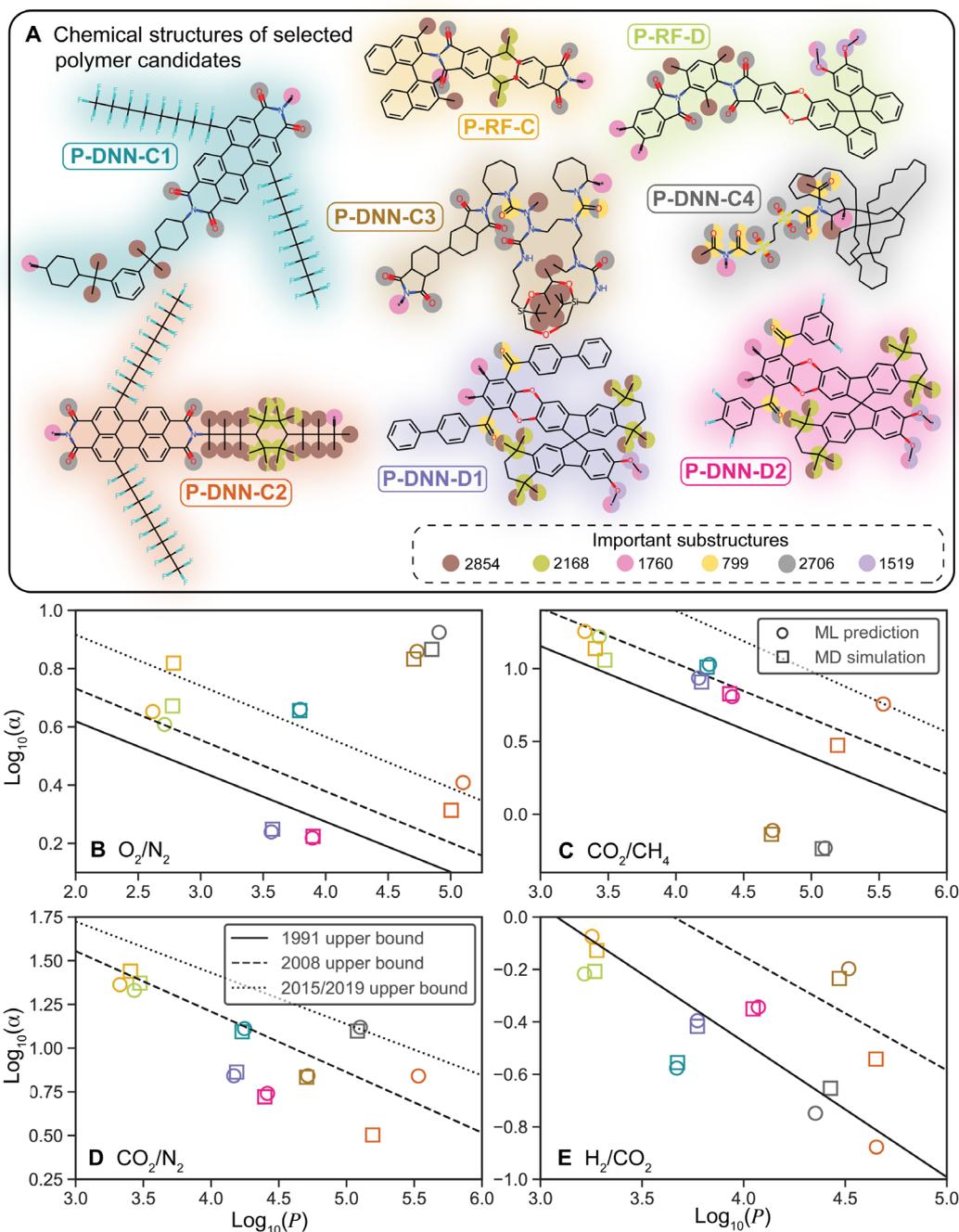
input features has improved metal-organic framework adsorption prediction (55), compared to using solely chemical descriptors. These microstructural features could be efficiently calculated via MD simulations, as demonstrated in this work. Alternatively, because high-throughput MD simulations can calculate gas permeabilities with reasonable accuracy, these simulations could also be used to augment the training set or be incorporated into active learning frameworks to reduce the uncertainty of ML models (56). Nevertheless, we find that fixed chemical features can sufficiently predict the gas permeabilities of the polymer membranes studied here.

We additionally gain insight into how the choice of ML model affects performance. At the same time, we demonstrate that ensembling is a powerful technique for improving prediction accuracy while simultaneously quantifying uncertainty. Traditionally, RF models are thought to work better on small datasets, while deep learning is reserved for large training sets. However, while decision trees are adequate for capturing simple relationships, neural networks can, in principle, approximate any function to arbitrary accuracy. In our study, we demonstrate that deep learning can be effectively applied to small training datasets on the order of a few hundred training samples. We believe that our DNN method is accurate and generalizable because of ensembling. Each DNN model, seeing a limited subsample of the data, captures complexities and nuances, which results in individual predictions with high variance; however, the overall model generalizes well when predictions are averaged together via ensembling. Training the 16 DNNs in our study and evaluating predictions for millions of samples are still computationally tractable on a personal computer. Although various other neural networks such as graph neural networks are garnering increased interest for certain molecular discovery and synthesis tasks (57), we do not observe notable performance gains from training graph convolutional, recurrent, or convolutional neural networks. We have reached a similar conclusion from our polymer informatics benchmark study on polymer glass transition (33). In short, we believe that deep learning techniques, even standard multilayer perceptrons, have much broader applicability to small datasets of chemical features than previously assumed.

We show that SHAP analysis can succinctly elucidate the impacts of input features, even for complex nonlinear models, which erodes the paradigm that ML models are black boxes (58). SHAP values can be calculated for nearly all supervised ML models, and we encourage future chemical and polymer informatics studies to take advantage of explainability in ML (59). A recent study also used coloring of substructures when training a graph neural network for interpretable ML (60), which suggests that feature importance analysis of ML models can be extended beyond fixed representations to learned chemical representations.

Our study of fixed feature importances solidifies many existing membrane design principles but additionally offers unique, generalized guidance for the molecular engineering of new polymers for gas separations. Broadly, SHAP analysis illuminates the chemical balancing act required for overcoming the permeability/selectivity trade-off. Polymers must juggle (i) the number of bulky chemical moieties—i.e., methyl groups and aliphatic rings—that increase microporosity (permeability at the expense of selectivity) with (ii) the number of polar groups—i.e., carbonyls and oxygens—that increase relative CO<sub>2</sub> and O<sub>2</sub> affinity (selectivity at the expense of permeability).

P-DNN-C3 and P-DNN-C4 are case studies into this balancing act. They achieve high permeability primarily through methyl groups



**Fig. 6. Validation, using MD simulations, of the performance of selected top polymer candidates from the ML models trained on MFFs with BLR-imputed permeabilities.** (A) Chemical structures for the selected polymer candidates with high performance. Important chemical substructures are highlighted in the molecules. Asterisks in chemical structures indicate connection points between repeating units. The predictions from ML models are shown as circles, while corresponding MD simulation values are shown as squares, for (B) O<sub>2</sub>/N<sub>2</sub>, (C) CO<sub>2</sub>/CH<sub>4</sub>, (D) CO<sub>2</sub>/N<sub>2</sub>, and (E) H<sub>2</sub>/CO<sub>2</sub> separations. Units of permeability are Barrers. P-RF-C is identified from the RF model, from dataset C. P-RF-D is identified from the RF model, from dataset D. P-DNN-C1 to P-DNN-C4 are identified from the DNN ensemble model, from dataset C. P-DNN-D1 and P-DNN-D2 are identified from the DNN ensemble model, from dataset D.

and large aliphatic rings, which is a relatively underexplored strategy in membrane design. At the same time, these polymers attain unprecedented O<sub>2</sub>/N<sub>2</sub> selectivity via the incorporation of polar groups, such as carbonyls and sulfonyls. Alternatively, P-DNN-C1 and P-DNN-C2 each feature an inflexible polycyclic backbone and two trifluoromethyl-containing side chains. Amazingly, they demonstrate

that our DNN model learns the importance of bulky spherical groups (such as trifluoromethyl groups) for creating steric frustration (61), which has been recognized in the gas separation community as favoring higher gas permeability (17, 62, 63). Restricted backbone mobility plus the presence of the bulky pendant groups disrupt polymer chain packing and lead to high fractional free volume and

ultrahigh permeability, while the polar polyimide backbone helps to maintain selectivity.

Differently, the discovered polymers from dataset D use a rigid ladder-type backbone with a spirobifluorene (SBF) unit, like many other ladder-type PIMs. The fused benzene rings in the SBF unit reduce the flexibility of the backbone around the spirocenter, and the two-bond ladder connections restrict the rotational ability of the backbone. The reduced chain flexibility may also prohibit chain motion to help resist physical aging (17, 64). To further increase permeability, P-DNN-D1 and P-DNN-D2 attach a fused tetramethyltetrahydronaphthalene (TMN) to the SBF unit, which incorporates additional aliphatic rings and methyl substituents (21).

Overall, the generalizable ML models presented here are capable of efficiently finding promising polymers with high performance, with thousands of candidates lying beyond the 2008 Robeson upper bound (7). In addition, the ultrahigh permeability polymers found in this work would allow for never-before-seen industrial gas separations with higher throughput while maintaining sufficient selectivity. Incredibly, the DNN model can extrapolate relatively accurately to high permeability predictions that it did not see in training. We believe that this amazing performance primarily arises from careful selection of diverse training samples and training with a neural network that can capture not only complexities but also generalizations through multitask parameter sharing and ensembling.

Our experimentally validated MD simulations of gas permeability confirm the ML predictions, which suggests that many of the polymer candidates found here can be translated to reality in experiments. Each of the promising polyimides identified here have a well-defined cross-linking formation from existing PubChem chemicals, which makes their syntheses feasible. However, the difficulty of synthesizing complex polymers in a solution-processable manner should not be underestimated. Therefore, to facilitate the overcoming of this challenge, we have tabulated the thousands of promising polymers that we have identified and included them in the GitHub repository associated with this work (<https://github.com/jsunn-y/PolymerGasMembraneML>), which we encourage experimental and computational researchers to explore further. While our models consider membrane performance to be constant, future efforts should also take into account how aging, plasticization, and swelling can degrade membrane performance over time, which is an important consideration in membrane design (5).

Ultimately, we provide the membrane design community with many previously unknown high-performance polymer candidates and key chemical features to consider when designing their molecular structures. Lessons from the workflow demonstrated in this study can likely serve as a guide for other materials discovery and design tasks, such as polymer membranes for desalination and water treatment (10), high-temperature fuel cells, and catalysis. With the continual improvement of ML techniques and an increase in computing power, we expect that ML-assisted design frameworks will only gain popularity and deliver increasingly substantial results in materials discovery for a wide range of applications.

## METHODS

### Calculation of chemical representations

The workflow for our ML method to learn synthesis-property relationships of gas separation membranes is shown in Fig. 1. In step 1, the training set consists of the single repeating units of 353 unique

homopolymers with at least one known gas permeability (among He, H<sub>2</sub>, O<sub>2</sub>, N<sub>2</sub>, CO<sub>2</sub>, and CH<sub>4</sub>), as obtained from the PoLyInfo and MSA databases. In the datasets, each polymer entry is identified on the basis of its unique SMILES string. ML models for prediction of gas permeability from chemistry must use a descriptive and appropriate input to represent the polymer (33). Thus, in step 2, RDKit (65) is used to calculate two different chemical representations of each polymer's repeating unit: its molecular descriptors and its MFF (41).

First, 146 relevant chemical descriptors are calculated, which generally includes information such as number of certain atom types, presence/absence of features, and number of rings, among other physical descriptors that can be calculated from the atom types and connectivity. A list of the available descriptors in RDKit is provided in table S1. Thus, the important chemical features of each polymer repeating unit are identified. We additionally use RDKit to generate the MFF for each repeating unit chemistry. In short, the fingerprinting process consists of the following (31): (i) assign each atom with an identifier, (ii) update each atom's identifiers based on its neighbors, (iii) remove duplicates, and (iv) fold list of identifiers into a bit vector (a Morgan fingerprint). In our case, the chemical substructures considered are up to three units in radius, where each atom or bond is one unit, resulting in 3209 different substructures. In the fingerprint vector of length 3209, each bucket indicates if a certain substructure is present, and we minimize information loss by accounting for frequency if a substructure is present multiple times in a single repeating unit, known as the MFF (31). Last, we shorten the fingerprint vector by only using the 114 most frequently occurring substructures in the training set as input features. Unlike group contribution methods, fingerprinting is dynamic and can evolve to include new chemical structures and connectivities between neighboring repeating units.

### Imputation of missing permeabilities

The calculated molecular features are then used as inputs for multitask ML models and trained to learn gas permeabilities (He, H<sub>2</sub>, O<sub>2</sub>, N<sub>2</sub>, CO<sub>2</sub>, and CH<sub>4</sub>), in step 3. For each of our supervised ML models, training is based on the log of the permeability measured in Barrers and, for a given polymer, the permeability values are averaged across multiple literature sources, if available. Many polymer entries in our training database have missing data, where gas permeabilities are not available for all six gases under study. Yuan *et al.* (49) have demonstrated effective imputation of missing gas permeability data using the MICE algorithm, if at least one gas permeability is available. We use their source code to impute missing gas permeabilities to augment our dataset. For MICE, we compare a linear predictive model, BLR, with a nonlinear predictive model, ERTs. By filling in missing gas permeabilities, imputation allows us to train multitask ML models. This improves our models via parameter sharing, which is physically reasonable, as permeabilities between different gases are correlated for a given membrane chemistry.

### Training of supervised ML models

We train multitask RF and DNN models to predict gas permeabilities based on chemical descriptors and fingerprints, with 20% of the data reserved for the test set and the remaining 80% used for training, selected randomly. RFs reduce variance in decision trees by making regression predictions based on the average of many decision trees: During the growth of each tree, each new decision

rule is made using only a random subset of data points and features. RFs thus build generalizable models of nonlinear relationships. We train each RF using 200 estimators, with training capped at the square root of the number of features for each decision tree and a max tree depth of 10.

In addition, we train DNNs. These DNNs have five hidden layers with 64, 64, 32, 16, and 8 nodes, respectively; ReLU activation; and dropout of 0.1. The multitask models output six permeabilities for the six gases of study. The DNNs are trained using minibatch gradient descent with a batch size of 64, the Adam optimizer, and mean squared error loss. Once the ML models are trained and achieve good performance, we then screen over 9 million hypothetical polymers (summarized in Table 1) to predict their gas permeabilities, in step 4. Our screening predictions are then used to identify promising polymer candidates with high permeability and selectivity. The code and datasets for our ML implementation can be found at <https://github.com/jsunn-y/PolymerGasMembraneML>.

### Ensembling

Because of their density and complexity, deep learning models can be susceptible to particularly high variance. Especially, when trained on a small dataset, there is inherent stochasticity resulting from the network's initialization and the order of data processing during training. There exist many ways to quantify and reduce uncertainty for problems with chemical inputs (56). One simple way to improve the predictive capacity of these models is through ensembling, or averaging together several models trained under different conditions. For example, given an ensemble of distinct models  $\mathcal{E} = \{M_1, M_2, \dots, M_n\}$  and inputs  $x$ , the ensemble prediction is given by the mean of all the model predictions

$$\bar{M}(x) = \sum_{M \in \mathcal{E}} \frac{M(x)}{n}$$

The uncertainty of the prediction can then be measured as the variance between model outputs.

$$U(x) = \sum_{M \in \mathcal{E}} \frac{(\bar{M}(x) - M(x))^2}{n}$$

While there are many ways to perform ensembling, we choose to use bootstrapping, or training each model in the ensemble with a different random subset of the training data. First, we randomly select 20% of the data to be the holdout set, which is used for performance scoring. Sixteen independent models are trained, using 80% of the entries in the non-holdout set each time, selected at random. The training curves of our DNN ensemble on MFFs with BLR imputation of permeabilities are given in fig. S14.

### Feature importance analysis using SHAP

Alongside the models trained in step 3 of our workflow, we can perform explainable ML. To strengthen our physical understanding of how chemical features are linked to performance in gas separation membranes, our primary tool involves assessing SHAP values from each model (42). In essence, the SHAP approach considers how well a model performs when each feature is neglected during training. By analyzing the quantitative impact of leaving out a feature on the model prediction, a feature importance can be assigned. Moreover, each sample's impact on the final model prediction can also be evaluated.

### In silico atomistic models of membranes

In step 5, all-atom MD simulations, using large-scale atomic/molecular massively parallel simulator (LAMMPS) (66), are performed. The polymer consistent force field (PCFF) (67) is used to describe the interatomic interactions of both polymer and gas, which has been widely used to calculate the mechanical properties, cohesive energies, heat capacities, and elastic constants of organic polymers.

We construct the polymeric membrane models via the multistep cross-linking of binary components in the polymers of interest, as most high-performance polymers are polyimides or ladder polymers. The reactive atoms are first assigned to each monomer, and 45 of each component are packed into a three-dimensional (3D)-periodic amorphous cell. Geometry optimization and five annealing cycles of the packed system are carried out. The optimized structure is then cross-linked under the constant temperature, constant volume (NVT) ensemble within an initial cutoff distance of 4.5 Å. Covalent bonds are formed between reactive atoms, and the cross-linked network is relaxed under the constant temperature, constant pressure (NPT) ensemble for 1 ns. After that, the next cross-linking step continues with an increased cutoff distance of 0.5 Å until the cross-linking degree reaches 90%. During the cross-linking process, extra hydrogen atoms are removed, and partial charges are updated to follow assignments from the force field and charge neutrality. The generated, cross-linked polymer structure is used for subsequent calculations. Details of the cross-linking results for selected ladder polymers and polyimides are presented in figs. S15 to S19.

### MD simulations for permeability validation

To validate the gas permeabilities of selected polymeric membranes, each gas's permeability is calculated as the product of its solubility and diffusivity (24). For solubility calculations, MD simulations are performed with a time step of 1.0 fs in the following sequence: (i) energy minimization, (ii) 0.5-ns NVT-MD simulation at 600 K, (iii) 0.5-ns NPT-MD simulation at 600 K and 1 bar, (iv) five thermal annealing cycles from 600 to 300 K with a temperature interval of 50 K at 1 bar, (v) 0.5-ns NPT-MD simulation at 300 K and 1 bar, and (vi) 0.5-ns NVT-MD simulation at 300 K. Last, the solubility coefficients of relevant gases are evaluated at infinite dilution, which are equal to their Henry's constants (68).

Before simulations of diffusivity, gas molecules, such as H<sub>2</sub>, CH<sub>4</sub>, CO<sub>2</sub>, O<sub>2</sub>, or N<sub>2</sub>, are inserted into the simulation box of the cross-linked polymer. The system is first equilibrated through a 21-step MD equilibration protocol (68). The system is then equilibrated for 1 ns under the NVT ensemble at 300 K, followed by 2 ns under the NPT ensemble at 300 K and 1 atm. Production runs are then performed for a duration of 7 ns. The first 2 ns are used for equilibration, and the remaining 5 ns for analysis. The diffusion coefficient of gas molecules in the cross-linked polymer is estimated by the mean squared displacement (MSD) defined as

$$\text{MSD}(t) = \frac{1}{6N} \frac{d}{dt} \lim_{t \rightarrow \infty} \sum_{i=0}^N \langle |r_i(t) - r_i(0)|^2 \rangle$$

where  $N$  is the number of gas molecules and  $r_i(t)$  is the position of molecule  $i$  at time  $t$ . MSD is calculated from the ensemble average  $\langle \dots \rangle$  of the trajectory, and we use the multiple-origin method to improve the statistical accuracy. In addition, to account for molecular adsorption to the polymer membrane at the saturation state, we consider the diffusivity of different numbers of gas molecules: 5, 10, 20, 30, 40, 50, and 100. We find that using 20, 30, or 40 gas

molecules result in similar diffusivities, which are averaged to give the calculated diffusivity. Last, on the basis of the solution-diffusion mechanism, gas permeability ( $P_i$ ) in a polymer membrane can be expressed as the product of the diffusivity ( $D_i$ ) and the solubility ( $S_i$ ). As summarized in table S14, our benchmark study on the solubility, diffusivity, and permeability of five pertinent gases ( $H_2$ ,  $N_2$ ,  $O_2$ ,  $CO_2$ , and  $CH_4$ ) in a PIM-1 membrane agrees well with available experimental data and simulation results (23, 46). Similarly, the permeabilities calculated using the method in this study for 10 other relevant polyimides and ladder polymers match with experimental measurements in the literature (fig. S12 and table S15).

## SUPPLEMENTARY MATERIALS

Supplementary material for this article is available at <https://science.org/doi/10.1126/sciadv.abp8823>

## REFERENCES AND NOTES

1. S. Basu, A. L. Khan, A. Cano-Odena, C. Liu, I. F. J. Vankelecom, Membrane-based technologies for biogas separations. *Chem. Soc. Rev.* **39**, 750–768 (2010).
2. S. Zhao, P. H. M. Feron, L. Deng, E. Favre, E. Chabanon, S. Yan, J. Hou, V. Chen, H. Qi, Status and progress of membrane contactors in post-combustion carbon capture: A state-of-the-art review of new developments. *J. Membr. Sci.* **511**, 180–206 (2016).
3. Y. Han, W. S. W. Ho, Polymeric membranes for  $CO_2$  separation and capture. *J. Membr. Sci.* **628**, 119244 (2021).
4. B. D. Freeman, Basis of permeability/selectivity tradeoff relations in polymeric gas separation membranes. *Macromolecules* **32**, 375–380 (1999).
5. D. F. Sanders, Z. P. Smith, R. Guo, L. M. Robeson, J. E. McGrath, D. R. Paul, B. D. Freeman, Energy-efficient polymeric gas separation membranes for a sustainable future: A review. *Polymer* **54**, 4729–4761 (2013).
6. H. B. Park, J. Kamcev, L. M. Robeson, M. Elimelech, B. D. Freeman, Maximizing the right stuff: The trade-off between membrane permeability and selectivity. *Science* **356**, eaab0530 (2017).
7. L. M. Robeson, The upper bound revisited. *J. Membr. Sci.* **320**, 390–400 (2008).
8. B. Comesaña-Gándara, J. Chen, C. G. Bezzu, M. Carta, I. Rose, M.-C. Ferrari, E. Esposito, A. Fuoco, J. C. Jansen, N. B. McKeown, Redefining the Robeson upper bounds for  $CO_2/CH_4$  and  $CO_2/N_2$  separations using a series of ultrapermeable benzotriptycene-based polymers of intrinsic microporosity. *Energ. Environ. Sci.* **12**, 2733–2740 (2019).
9. R. Swaidan, B. Ghanem, I. Pinnau, Fine-tuned intrinsically ultramicroporous polymers redefine the permeability/selectivity upper bounds of membrane-based air and hydrogen separations. *ACS Macro Lett.* **4**, 947–951 (2015).
10. J. R. Werber, C. O. Osuji, M. Elimelech, Materials for next-generation desalination and water purification membranes. *Nat. Rev. Mater.* **1**, 1–15 (2016).
11. S. Wang, X. Li, H. Wu, Z. Tian, Q. Xin, G. He, D. Peng, S. Chen, Y. Yin, Z. Jiang, M. D. Guiver, Advances in high permeability polymer-based membrane materials for  $CO_2$  separations. *Energ. Environ. Sci.* **9**, 1863–1890 (2016).
12. T. J. Corrado, Z. Huang, D. Huang, N. Wamble, T. Luo, R. Guo, Pentiptycene-based ladder polymers with configurational free volume for enhanced gas separation performance and physical aging resistance. *Proc. Natl. Acad. Sci. U.S.A.* **118**, e2022204118 (2021).
13. H. Sanaeepur, A. Ebadi Amooghini, S. Bandehali, A. Moghadassi, T. Matsuura, B. V. der Bruggen, Polyimides in membrane gas separation: Monomer's molecular design and structural engineering. *Prog. Polym. Sci.* **91**, 80–125 (2019).
14. J. Wang, Z. Shi, Y. Zang, H. Jia, M. Teraguchi, T. Aoki, Macromolecular design for oxygen/nitrogen permselective membranes—Top-performing polymers in 2020. *Polymers* **13**, 3012 (2021).
15. J. Liu, X. Hou, H. B. Park, H. Lin, High-performance polymers for membrane  $CO_2/N_2$  separation. *Chem. A Eur. J.* **22**, 15980–15990 (2016).
16. N. B. McKeown, Polymers of intrinsic microporosity (PIMs). *Polymer* **202**, 122736 (2020).
17. T. Corrado, R. Guo, Macromolecular design strategies toward tailoring free volume in glassy polymers for high performance gas separation membranes. *Mol. Syst. Des. Eng.* **5**, 22–48 (2020).
18. M. F. Jimenez-Solomon, Q. Song, K. E. Jelfs, M. Munoz-Ibanez, A. G. Livingston, Polymer nanofilms with enhanced microporosity by interfacial polymerization. *Nat. Mater.* **15**, 760–767 (2016).
19. B. S. Ghanem, R. Swaidan, E. Litwiller, I. Pinnau, Ultra-microporous triptycene-based polyimide membranes for high-performance gas separation. *Adv. Mater.* **26**, 3688–3692 (2014).
20. B. S. Ghanem, R. Swaidan, X. Ma, E. Litwiller, I. Pinnau, Energy-efficient hydrogen separation by AB-type ladder-polymer molecular sieves. *Adv. Mater.* **26**, 6696–6700 (2014).
21. I. Rose, C. G. Bezzu, M. Carta, B. Comesaña-Gándara, E. Lasseguette, M. C. Ferrari, P. Bernardo, G. Clarizia, A. Fuoco, J. C. Jansen, K. E. Hart, T. P. Liyana-Arachchi, C. M. Colina, N. B. McKeown, Polymer ultrapermeability from the inefficient packing of 2D chains. *Nat. Mater.* **16**, 932–937 (2017).
22. R. C. Dutta, S. K. Bhatia, Atomistic investigation of mixed-gas separation in a fluorinated polyimide membrane. *ACS Appl. Polym. Mater.* **1**, 1359–1371 (2019).
23. W. Fang, L. Zhang, J. Jiang, Polymers of intrinsic microporosity for gas permeation: A molecular simulation study. *Mol. Simul.* **36**, 992–1003 (2010).
24. R. M. Venable, A. Krämer, R. W. Pastor, Molecular dynamics simulations of membrane permeability. *Chem. Rev.* **119**, 5954–5997 (2019).
25. S. Yi, B. Ghanem, Y. Liu, I. Pinnau, W. J. Koros, Ultraselective glassy polymer membranes with unprecedented performance for energy-efficient sour gas separation. *Sci. Adv.* **5**, eaaw5459 (2019).
26. L. M. Robeson, C. D. Smith, M. Langsam, A group contribution approach to predict permeability and permselectivity of aromatic polymers. *J. Membr. Sci.* **132**, 33–54 (1997).
27. E. R. Hensema, M. H. V. Mulder, C. A. Smolders, C. A. Smolders, On the mechanism of gas transport in rigid polymer membranes. *J. Appl. Polym. Sci.* **49**, 2081–2090 (1993).
28. M. H. Cohen, D. Turnbull, Molecular transport in liquids and glasses. *J. Chem. Phys.* **31**, 1164–1169 (1959).
29. G. Chen, Z. Shen, A. Iyer, U. F. Ghuman, S. Tang, J. Bi, W. Chen, Y. Li, Machine-learning-assisted de novo design of organic molecules and polymers: Opportunities and challenges. *Polymers* **12**, 163 (2020).
30. D. J. Audus, J. J. de Pablo, Polymer informatics: Opportunities and challenges. *ACS Macro Lett.* **6**, 1078–1082 (2017).
31. L. Tao, G. Chen, Y. Li, Machine learning discovery of high-temperature polymers. *Patterns* **2**, 100225 (2021).
32. G. Chen, L. Tao, Y. Li, Predicting polymers' glass transition temperature by a chemical language processing model. *Polymers* **13**, 1898 (2021).
33. L. Tao, V. Varshney, Y. Li, Benchmarking machine learning models for polymer informatics: An example of glass transition temperature. *J. Chem. Inf. Model.* **61**, 5395–5413 (2021).
34. S. Wu, Y. Kondo, M. Kakimoto, B. Yang, H. Yamada, I. Kuwajima, G. Lambard, K. Hongo, Y. Xu, J. Shiomi, C. Schick, J. Morikawa, R. Yoshida, Machine-learning-assisted discovery of polymers with high thermal conductivity using a molecular design algorithm. *Npj Comput. Mater.* **5**, 1–11 (2019).
35. A. Mannodi-Kanakithodi, G. Pilania, T. D. Huan, T. Lookman, R. Ramprasad, Machine learning strategy for accelerated design of polymer dielectrics. *Sci. Rep.* **6**, 20952 (2016).
36. W. Sun, Y. Zheng, K. Yang, Q. Zhang, A. A. Shah, Z. Wu, Y. Sun, L. Feng, D. Chen, Z. Xiao, S. Lu, Y. Li, K. Sun, Machine learning-assisted molecular design and efficiency prediction for high-performance organic photovoltaic materials. *Sci. Adv.* **5**, eaay4275 (2019).
37. R. Gómez-Bombarelli, J. Aguilera-Iparraguirre, T. D. Hirzel, D. Duvenaud, D. Maclaurin, M. A. Blood-Forsythe, H. S. Chae, M. Einzinger, D.-G. Ha, T. Wu, G. Markopoulos, S. Jeon, H. Kang, H. Miyazaki, M. Numata, S. Kim, W. Huang, S. I. Hong, M. Baldo, R. P. Adams, A. Aspuru-Guzik, Design of efficient molecular organic light-emitting diodes by a high-throughput virtual screening and experimental approach. *Nat. Mater.* **15**, 1120–1127 (2016).
38. C. L. Ritt, M. Liu, T. A. Pham, R. Epsztajn, H. J. Kulik, M. Elimelech, Machine learning reveals key ion selectivity mechanisms in polymeric membranes with subnanometer pores. *Sci. Adv.* **8**, eabl5771 (2022).
39. J. W. Barnett, C. R. Bilchak, Y. Wang, B. C. Benicewicz, L. A. Murdock, T. Bureau, S. K. Kumar, Designing exceptional gas-separation polymer membranes using machine learning. *Sci. Adv.* **6**, eaaz4301 (2020).
40. T. Liu, L. Liu, F. Cui, F. Ding, Q. Zhang, Y. Li, Predicting the performance of polyvinylidene fluoride, polyethersulfone and polysulfone filtration membranes using machine learning. *J. Mater. Chem. A* **8**, 21862–21871 (2020).
41. D. Rogers, M. Hahn, Extended-connectivity fingerprints. *J. Chem. Inf. Model.* **50**, 742–754 (2010).
42. S. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions. arXiv:1705.07874 [cs.LG] (22 May 2017).
43. R. Ma, T. Luo, P11M: A benchmark database for polymer informatics. *J. Chem. Inf. Model.* **60**, 4684–4690 (2020).
44. A. Ghosh, S. K. Sen, S. Banerjee, B. Voit, Solubility improvements in aromatic polyimides by macromolecular engineering. *RSC Adv.* **2**, 5900–5926 (2012).
45. S. Kim, J. Chen, T. Cheng, A. Gindulyte, J. He, S. He, Q. Li, B. A. Shoemaker, P. A. Thiessen, B. Yu, L. Zaslavsky, J. Zhang, E. E. Bolton, PubChem 2019 update: Improved access to chemical data. *Nucleic Acids Res.* **47**, D1102–D1109 (2019).
46. P. M. Budd, K. J. Msayib, C. E. Tattershall, B. S. Ghanem, K. J. Reynolds, N. B. McKeown, D. Fritsch, Gas separation membranes from polymers of intrinsic microporosity. *J. Membr. Sci.* **251**, 263–269 (2005).
47. N. Du, M. D. Guiver, G. P. Robertson, Ladder polymers with intrinsic microporosity and process for production thereof. U.S. Patent 9,371,429 (2016).
48. L. McInnes, J. Healy, J. Melville, UMAP: Uniform manifold approximation and projection for dimension reduction. arXiv:1802.03426 [stat.ML] (9 February 2018).

49. Q. Yuan, M. Longo, A. W. Thornton, N. B. McKeown, B. Comesaña-Gándara, J. C. Jansen, K. E. Jeffs, Imputation of missing gas permeability data for polymer membranes using machine learning. *J. Membr. Sci.* **627**, 119207 (2021).
50. T. H. Kim, W. J. Koros, G. R. Husk, K. C. O'Brien, Relationship between gas separation properties and chemical structure in a series of aromatic polyimides. *J. Membr. Sci.* **37**, 45–62 (1988).
51. Y. Hu, M. Shiotsuki, F. Sanda, B. D. Freeman, T. Masuda, Synthesis and properties of indan-based polyacetylenes that feature the highest gas permeability among all the existing polymers. *Macromolecules* **41**, 8525–8532 (2008).
52. H. Lin, B. D. Freeman, Gas permeation and diffusion in cross-linked poly(ethylene glycol diacrylate). *Macromolecules* **39**, 3568–3580 (2006).
53. P. Ertl, A. Schuffenhauer, Estimation of synthetic accessibility score of drug-like molecules based on molecular complexity and fragment contributions. *J. Chem. I.* **8** (2009).
54. J. G. Wijmans, R. W. Baker, The solution-diffusion model: A review. *J. Membr. Sci.* **107**, 1–21 (1995).
55. M. Pardakhti, E. Moharrer, D. Wanik, S. L. Suib, R. Srivastava, Machine learning using combined structural and chemical descriptors for prediction of methane adsorption performance of metal organic frameworks (MOFs). *ACS Comb. Sci.* **19**, 640–645 (2017).
56. L. Hirschfeld, K. Swanson, K. Yang, R. Barzilay, C. W. Coley, Uncertainty quantification using neural networks for molecular property prediction. *J. Chem. Inf. Model.* **60**, 3770–3780 (2020).
57. K. Yang, K. Swanson, W. Jin, C. Coley, P. Eiden, H. Gao, A. Guzman-perez, T. Hopper, B. Kelley, M. Mathea, A. Palmer, V. Settels, T. Jaakkola, K. Jensen, R. Barzilay, Analyzing learned molecular representations for property prediction. *J. Chem. Inf. Model.* **59**, 3370–3388 (2019).
58. C. Rudin, Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nat. Mach. Intell.* **1**, 206–215 (2019).
59. S. M. Lundberg, G. Erion, H. Chen, A. DeGrave, J. M. Prutkin, B. Nair, R. Katz, J. Himmelfarb, N. Bansal, S.-I. Lee, From local explanations to global understanding with explainable AI for trees. *Nat. Mach. Intell.* **2**, 56–67 (2020).
60. J. Jiménez-Luna, M. Skalic, N. Weskamp, G. Schneider, Coloring molecules with explainable artificial intelligence for preclinical relevance assessment. *J. Chem. Inf. Model.* **61**, 1083–1094 (2021).
61. N. Du, G. P. Robertson, J. Song, I. Pinnau, S. Thomas, M. D. Guiver, Polymers of intrinsic microporosity containing trifluoromethyl and phenylsulfone groups as materials for membrane gas separation. *Macromolecules* **41**, 9656–9662 (2008).
62. Y. He, F. M. Benedetti, S. Lin, C. Liu, Y. Zhao, H.-Z. Ye, T. Van Voorhis, M. G. De Angelis, T. M. Swager, Z. P. Smith, Polymers with side chain porosity for ultrapermeable and plasticization resistant materials for gas separations. *Adv. Mater.* **31**, 1807871 (2019).
63. Y. Zhao, Y. He, T. M. Swager, Porous organic polymers via ring opening metathesis polymerization. *ACS Macro Lett.* **7**, 300–304 (2018).
64. R. Swaidan, B. Ghanem, E. Litwiller, I. Pinnau, Physical aging, plasticization and their effects on gas permeation in "rigid" polymers of intrinsic microporosity. *Macromolecules* **48**, 6553–6561 (2015).
65. G. Landrum, *RDKit: A Software Suite for Cheminformatics, Computational Chemistry, and Predictive Modeling* (Academic Press, 2013).
66. S. Plimpton, Fast parallel algorithms for short-range molecular dynamics. *J. Comput. Phys.* **117**, 1–19 (1995).
67. H. Sun, S. J. Mumby, J. R. Maple, A. T. Hagler, An ab initio CFF93 all-atom force field for polycarbonates. *J. Am. Chem. Soc.* **116**, 2978–2987 (1994).
68. V. Varshney, S. S. Patnaik, A. K. Roy, B. L. Farmer, A molecular dynamics study of epoxy-based networks: Cross-linking procedure and prediction of molecular and material properties. *Macromolecules* **41**, 6837–6842 (2008).
69. L. H. Hall, B. Mohney, L. B. Kier, The electrotopological state: Structure information at the atomic level for molecular graphs. *J. Chem. Inf. Comput. Sci.* **31**, 76–82 (1991).
70. M. C. Sorkun, J. M. V. A. Koelman, S. Er, Pushing the limits of solubility prediction via quality-oriented data selection. *iScience* **24**, 101961 (2021).
71. P. G. Francoeur, D. R. Koes, SolTranNet—A machine learning tool for fast aqueous solubility prediction. *J. Chem. Inf. Model.* **61**, 2530–2536 (2021).
72. A. Chandrasekaran, C. Kim, S. Venkatram, R. Ramprasad, A deep learning solvent-selection paradigm powered by a massive solvent/nonsolvent database for polymers. *Macromolecules* **53**, 4764–4769 (2020).
73. H. Sun, COMPASS: An ab initio force-field optimized for condensed-phase applications overview with details on alkane and benzene compounds. *J. Phys. Chem. B* **102**, 7338–7364 (1998).
74. S.-H. Park, K.-J. Kim, W.-W. So, S.-J. Moon, S.-B. Lee, Gas separation properties of 6FDA-based polyimide membranes with a polar group. *Macromol. Res.* **11**, 157–162 (2003).
75. C. G. Bezzu, M. Carta, M.-C. Ferrari, J. C. Jansen, M. Monteleone, E. Esposito, A. Fuoco, K. Hart, T. P. Liyana-Arachchi, C. M. Colina, N. B. McKeown, The synthesis, chain-packing simulation and long-term gas permeability of highly selective spirofluorene-based polymers of intrinsic microporosity. *J. Mater. Chem. A* **6**, 10507–10514 (2018).
76. B. Satilmis, M. Lanč, A. Fuoco, C. Rizzuto, E. Tocci, P. Bernardo, G. Clarizia, E. Esposito, M. Monteleone, M. Dendisová, K. Friess, P. M. Budd, J. C. Jansen, Temperature and pressure dependence of gas permeation in amine-modified PIM-1. *J. Membr. Sci.* **555**, 483–496 (2018).
77. B. S. Ghanem, N. B. McKeown, P. M. Budd, J. D. Selbie, D. Fritsch, High-performance membranes from polyimides with intrinsic microporosity. *Adv. Mater. Deerfield Beach Fla.* **20**, 2766–2771 (2008).
78. M. Carta, P. Bernardo, G. Clarizia, J. C. Jansen, N. B. McKeown, Gas permeability of hexaphenylbenzene based polymers of intrinsic microporosity. *Macromolecules* **47**, 8320–8327 (2014).
79. J. Wu, S. Japip, T.-S. Chung, Infiltrating molecular gatekeepers with coexisting molecular solubility and 3D-intrinsic porosity into a microporous polymer scaffold for gas separation. *J. Mater. Chem. A* **8**, 6196–6209 (2020).
80. Y. Liu, A. Tang, J. Tan, C. Chen, D. Wu, H. Zhang, Structure and gas barrier properties of polyimide containing a rigid planar fluorene moiety and an amide group: Insights from molecular simulations. *ACS Omega* **6**, 4273–4281 (2021).

**Acknowledgments:** We acknowledge the Texas Advanced Computing Center (TACC) at The University of Texas at Austin (Frontera project and National Science Foundation Award 1818253) and National Renewable Energy Laboratory (Eagle Computing System) for providing HPC resources that have contributed to the research results reported within this paper. J.Y. would like to thank S. Krishnaswamy, M. Amodio, and A. Haji-Akbari for their guidance on efforts related to the project. We would like to thank M. Ostwal for helpful comments on the manuscript. **Funding:** We gratefully acknowledge financial support from the Air Force Office of Scientific Research through the Air Force's Young Investigator Research Program (FA9550-20-1-0183; program manager: M.-J. Pan) and the National Science Foundation (CMMI-1934829 and CAREER Award CMMI-2046751). Y.L. would like to express thanks for the support from 3M's Non-Tenured Faculty Award. Y.L. and J.R.M. would like to thank the support from the National Alliance for Water Innovation (NAWI), funded by the U.S. Department of Energy, Office of Energy Efficiency and Renewable Energy (EEERE), Advanced Manufacturing Office, under Funding Opportunity Announcement Number DE-FOA-0001905. J.Y. was supported by the National Science Foundation Graduate Research Fellowship under fellow ID 2021309491. This research also benefited in part from the computational resources and staff contributions provided by the Booth Engineering Center for Advanced Technology (BECAT) at UConn. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the U.S. Department of Defense. **Author contributions:** Y.L. and L.T. conceived the idea and supervised the research. J.Y. and L.T. collected and analyzed the data and implemented the ML models. J.H. and Y.L. developed and analyzed the molecular simulations. J.Y., L.T., J.R.M., and Y.L. contributed to the design of the project and data analysis. J.Y., L.T., and J.H. wrote the first draft of the manuscript, and all authors contributed to revising the manuscript. **Competing interests:** The authors declare that they have no competing interests. **Data and materials availability:** All data needed to evaluate the conclusions in the paper are present in the paper and/or the Supplementary Materials. The datasets used in this work can be accessed from <http://dx.doi.org/10.22002/D1.20048>.

Submitted 4 January 2022

Accepted 7 June 2022

Published 20 July 2022

10.1126/sciadv.abn9545

## Machine learning enables interpretable discovery of innovative polymers for gas separation membranes

Jason YangLei TaoJinlong HeJeffrey R. McCutcheonYing Li

*Sci. Adv.*, 8 (29), eabn9545. • DOI: 10.1126/sciadv.abn9545

### View the article online

<https://www.science.org/doi/10.1126/sciadv.abn9545>

### Permissions

<https://www.science.org/help/reprints-and-permissions>

Use of this article is subject to the [Terms of service](#)