



# Multi-property prediction and high-throughput screening of polyimides: An application case for interpretable machine learning

Bo Zhang<sup>a</sup>, Xueqing Li<sup>c</sup>, Xinxin Xu<sup>b</sup>, Jingguo Cao<sup>c,\*</sup>, Ming Zeng<sup>a,\*\*</sup>, Wu Zhang<sup>a</sup>

<sup>a</sup> College of Marine and Environmental Sciences, Tianjin University of Science & Technology, 300457, Tianjin, China

<sup>b</sup> Department of Environmental Engineering, College of Environmental & Resource Sciences, Zhejiang University, 310058, Hangzhou, China

<sup>c</sup> College of Chemical Engineering and Materials Science, Tianjin University of Science & Technology, 300457, Tianjin, China

## ARTICLE INFO

### Keywords:

Polyimide  
Machine learning  
Heat resistance  
Mechanical properties  
Dielectricity  
Optical properties

## ABSTRACT

Polyimide (PI), as a high-performance polymer widely used in aerospace, optoelectronics, microelectronics, etc., the properties it focused on in different areas of application were diverse. However, most of the past machine learning studies in polyimide property prediction were focused only on the prediction of a single property. This study focused on four major categories of properties of polyimide, including thermal (T<sub>g</sub>, T<sub>d5</sub>, T<sub>d10</sub>, and CTE), mechanical (T<sub>s</sub> and T<sub>M</sub>), dielectric ( $\epsilon$ ), and optical ( $\lambda_{\text{cutoff}}$ , T<sub>400</sub>, n<sub>av</sub>, and  $\Delta n$ ), totaling eleven properties. PI data synthesized from previous studies were collected. MorGan fingerprints, improved MorGan fingerprints, RDKit and Mordred descriptors were selected as feature representations. Four ML models such as DNN, RF, XGBoost and BT were also built. Collectively, 176 machine learning models were trained for 11 predictions of properties. The performance and generalization ability of the models were confirmed by experimental validation, external validation and leave-one-out cross-validation. SHAP analysis was used to explain the optimal model for each property prediction from a physicochemical point of view and structural aspects, and three PIs with different structures were designed accordingly. Finally, a high-throughput virtual screening of nearly 7.6 million polyimides was performed based on the trained model. SA scores was used to evaluate the ease of synthesis of each PI, and finally high-performance PIs with potential and easy to synthesize were selected for each field. This study could be expected to provide a guideline and a design framework for the application of PIs in various fields in the future.

## 1. Introduction

Polyimide (PI), the top material in industrial plastics, has excellent chemical stability, outstanding heat resistance and remarkable mechanical properties [1,2]. Due to these remarkable thermal and mechanical properties, PI was widely used in aerospace [3], optoelectronics [4,5], microelectronics [6,7] and automotive industries [8]. In commercial production, polyimides can be polymerized by reacting diamines with dianhydrides and removing water molecules [9]. In addition, they can also be polymerized by reacting diamines with diisocyanates and removing carbon dioxide molecules [10]. Over the past decades, researchers have successfully synthesized many polyimides with specific structures by selecting suitable diamines and dianhydrides or diisocyanates [11–13]. Nevertheless, the traditional polyimide research process was highly dependent on extensive experimental

studies and the selection of appropriate diamine, dianhydride or diisocyanate structures. These traditional experimental methods, both in the past and present, have faced the challenges of repetitive experimentation, inefficiency, and unavoidable repeated measurements [14–16]. In order to accelerate material design and effectively predict the performance of novel PIs, researchers proposed molecular simulation and computational methods such as MD simulation, MC simulation and mathematical modeling [17–19]. In the field of MD simulation, there had been many notable research results. For example, to predict structures with high glass transition temperatures for guiding the design of polyimides, Liang et al. [20] used molecular dynamics simulations to predict the T<sub>g</sub> of polyimides. With the help of the Dreding II force field, based on the analysis of the three aromatic PI barriers, they found a correlation between T<sub>g</sub> and the barrier height. Upon comparison, the T<sub>g</sub> values of the three aromatic polyimides simulated by molecular

\* Corresponding author.

\*\* Corresponding author.

E-mail addresses: [cjg@tust.edu.cn](mailto:cjg@tust.edu.cn) (J. Cao), [ming.zeng@tust.edu.cn](mailto:ming.zeng@tust.edu.cn) (M. Zeng).

<https://doi.org/10.1016/j.polymer.2024.127603>

Received 22 June 2024; Received in revised form 27 August 2024; Accepted 5 September 2024

Available online 7 September 2024

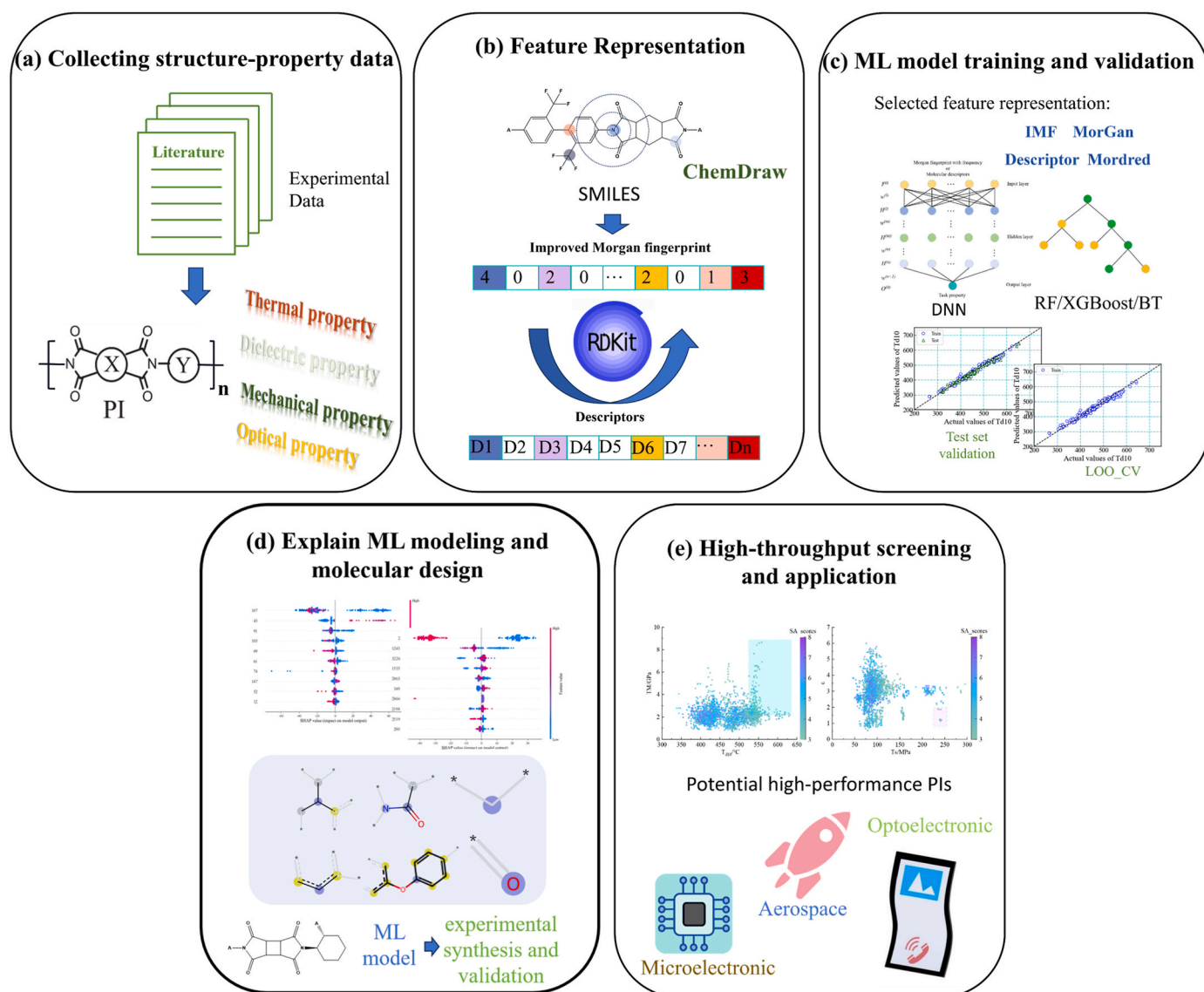
0032-3861/© 2024 Elsevier Ltd. All rights are reserved, including those for text and data mining, AI training, and similar technologies.

dynamics methods were in good agreement with the experimentally tested values. Although this study provided a new insight into PI design, however, the conclusions drawn were based on a small number of polyimides and their generalizability was unknown. H. Lei et al. [21] analyzed the thermodynamic properties of six polyimides by MD simulation, and they found the physicochemical and structural factors affecting the mechanical properties of aromatic polyimides from the low-scale hydrogen bonding interactions, chain conformation, etc. as well as the larger-scale XRD analysis. Compared with the past, their analysis was more in depth and comprehensive, but still suffered from the limitations of time-consuming calculations and limited analytical structures. Currently, molecular simulation and computational methods still face challenges such as time-consuming and expensive all-atom dynamics simulations, limited scope of polymer studies, and difficulties in parameter selection [22–24].

In this context, machine learning (ML) techniques driven by large traditional experimental databases brought a new light to the field of polymer materials design. Compared with traditional experimental synthesis and MD simulations, ML has the advantage of being computationally fast and able to efficiently establish the relationship between

the target polymer structure and its properties [25–27]. Machine learning techniques have been widely used for the prediction of thermal [28–30], mechanical [31], dielectric [32] and optical properties of polymers [33,34]. For example, S.Y. Zhang et al. [35] collected Tg data of 878 polyimides, calculated 208 descriptors using RDKit, and finally obtained 8 key descriptors for training ML models using multiple dimensionality reduction methods. After comparing multiple machine learning models, their optimal model showed excellent performance with a root mean square error of only 11K. On the other hand, Hong Zhang et al. [36] collected data from 652 polyimide molecules, extracted the repeating unit structure to compute the molecular descriptors, and compared seven machine learning algorithms to predict the Tg and cutoff wavelength. The root-mean-square error of the final glass transition temperature prediction model was 33.9 K, and the root-mean-square error of the cutoff wavelength prediction model was 17.1 nm. The excellent research work of countless predecessors has laid a solid foundation for the development of machine learning techniques in the materials domain.

From a general perspective, machine learning (ML)-driven design of polymer materials currently faces the following challenges: 1) most



**Fig. 1.** workflow: (a) collected structure-property data from the literature; (b) represented the structure in SMILES and calculated Fingerprint or Descriptor; (c) trained ML model and validated it; (d) interpreted the trained model and used it to design molecules and validated it; (e) performed high-throughput screening using the trained model to select high-performance polymers for various fields.

studies were mainly focused on the prediction of a single property, or only a few studies considered two properties at the same time. However, for polyimides, which were widely used, there were relatively few studies that could predict multiple properties at the same time; 2) In these few studies that predicted two or more properties, although the number of property predictions increased, the model performance was difficult to guarantee; 3) After training the ML models with excellent performance, it was still a challenge to select the appropriate polymer chemical space and provide guidance for material design in different fields; 4) How to understand the effect of structure on properties from a physicochemical point of view based on the established models?

To meet the above challenges, several types of properties that were the main concern for PI application in aerospace, optoelectronics, and microelectronics were selected as the targets for machine learning prediction in this study. These properties included thermal ( $T_g$ ,  $T_{d5}$ ,  $T_{d10}$ , and CTE), mechanical (Ts and TM), dielectric ( $\epsilon$ ), and optical ( $\lambda_{\text{cutoff}}$ , T400,  $n_{\text{av}}$ , and  $\Delta n$ ) properties, totaling 11. The workflow was shown in Fig. 1. 1018 polyimides were collected from the literature along with the experimental values for each property. Descriptors and molecular fingerprints were chosen as feature representations and four ML models such as DNN, RF, XGBoost and BT were compared. External validation, experimental validation, and leave-one-out cross-validation were used to determine the performance and generalization ability of the models. Interpretable SHAP [37] analysis was used to reveal the structural and physicochemical principles behind the model outputs, and three different PIs were designed accordingly. Finally, a high-throughput virtual screening of nearly 7.6 million polyimides was performed in this study. Meanwhile, the ease of synthesis of each PI was evaluated using SA scores [38], and potential high-performance polyimides with SA scores lower than 4 were selected for the three main areas (aerospace, optoelectronics, and microelectronics). Overall, the models developed in this study had the ability of high performance, multi-property prediction, interpretability, and high-throughput screening, which was expected to facilitate the synthesis of materials in the future when polyimides were applied in different fields.

## 2. Materials and methods

### 2.1. Datasets

In order to construct a structure-property relationship model covering a multitude of properties of polyimides, the structure-property data of 1018 polyimides were extracted from 119 published articles and constituted as dataset I (see Table S4 for the source literature of dataset I). Dataset I contained data on 11 properties of polyimide, including four thermal properties: glass transition temperature ( $T_g$ ), thermal decomposition temperature (5 %) ( $T_{d5}$ ), thermal decomposition temperature (10 %) ( $T_{d10}$ ), and coefficient of thermal expansion (CTE); three physical properties: tensile strength (Ts), tensile modulus (TM), and dielectric constant ( $\epsilon$ ); and four optical properties: transparency (400 nm) (T400), cutoff wavelength ( $\lambda_{\text{cutoff}}$ ), average refractive index ( $n_{\text{av}}$ ) and birefringence ( $\Delta n$ ). Specific data amounts and data ranges for these 11 properties were listed in Table 1. Considering the effects of different testing methods, conditions, and synthesis methods of PIs on the values of each property, when encountering different property values of the same PI,

we all take the mean value as a representative. Each PI in dataset 1 requires number-average molecular weight ( $M_n$ ) > 6000 g/mol and weight-average molecular weight ( $M_w$ ) > 10,000 g/mol. This was done to limit the effect of processing conditions and other factors on the performance of the PI. The distribution of data for each property was shown in Fig. 2.

Inspired by the study of Yang et al. [23], a series of diamine, dianhydride and diisocyanate monomers that had already been successfully synthesized in previous studies were downloaded from the PubChem database. Based on the conventional synthesis methods of polyimides, about 7.6 million hypothetical polyimides were constructed in this study to form data set II. These hypothetical polyimides, although not actually synthesized, provided a broad space of polymer chemistry that could be useful in guiding the design of polyimides applicable to different fields.

### 2.2. Feature representations

**SMILES** The repeating units of a polymer were used to represent its overall molecular structure. The presence of junctions of duplicate units was considered and denoted by “\*”. For each polyimide sample in the dataset, the repetitive units of its molecular structure were drawn in this study using ChemDraw 19.0 software and converted into SMILES strings. In order to avoid different SMILES encoding forms for the same molecular structure, this study used the MolStandardize module of the RDKit package to standardize the converted SMILES characters in a uniform manner.

**Descriptor** A total of two types of descriptors were computed, RDKit and Mordred. The SMILES of each PI molecule were converted into descriptors using RDKit [39], an open-source chemical information Python toolkit that calculated 200 2D molecular descriptors. These descriptors respond to the physicochemical, topological, and charged properties of the polymers as well as the intermolecular forces (Table S13). To simplify the model inputs, all descriptor columns with all zeroes in the calculation were eliminated, resulting in the 141 most valuable descriptors (Table S2). 1825 descriptors, containing 2D and 3D, were computed using Mordred. The number of Mordred descriptors was reduced to 141 using variance-based dimensionality reduction [40] and random forest feature importance scoring methods [35].

**Fingerprint** Based on the SMILES of each PI, the Morgan Fingerprint (MF) [41] was computed by using RDKit. The size of the iteration radius of the Morgan fingerprint was crucial to its effectiveness and needed to be adjusted for different molecule types. An iteration radius that was too small would mainly describe the local structure with poor molecular differentiation, while an iteration radius that was too large might increase the computation time and capture a structure that was more globally focused and missing some local information. After continuous trials, 3 was chosen as the iteration radius for Morgan fingerprints. Meanwhile, learning from the previous experience [42], in order to supplement the number of occurrences of a substructure in the repetition unit, the algorithm was fine-tuned in this study, so that the Morgan fingerprint computed by the algorithm contains information on the frequency of occurrence of substructures. Based on the Improved Morgan Fingerprint (IMF) (Fig. S2), 3331-bit vectors were computed, and similarly, to simplify the model inputs, the IMF was screened for the number of features. This was done by summing the frequency of each

**Table 1**

An introduction to the 11 properties of the data in Dataset I.

Property	$T_g$	$T_{d5}$	$T_{d10}$	CTE	Ts	TM
The number of sample points in the dataset	830(787 unique)	589(561 unique)	353(336 unique)	311(297 unique)	332(318 unique)	296(283 unique)
Property range	70~480°C	210.8~631°C	264~652°C	(-5)~129 ppmK <sup>-1</sup>	26.9~390MPa	1~8.83 GPa
Property	$\epsilon$	T400	$\lambda_{\text{cutoff}}$	$n_{\text{av}}$	$\Delta n$	
The number of sample points in the dataset	192(190 unique)	243(232 unique)	346(335 unique)	157(145 unique)	245(240 unique)	
Property range	2.14~3.94	0~90.1 %	226~447 nm	1.4977~1.7411	0~0.2265	

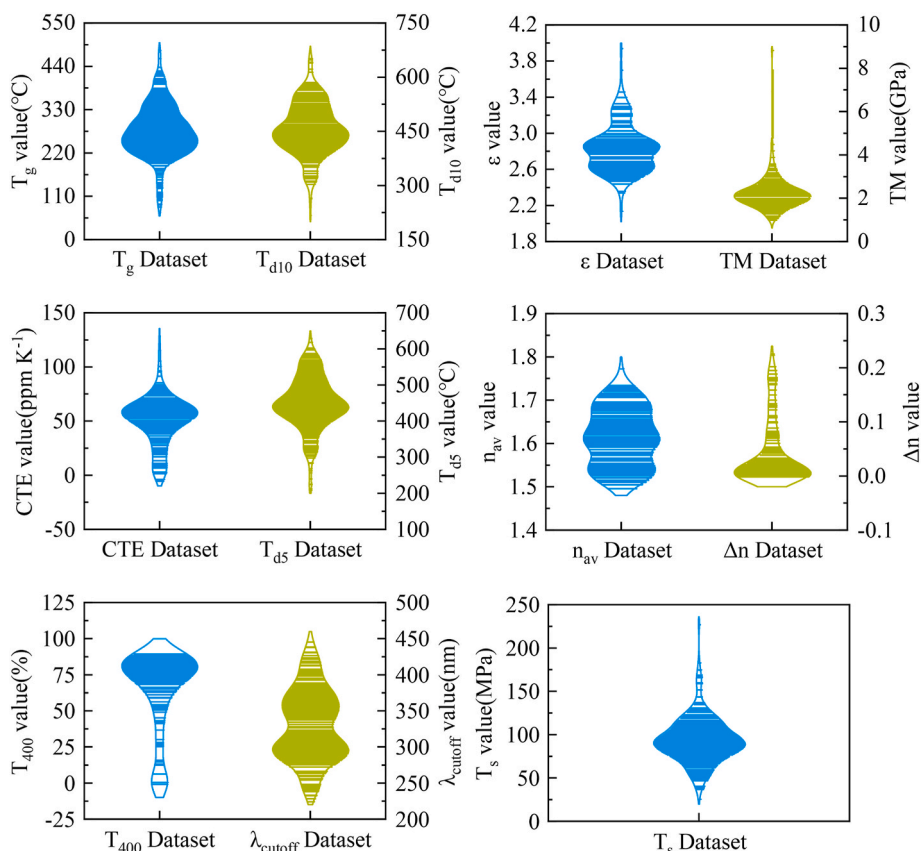


Fig. 2. Distribution of data for the 11 properties in data set I.

column of substructures, qualifying the sum, and finally obtaining the 108 most frequently occurring substructures. The improved MorGan fingerprint minimizes the loss of information in fingerprint extraction, for which the original MorGan fingerprint was similarly created for comparison.

### 2.3. Construction and validation of machine learning model

To model the structure-property relationship of polyimide, four machine learning models, that is, Deep Neural Network (DNN), Random Forest (RF), eXtreme Gradient Boosting (XGBoost), and Boosted tree (BT) model, were employed (Fig. S1). Compared to ANN, DNN has more hidden layers and is more suitable for discovering the complex nonlinear relationships between polymer structures. The DNN model in this study used RELU as the activation function and Adam as the optimizer. Also, three decision tree-based models, i.e., RF, XGBoost and BT models, were also trained which are well suited for multi-input feature samples and are friendly to improved fingerprints and descriptors.

ML models for single-task prediction were built for each of the 11 properties in dataset I. To select the optimal structure-property relationship model for each property, different combinations of input feature representations and model choices were compared. The complete steps of model building were as follows: 1) each polyimide structure was represented as a repeating unit and computed as SMILES; 2) Descriptors and Fingerprints were computed based on SMILES; and 3) DNN, RF, XGBoost and BT models were built based on descriptors or fingerprints, respectively. Thus, for each prediction of PI properties, we compared four ML models and four descriptors or molecular fingerprints. In other words, sixteen models were compared for each property, e.g., Descriptor-DNN-Tg, IMF-XGBoost-Tg, Mordred-BT-Tg, and MorGan-RF-Tg. For the 11 properties, a total of 176 single-task ML models were built and trained.

During model training, the division ratio between the training set and the test set was regulated as a hyperparameter, and two types of divisions, 7:3 and 8:2, were considered. In order to capture as much of the complex nonlinear relationship between polymer structure and performance as possible, but without overfitting the model, our DNN models all contain 4 hidden layers and a dropout of 0.1, which was informed by previous tuning experience with a large number of similar tasks. The number of neuron nodes within each hidden layer of the DNN model was adjusted using manual parameter tuning. The hyperparameters of the remaining models were tuned by a combination of Bayesian optimization and manual parameter tuning (The parameter information of the optimal models for the 11 properties was shown in Table S1).

To prevent model overfitting, the optimal model for each property was subjected to external validation, leave-one-out cross validation (LOO\_CV), and experimental validation. External validation referred to the validation of the model using the test set, which could be used to judge the performance of the model since this part of the data was not involved in the training of the model. Leave-one-out cross-validation was a type of cross-validation method, which involved extracting N-1 data from a sample containing N data to put into model training, leaving one data as validation, and taking one data from the sample as validation each time, running it for N times, and finally taking the mean of the N validation errors. Finally, three resins were prepared for further validation of the model performance. The coefficient of determination ( $R^2$ ), root mean square error (RMSE), the coefficient of determination of leave-one-out cross-validation ( $Q^2_{LOO\_CV}$ ) and the mean root mean square error (Mean\_RMSE) were used to evaluate the performance and generalization ability of the model.

### 3. Results and discussion

#### 3.1. Performance evaluation and validation of the ML models

The objective of this study was to establish a stable and reliable structure-property relationship model for 11 properties of polyimides, and at the same time to utilize this ML model in conjunction with interpretable SHAP analyses to guide the molecular structure design and screening of high-performance polymers for various fields. To this end, the aerospace, optoelectronics and microelectronics fields were taken as examples, and six properties commonly concerned in these three fields were highlighted and analyzed, including: thermal decomposition temperature (10 %) ( $T_{d10}$ ), tensile modulus (TM), transparency (400 nm) (T400), coefficient of thermal expansion (CTE), dielectric constant ( $\epsilon$ ), and tensile strength (Ts). The remaining five properties of the model were also discussed in full, as shown in the Appendix.

For each property, sixteen models were built and compared from the perspective of feature representation and model selection, and finally, the optimal model for a single property was selected from them. Fig. 3 showed the performance of the optimal models for the six properties in focus, such as  $T_{d10}$ , TM, T400, CTE,  $\epsilon$ , and Ts. The results showed that all models had  $R^2$  above 0.9 (the performance graphs of the optimal models for the rest of the properties were shown in Figs. S6–S9). For  $T_{d10}$ , the optimal model was the RF model based on descriptors, and the  $R^2$  of the training set and test set were 0.979 and 0.980, respectively; the RMSE of the training set and test set were 9.771 and 10.057, respectively (Table 2); the average error of about 9.9 °C could be approximated to the error caused by the differences in polyimide synthesis conditions and methods, which demonstrated the excellent performance of the model. For TM, its optimal model was the DNN model built based on IMF, and its  $R^2$  for the training set and test set were 0.948 and 0.940, respectively; and the RMSE for the training set and test set were 0.168 and 0.119, respectively. The rest of them, the optimal models for the properties of T400, CTE,  $\epsilon$ , and Ts had the  $R^2$  for the training set/test set of 0.993/0.991, 0.988/0.987, 0.953/0.924 and 0.952/0.929, respectively; and the RMSEs of the training set/test set were 2.253/2.288, 2.496/2.272, 0.058/0.060 and 6.938/6.584, respectively. Table 2 shows the performance of the optimal ML model for 11 properties. It could be concluded that among the optimal models built for predicting the 11 properties of polyimide, the  $R^2$  of all the models for both the training and test sets were above 0.9, which fully demonstrated the excellent performance and good generalization ability of the built models.

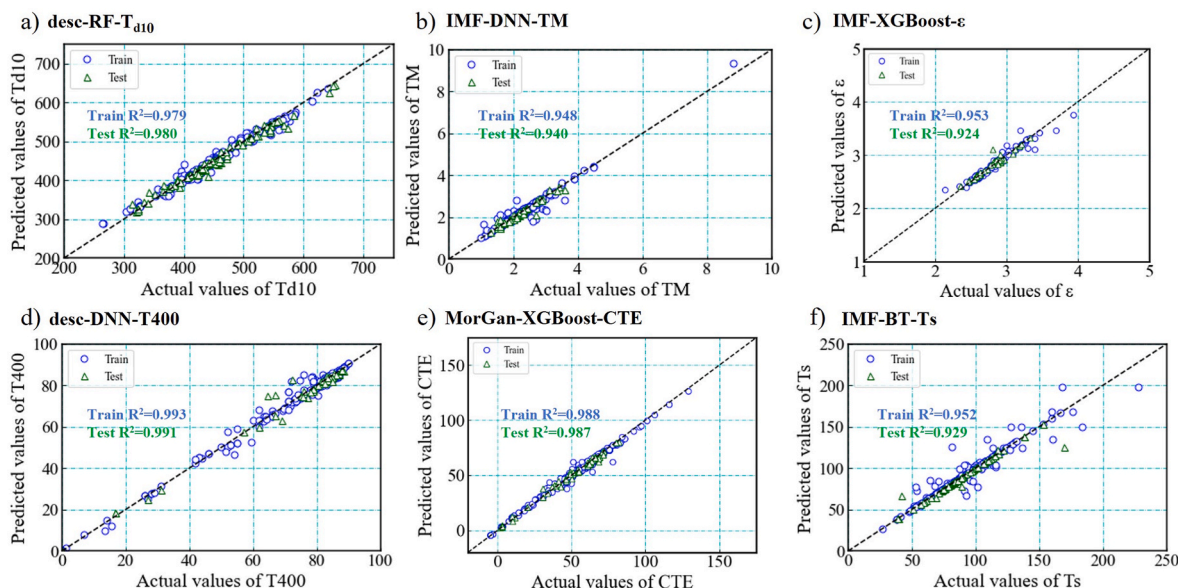
**Table 2**

Performance statistics of the optimal model for 11 properties prediction.

Model	Train		Test		Leave-one-out Cross-Validation	
	$R^2$	RMSE	$R^2$	RMSE	$Q^2(\text{LOO\_CV})$	Mean RMSE
IMF-RF-Tg	0.962	12.520	0.951	14.506	0.958	20.200
Mord-BT- $T_{d5}$	0.983	9.157	0.986	8.394	0.984	9.007
Desc-RF- $T_{d10}$	0.979	9.771	0.980	10.057	0.981	18.996
MorGan-XGBoost-CTE	0.988	2.496	0.987	2.272	0.988	2.462
IMF-BT- $T_s$	0.952	6.938	0.929	6.584	0.949	6.858
IMF-DNN- $T_M$	0.948	0.168	0.940	0.119	0.939	0.375
IMF-XGBoost- $\epsilon$	0.953	0.058	0.924	0.060	0.948	0.060
Desc-DNN-T400	0.993	2.253	0.991	2.288	0.986	6.619
IMF-RF- $\lambda_{\text{cutoff}}$	0.981	6.452	0.980	7.484	0.981	10.640
Mord-XGBoost- $n_{\text{av}}$	0.975	0.0111	0.957	0.0111	0.973	0.0110
Desc-DNN- $\Delta n$	0.989	0.0064	0.985	0.0062	0.984	0.0114

According to Table 2 we could intuitively determine the impact of feature representation and model selection on model performance. Among the 11 optimal models, 5 models adopted the IMF, 1 model adopted the original MorGan fingerprints, 3 models adopted the RDKit descriptor, and 2 models adopted the Mordred descriptor. Therefore, in terms of the number of optimal models, the IMF undoubtedly performed optimally well. Table S5 presents a complete picture of the impact of feature representation and model selection on model performance. For example, for the prediction of  $\epsilon$ , the optimal model was the XGBoost model built based on IMF, whose  $R^2$  exceeded 0.9 for both the training and test sets. In terms of fingerprint comparison, all the models built based on IMF outperform those built based on MorGan fingerprints. This was due to the fact that in the original MorGan, different substructures may be hashed into the same vectors, i.e., there was a collision problem in the original MorGan fingerprints, whereas the improved MorGan could greatly reduce the collision probability, which in turn could minimize the information loss.

The results of leave-one-out cross-validation (LOO\_CV) of the optimal model for the six properties were shown in Fig. 4, which



**Fig. 3.** (a)–(f) Plots of the optimal model performance for  $T_{d10}$ , TM,  $\epsilon$ , T400, CTE, and Ts, respectively.

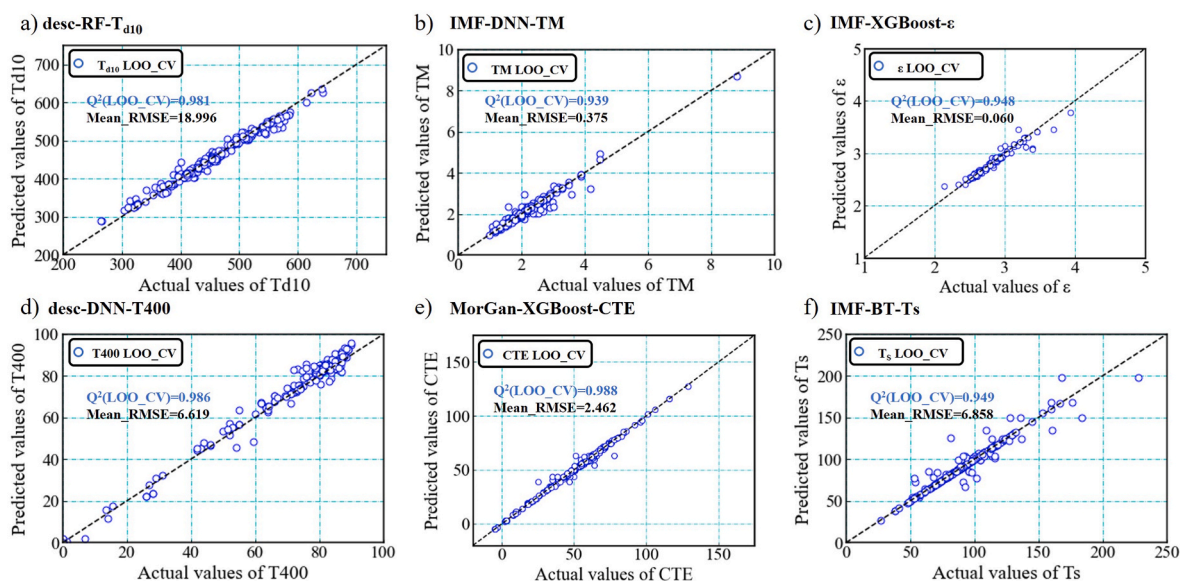


Fig. 4. (a)–(f) Distribution plots of the results of leave-one-out cross-validation for the optimal models of T<sub>d10</sub>, TM, ε, T400, CTE, and T<sub>s</sub>, respectively.

indicated that after undergoing cross-validation, the performance of each model was still excellent although there was a slight decrease in the performance ( $Q^2(\text{LOO\_CV})$  of T<sub>d10</sub>, TM, T400, CTE, ε, and T<sub>s</sub> were 0.981, 0.939, 0.986, 0.988, 0.948, and 0.949). In addition, the Mean\_RMSE for leave-one-out cross-validation was calculated. The results show that although the Mean\_RMSE for leave-one-out cross-validation has increased when compared to the training and test sets, the error was still within a reasonable range and the performance of the model was still reliable. For example, the Mean\_RMSE of T<sub>d10</sub> was 18.996 °C. This validation error of about 19 °C was excellent as it was the average RMSE of multiple cross-validations, proving the good generalization ability of the model and that there was no overfitting (See Figs. S10–S13 for details of the optimal model LOO\_CV for the remaining properties). We

collected 15–30 experimental values for each property from 24 recently published articles and further compared them to the ML predictions, which showed that the model's prediction errors were within a good range (Tables S14–16).

### 3.2. Interpretability analysis of the model using SHAP

After establishing machine learning regression relationship models for polyimide, we need a way to explain the key influencing factors in this regression relationship. In other words, we need to build a set of interpretable machine learning models that could be used to guide the design of future materials. The SHAP method is one of the most powerful machine learning interpretable analytical tools available today, and has

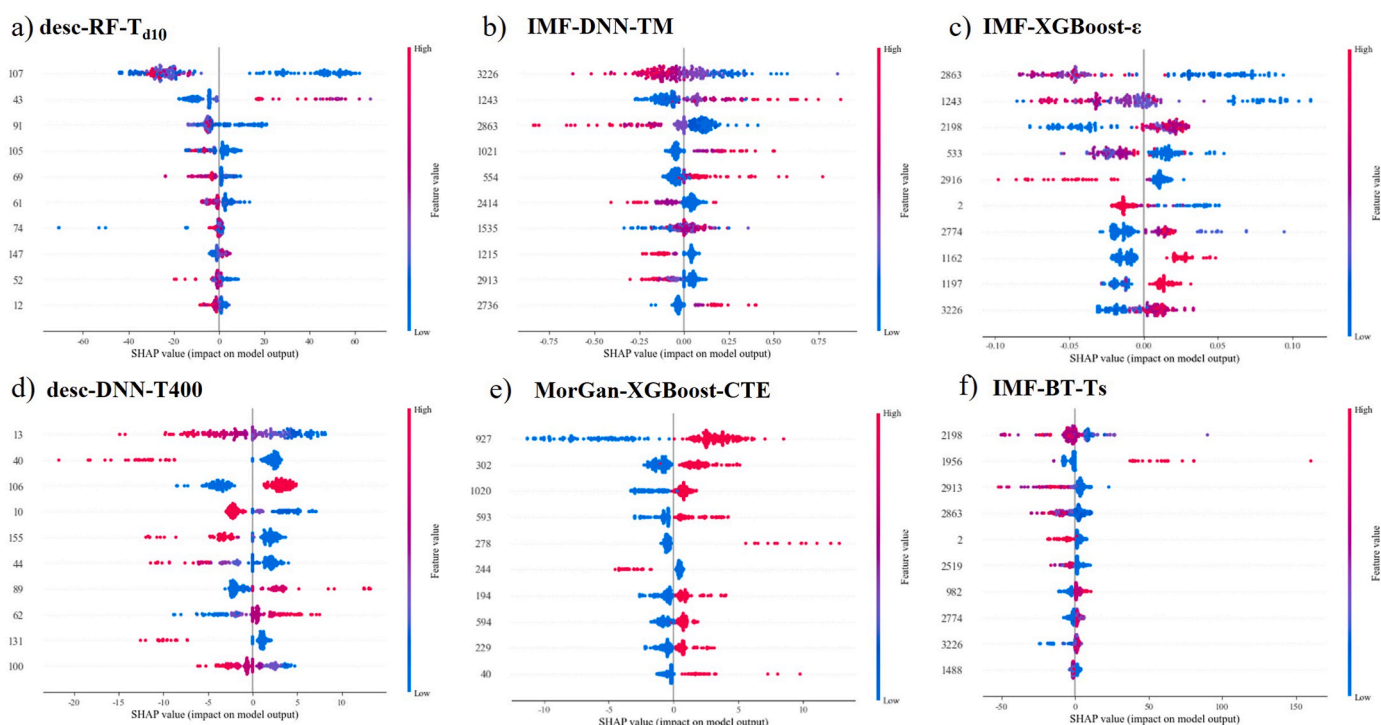


Fig. 5. (a)–(f) SHAP analysis plots of the optimal models for T<sub>d10</sub>, TM, ε, T400, CTE and T<sub>s</sub>, respectively.

been used in a number of research studies. The Shapely values computed by the SHAP analysis stem from game theory, and it could quantify each input variable's marginal contribution, which in turn quantifies the impact of each input variable on the final prediction result. In this study, the SHAP method is used to analyze the output of the model. According to SHAP, the significant contribution of each physicochemical descriptor or substructure to the model output results could be obtained.

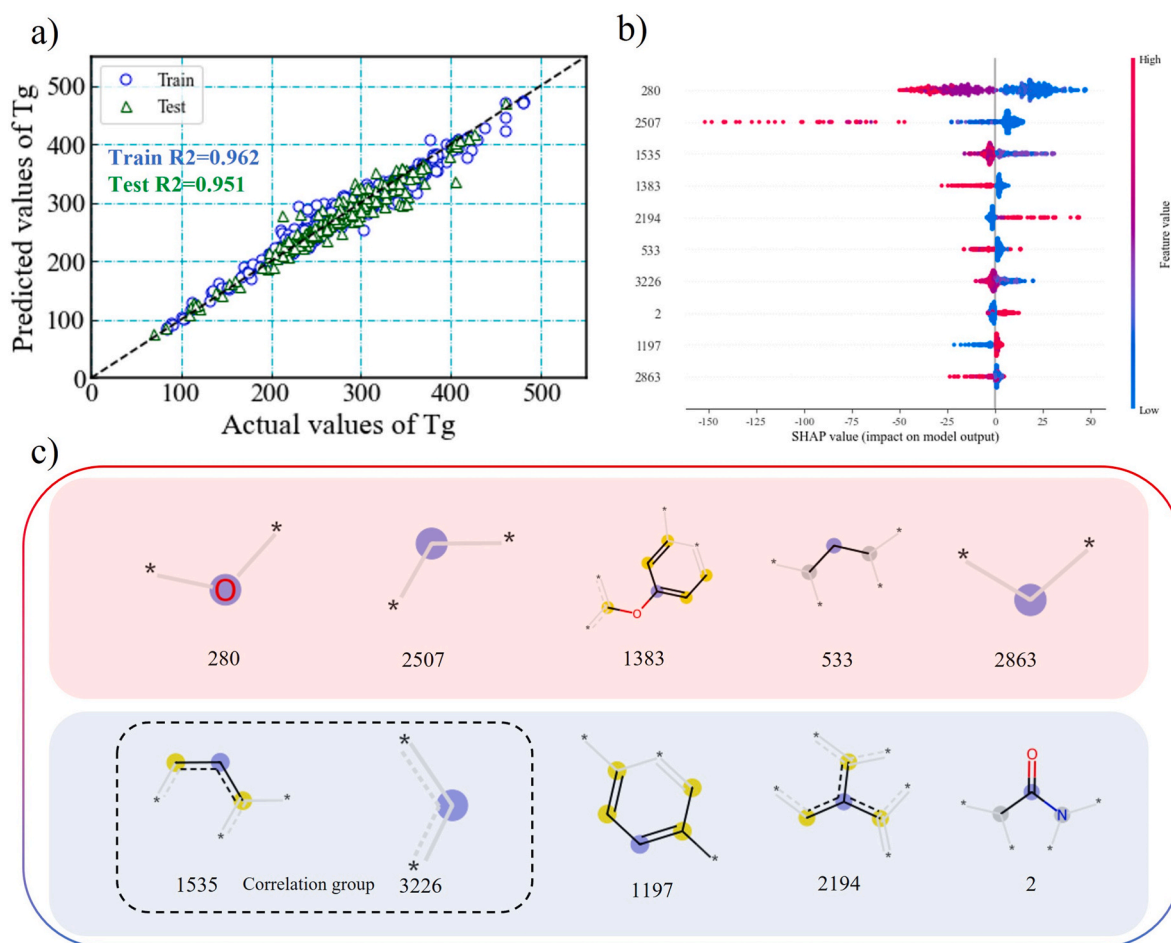
Here, the SHAP results of the optimal model based on the predictions of the six properties of polyimide were highlighted (see Figs. S18–S21 for the SHAP analysis of the remaining five properties). Fig. 5 displayed the SHAP values of the 10 most important input variables in the machine learning predictive models of the six properties. Based on these SHAP analyses, combined with the physicochemical meanings represented behind each descriptor (Table S3) or the substructures represented behind each fingerprint, a complete explanation of the predictive model for the target properties could be obtained. Meanwhile, by observing the commonalities between the most important input features of the property prediction models built based on the same feature representation, common determinants between multiple properties could be discovered. For example, there were common substructures 3226 and 2863 for the three models  $T_s$ ,  $T_M$ , and  $\epsilon$  built based on the IMF, 2863 was on behalf of the flexible structure  $-CH_2-$  (Fig. S3), and its existence could enhance the solubility and processability of polyimide, which could be beneficial to enhance the elongation at break of polyimide but at the same time it would lead to a decrease in the strength and modulus of polyimide [43]. In addition, the introduction of flexible units usually increased the free volume and intermolecular distance of the material, which could lead to a decrease in the dielectric constant [44]. Therefore, it could also be seen that 2863 has a negative contribution to  $T_s$ ,  $T_M$  and  $\epsilon$ . Through the fingerprint visualization analysis (Figs. S3–S5) and the molecular structure arrangement law of polyimide, it could be judged that 3226 was the partial molecular structure of benzene ring. This explained why 3226 was important in the prediction of all three properties - because for polymers, the presence or absence or number of benzene rings in the molecular structure would affect their physicochemical properties (Fig. S24 showed a detailed impact analysis for Feature 3226).

Similarly, a particular property could be analyzed separately in conjunction with SHAP. Firstly, for  $T_{d10}$ , descriptors 43 and 107 were its two important features. Descriptor 107 was NumSaturatedRings (the number of saturated rings in the molecule), which demonstrated a negative contribution to  $T_{d10}$ . Compared to more stable structures such as aromatic rings, saturated carbon-carbon single bonds could be prone to breakage at high temperatures, and therefore may lead to a decrease in the thermal stability of the material. Descriptor 43 was PEOE\_VSA6, which showed some positive contribution to  $T_{d10}$ . PEOE\_VSA6 responds to the size of the polar surface region in the molecule within a specific charge range, and higher polarity was usually accompanied by stronger intermolecular forces, and thus may improve the thermal stability of the material. For  $T_{400}$ , descriptor 40 (PEOE\_VSA3) was an important feature, reflecting the size of the polar surface regions in the molecule. The more polar regions may cause stronger intermolecular interactions, leading to tighter molecular packing. This may lead to a decrease in the optical transparency of the polyimide film. Thus, a negative contribution of descriptor 40 for  $T_{400}$  could also be seen. For CTE, substructure 927 was an important feature, which represented the flexible linkage  $-CH_2-$ , which contributes positively to the CTE. Flexible linkage units tend to increase the flexibility of the molecular chains of polyimides, which makes it easier for the molecular chains to stretch and move when subjected to heat. As a result, flexible linkage units tend to increase the linear coefficient of thermal expansion of polyimide.

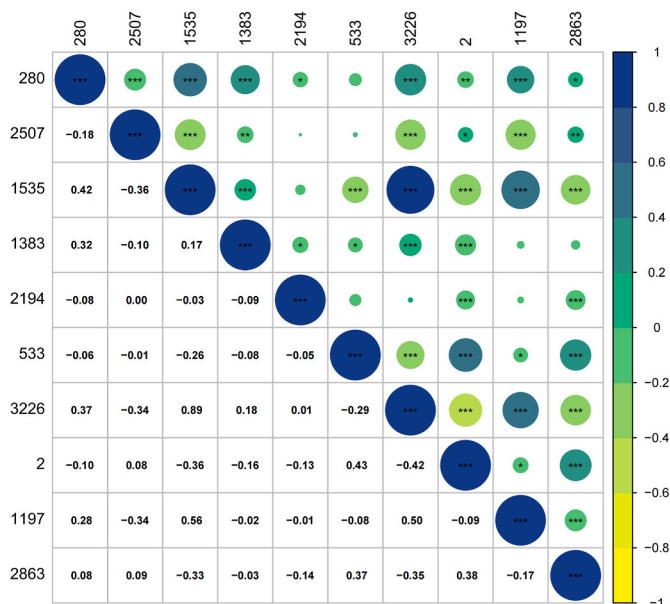
The glass transition temperature ( $T_g$ ) was a key property in the evaluation of polymers, reflecting the rigidity and heat resistance of the polymer molecule, and was an important criterion for evaluating polymer performance in several applications. Therefore, here the focus would be exclusively on the  $T_g$  model for polyimides and its interpretation.

Fig. 6 presented the  $T_g$  prediction model and its SHAP analysis results, which was the RF model built based on IMF. Fig. 6a showed the performance of the  $T_g$  optimal model, and the results showed that the  $R^2$  of this model was 0.962 and 0.951, and the RMSE was 12.520 and 14.506 on the training and test sets, respectively (see Table S5). Combined with the model's leave-one-out cross validation performance ( $Q^2 = 0.958$ , Mean\_RMSE = 20.2) (see Fig. S10), this demonstrated that the model for  $T_g$  prediction had excellent performance and good generalization ability. Based on this, Fig. 6b displayed the SHAP analysis results of the model, and the 10 most important substructures related to  $T_g$  were shown; for better understanding, these 10 substructures were visualized in Fig. 6c. The introduction of flexible chains such as  $-O-$ ,  $-SO_2-$ ,  $-S-$  and  $-CH_2-$  could lead to a decrease in the glass transition temperature of polyimide [43]. Combined with Fig. 6b–c, the associated substructures were represented in 280, 2507, 1383, and 2863, and all of their SHAP results were also negatively contributing to  $T_g$ . Substructure 533 belonged to the aliphatic ring, and according to SHAP, it had a negative contribution to  $T_g$ , which was consistent with the known theory [43]. Groups with positive contributions to  $T_g$  included 1197, 2194, and 2, the first two of which are aromatic units, and high glass transition temperatures were usually found in aromatic structures [45]. These interpretations supported the existing polyimide design theories, i.e., that  $T_g$  could be increased by the introduction of aromatic and hetero-aromatic rings, and that aliphatic structures or flexible linkage bonds could be introduced into the PI backbone to adjust the  $T_g$  to the desired size. Meanwhile, the SHAP results indicated that substructure 2 had a positive contribution to  $T_g$ . Analyzing the substructure from the functional group point of view, the substructure belonged to amide group, therefore, we introduced the idea that the presence of amide group might contribute to increase the glass transition temperature. This could be analyzed and understood from the following perspective: the amide group contained N and O, atoms that highly influenced the  $T_g$ , providing the possibility of the introduction of molecular polarity and the formation of hydrogen bonding, which contributed to the increase of intermolecular forces, which in turn might lead to an increase in the  $T_g$ . To explore whether these 10 substructures independently contributed to  $T_g$ , correlation thermograms were produced (see Fig. 7, and for the rest of the models see Figs. S14–S17). According to Fig. 7, most of the substructures were weakly correlated with each other, suggesting an independent effect on the  $T_g$  contribution. Although the substructures 1535 and 3226 showed a high correlation with each other, a deeper analysis revealed that 3226 was part of 1535 and they both belonged to the benzene ring. Substructures 3226 and 1535, identified as partially aromatic units in Fig. 8b, demonstrated a negative contribution to  $T_g$ . In fact, the reason for such a result was that the effect of the substructures on the properties was not invariant. Because each substructure does not affect the properties independently, there are synergistic effects between the substructures. The presence of synergistic effects with other substructures may lead the same substructure to make two opposite contributions to property in different polymers. In this case, the contribution should be approximated by combining the influential mean value of the feature over the entire data set.

Here, an alternative SHAP plot-force plot was used to further explain the effect of individual features on the performance of polyimide in the dataset. Figs. S26–27 showed the effects of features 3226 and 1535 on  $T_g$  in the whole dataset, respectively, and the results showed that these two features contributed positively to  $T_g$  at some of the frequencies of occurrence. And due to the interaction between the substructures, it also showed a negative contribution at some frequencies. This indicated that the effect of features 3226 and 1535 on the  $T_g$  of polyimide was also related to their frequency numbers of occurrence (NoteS4 shows a more detailed SHAP analysis). In addition, in Fig. 8b, features 3226 and 1535 contribute positively to PI-2, while feature 1535 contributes negatively to PI-3. This phenomenon further supports the conclusion that the effect of substructure on properties was not invariant.



**Fig. 6.** Interpretability analysis of the Tg model. (a) Optimal model for Tg prediction; (b) SHAP plot of the optimal model for Tg prediction; (c) visualization of the 10 most important substructures affecting the Tg model. Where blue was the central atom when the MorGan algorithm was executed, yellow was the aromatic atom, and gray was the aliphatic atom. (For interpretation of the references to color in this figure legend, the reader is referred to the Web version of this article.)



**Fig. 7.** Correlation analysis of the most important 10 substructures of the Tg model.

### 3.3. Molecular design according to SHAP analysis and experimental verification

Based on the interpretation of the above SHAP analysis and the guidance for molecular structure design, we experimentally designed three polyimides with different structures. Among these three PIs, PI-3 showed the lowest Tg due to the negative contribution of substructure 1535 to it (Fig. 8b). In contrast, PI-2 exhibited a higher Tg due to the positive contribution of both substructures 1535 and 3226 to it and the presence of no significant negatively contributing groups in PI-2. The design value of PI-1 was the highest, and although it did not have the presence of features 3226 and 1535 (3226 = 0 and 1535 = 0, i.e., not present), it detected the contribution of substructure 2 in the SHAP analysis, which was positively contributing to the increase in Tg, compared to PI-2 and PI-3 (Fig. 8a–b). Also, there was no significant negative contributing group present in PI-1 and it showed the presence of other positive contributing groups in its SHAP plot such as features 2223, 2392 and 2026. The predicted values of PI-1, PI-2 and PI-3 were 350.62 °C, 328.56 °C and 269.44 °C, respectively, and the syntheses of the three PIs were detailed in Note S1. DSC was tested on a NETZSCH DSC(TA) instrument between 30 and 400 °C at a heating rate of 5 °C/min. According to the DSC curves (Fig. 9a), the tested values of PI-1, PI-2 and PI-3 were 313.57 °C, 327.24 °C and 243.01 °C, respectively. Among them, PI-2 had the best prediction with an error of only 1.32 °C (Fig. 9b), while PI-1 had the largest error of 37.05 °C. This might be due to the fact that the present model was based only on the structure to predict the properties, and did not take into account the effects of different synthesis

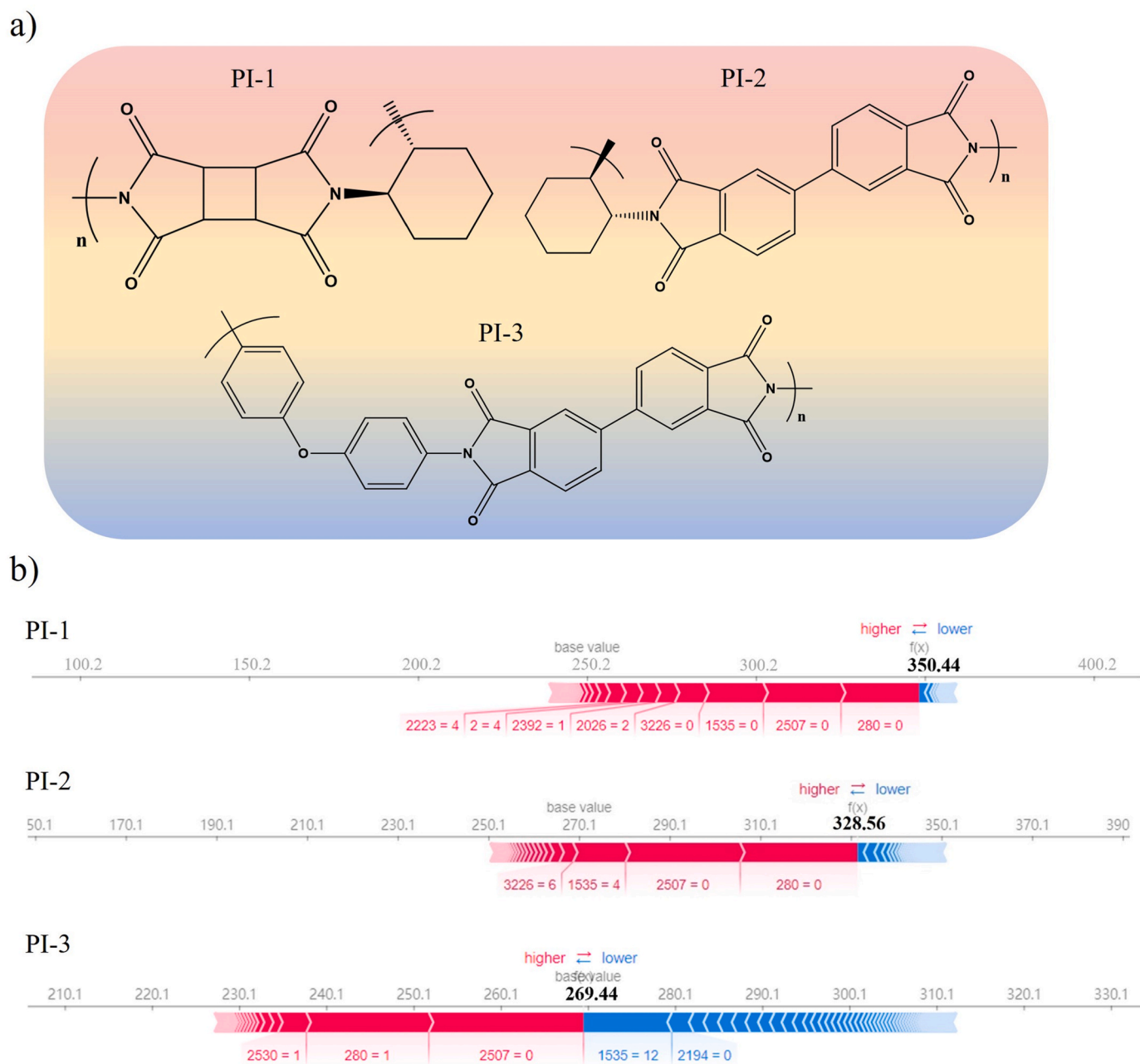
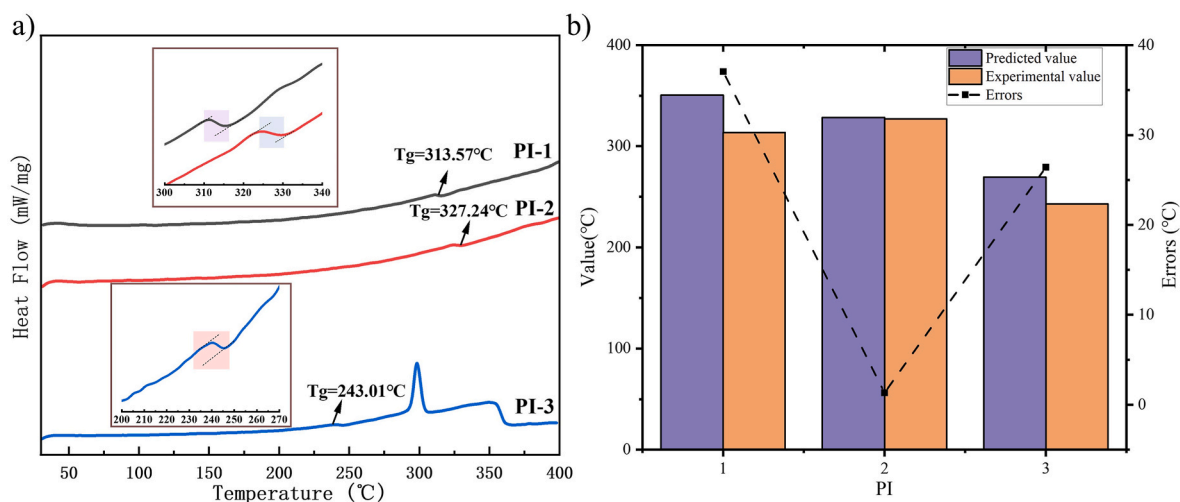


Fig. 8. Analysis of three PIs used for experimental synthesis and validation. (a) Structures of the three PIs. (b) SHAP analysis of the three PIs.

conditions and testing methods on the properties of polyimide, e.g., if the molecular weight was relatively low, it would probably lead to the synthesized PIs' T<sub>g</sub> values being low. The average error of the three PIs was about 21.6 °C, which was close to Q<sup>2</sup> (LOO.CV) (20.2 °C). Experimental tests of the remaining 10 properties and their comparisons with ML predictions are shown in Note S3, Fig. S23 and Tables S17–19. Overall, the prediction error of the present model was within a good range and therefore could be used for the development of new materials. In practical applications, when designing new PI resins, the expected design structure could be input into the model to quickly predict the performance of the PI to save time and experimental cost, and at the same time, the model could be utilized to explore more potential molecular design guides to accelerate the discovery of new PIs.

### 3.4. High throughput screening of resins for potential applications in various fields

After training a well-performing ML model, we would prefer to apply the model to select potential high-performance polymers for various fields in a larger chemical space. Therefore, in this study, a high-throughput screening was conducted with PI applications in aerospace, optoelectronics and microelectronics as an example. For better visualization, it was first shown as an example of 3263 data, which contained 973 PIs from dataset 1 and 2290 hypothetical polyimides from dataset 2. When PIs were applied in aerospace [46], the requirements of heat resistance and mechanical strength were considered first, so two properties, T<sub>d10</sub> and TM, were selected as the screening, and the requirements of the two properties were T<sub>d10</sub> > 525 °C and TM > 2.2 GPa, respectively; when PIs were applied in the field of optoelectronics [47], the requirements of their light transmittance and coefficient of



**Fig. 9.** (a) DSC test results; (b) comparison of predicted and experimental values. Note: PI-3 already showed >5 % pyrolytic loss at 300 °C, so the fluctuation of the curve at this temperature was not considered as a glass transition.

thermal expansion were taken into account, and thus T400 and CTE were used as screening properties, requiring  $T_{400} > 85\%$  and  $CTE < 50 \text{ ppm K}^{-1}$ ; similarly, considering the requirements of dielectricity and mechanical strength when PI was applied in the field of microelectronics [48],  $T_s$  and  $\epsilon$  were selected, requiring  $T_s > 200 \text{ MPa}$  and  $\epsilon < 2$ . Finally, considering the ease of synthesis of these real or hypothetical polyimides, the assessment of the synthesized accessibility score (SAscore) was chosen as a judgment (SAscore definitions and calculation codes presented in [Appendix Note S2](#)). The results of the high-throughput screening of the three domains were shown in [Fig. 10](#), where the polymers screened for the desired properties for each domain were framed, and the color response represented as the SA\_scores value. As well, repeating units of high-performance PIs with SA\_scores <5.5 were listed. The high-performance PIs screened here were all identified based on real-life dianhydrides, diamines, or diisocyanates based on the rules of PI synthesis. However, the synthesis of polymers was a very complex process, not only relying on the screening of structures, but also the selection of suitable solvents and the synthesis process were very important. In this study, we have not focused on exploring the synthesis process of these potential polymers, but we have summarized the structures of these high-performance polymers. In order to accelerate the synthesis process for future researchers, we have screened potential solvents for these high-performance polymers using Polymer Genome, another large machine learning database. (Monomers and solvents for high-performance PIs are listed in Supporting Information Tables S6-S12). These numbers of high-performance PIs were very sparse when screened only on a small scale. [Fig. S22](#) displayed the results of high-throughput screening of nearly 7.6 million hypothetical PIs in dataset II, clearly reflecting that the number of high-performance PIs in each domain was greatly increased, and thus nearly two hundred potential high-performance PIs with SA\_scores <4 that were easy to synthesize were selected for each domain separately in this study (excel). This work could greatly reduce the repetitive and time-consuming work of experimentalists, and would be expected to help the synthesis and application of high-performance PI materials in various fields.

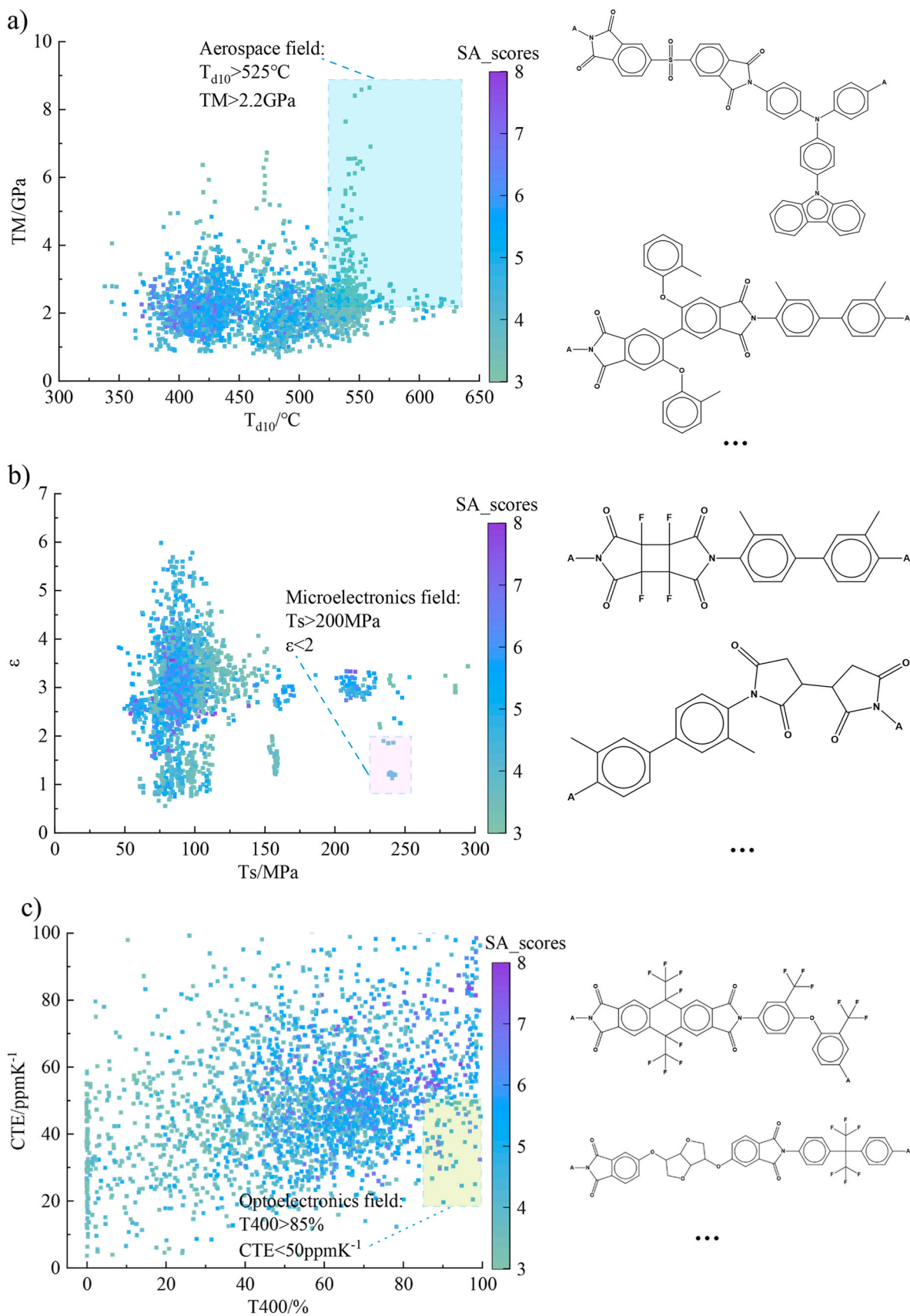
#### 4. Conclusion

In this study, we focused on several representative types of properties of polyimides that were of importance in various fields of application, including thermal ( $T_g$ ,  $T_{d5}$ ,  $T_{d10}$ , and CTE), mechanical ( $T_s$  and  $TM$ ), dielectric ( $\epsilon$ ), and optical ( $\lambda_{\text{cutoff}}$ , T400,  $n_{\text{av}}$ , and  $\Delta n$ ) properties, for a total of 11 properties. Four feature representations such as MorGan

fingerprint, IMF, RDKit and Mordred descriptor were considered and the effects of four models such as DNN, RF, XGBoost and BT were compared. External validation, leave-one-out cross-validation, and experimental validation were used to determine the performance and generalization ability of each property prediction model. For a focused discussion, 1–2 representative properties ( $T_g$ , CTE,  $T_{d10}$ ,  $T_s$ ,  $TM$ ,  $\epsilon$ , and T400) from each of the four categories of properties were selected for discussion and analysis in the main text. The results showed that the  $R^2$  of the optimal prediction model for all 11 properties was higher than 0.90. SHAP was used to interpret each optimal model from the physical-chemical point of view as well as from the substructural point of view. In analyzing  $T_g$ , we found that the amide substructure may contribute positively to  $T_g$ . Therefore, in this study, from the experimental point of view, we designed three PIs with different structures and properties of the three polyimides, which contained the key elements found in SHAP analysis, such as the aromatic ring, aliphatic units, flexible chain segments, and the newly discovered amide substructure, respectively. The final results of the experimental testing of the three resins were in good agreement with the ML predictions. Based on the existing real and hypothetical polyimides (nearly 7.6 million), a high-throughput screening of six properties of interest to the aerospace, optoelectronics, and microelectronics fields was performed. The synthesis of these potential polymers was also taken into account, and the ease of synthesis of each polymer was evaluated in terms of SA\_scores, resulting in the selection of high-performance potential polymers with SA\_scores <4 for each field. In conclusion, the present PI multi-property ML prediction model combined high performance with interpretability, which could significantly reduce the time cost of experimental trial-and-error and simulation. For the future, the high-performance model and its interpretable analysis could be used to rationally regulate the design and selection of PI molecules, which could contribute to the synthesis of high-performance polyimides and expand their applications in various fields. In future studies we should collect data on polyimide processing conditions, heat treatment, crystallinity, etc., and consider predicting the properties of different polyimide aggregation state structures in more depth.

#### Data availability statement

The DOIs of the data source literature used for machine learning model training, testing, and validation in this paper are included in the Supporting Information. Hyperparameters for machine learning model training are also detailed in the Supporting Information. All relevant and specific data in the paper, including the raw data for model training, testing and validation, 7.6 million high-throughput screened data, and



**Fig. 10.** High-throughput screening of the ML model in three fields. (a)–(c) The results of high-throughput screening for aerospace, microelectronics and optoelectronics fields with some of their potential high-performance polymers, respectively.



- magnesium dissolution modulators using sparse machine learning models, *npj Comput. Mater.* 7 (1) (2021).
- [41] D. Rogers, M. Hahn, Extended-connectivity fingerprints, *J. Chem. Inf. Model.* 50 (5) (2010) 742–754.
- [42] L. Tao, V. Varshney, Y. Li, Benchmarking machine learning models for polymer Informatics: an example of glass transition temperature, *J. Chem. Inf. Model.* 61 (11) (2021) 5395–5413.
- [43] Y.B. Zhuang, J.G. Seong, Y.M. Lee, Polyimides containing aliphatic/alicyclic segments in the main chains, *Prog. Polym. Sci.* 92 (2019) 35–88.
- [44] C. Qian, Z.-G. Fan, W.-W. Zheng, R.-X. Bei, T.-W. Zhu, S.-W. Liu, Z.-G. Chi, M. P. Aldred, X.-D. Chen, Y. Zhang, J.-R. Xu, A facile Strategy for non-fluorinated intrinsic low-k and low-loss dielectric polymers: Valid exploitation of secondary relaxation Behaviors, *Chin. J. Polym. Sci.* 38 (3) (2020) 213–219.
- [45] D.J. Liaw, K.L. Wang, Y.C. Huang, K.R. Lee, J.Y. Lai, C.S. Ha, Advanced polyimide materials: syntheses, physical properties and applications, *Prog. Polym. Sci.* 37 (7) (2012) 907–974.
- [46] S.J. Rinehart, B.N. Nguyen, R.P. Viggiano, M.A.B. Meador, M.D. Dadmun, Quantitative evaluation of the hierarchical porosity in polyimide aerogels and corresponding solvated Gels, *ACS Appl. Mater. Interfaces* 12 (27) (2020) 30457–30465.
- [47] S.L. Vivod, M.A.B. Meador, C. Pugh, M. Wilkosz, K. Calomino, L. McCorkle, Toward improved optical transparency of polyimide aerogels, *ACS Appl. Mater. Interfaces* 12 (7) (2020) 8622–8633.
- [48] S.Z. Li, Z.B. Zheng, S.W. Liu, Z.G. Chi, X.D. Chen, Y. Zhang, J.R. Xu, Ultrahigh thermal and electric conductive graphite films prepared by g-C<sub>3</sub>N<sub>4</sub> catalyzed graphitization of polyimide films, *Chem. Eng. J.* 430 (2022).