

---

# Graph Convolutional Neural Networks for Polymers Property Prediction

---

**Mingang Zeng**

Deep Learning Department, Institute for Infocomm Research,  
A\*STAR (Agency for Science, Technology and Research)  
1 Fusionopolis Way, #21-01 Connexis, Singapore 138632

**Jatin Nitin Kumar**

Soft Materials, Institute of Materials Research and Engineering  
2 Fusionopolis Way, #08-03 Innovis, Singapore 138634

**Zeng Zeng**

Deep Learning Department, Institute for Infocomm Research,  
A\*STAR (Agency for Science, Technology and Research)  
1 Fusionopolis Way, #21-01 Connexis, Singapore 138632

**Ramasamy Savitha**

Deep Learning Department, Institute for Infocomm Research,  
A\*STAR (Agency for Science, Technology and Research)  
1 Fusionopolis Way, #21-01 Connexis, Singapore 138632

**Vijay Ramaseshan Chandrasekhar \***

Deep Learning Department, Institute for Infocomm Research,  
A\*STAR (Agency for Science, Technology and Research)  
1 Fusionopolis Way, #21-01 Connexis, Singapore 138632  
vijay@i2r.a-star.edu.sg

**Kedar Hippalgaonkar \***

Electronic Materials, Institute of Materials Research and Engineering  
2 Fusionopolis Way, #08-03 Innovis, Singapore 138634  
kedarh@imre.a-star.edu.sg

## Abstract

A fast and accurate predictive tool for polymer properties is demanding and will pave the way to iterative inverse design. In this work, we apply graph convolutional neural networks (GCNN) to predict the dielectric constant and energy bandgap of polymers. Using density functional theory (DFT) calculated properties as the ground truth, GCNN can achieve remarkable agreement with DFT results. Moreover, we show that GCNN outperforms other machine learning algorithms. Our work proves that GCNN relies only on morphological data of polymers and removes the requirement for complicated hand-crafted descriptors, while still offering accuracy in fast predictions.

---

\*equal corresponding author

## 1 Introduction

Polymers are materials with tunable structures and chemical functionality that influence their physical and chemical properties.[1] A deep understanding of the relationship between structure and properties is required for innovating novel materials. However, this relationship is complex and often not well understood.[2] Density Functional Theory (DFT) is useful in estimating bulk polymer properties, but is computationally expensive. It has been applied on small molecules, and a 20,000 size dataset was generated as a part of the Harvard Clean Energy Project [3]. Using a subset from this data, by representing molecules as graphs (with individual atoms as vertices and bonds as edges), a neural fingerprinting approach expanding upon a Simplified Molecular-Input Line-Entry System (SMILES) enabled high predictive capability of solubility, drug efficacy and photovoltaic efficiency [4]. Following this, a reinforced adversarial neural computer based on reinforcement learning allowed access to a broad chemical space towards predictive synthesis of small molecules [5]. Moving beyond small molecules, to enable property prediction for bulk polymers, high-throughput DFT calculations were performed on different repeat units to provide a “fingerprint”, from which the final polymer properties were derived algorithmically.[6] Bypassing DFT for fingerprinting and instead using direct morphological information would be a significant advance in the discovery of new materials and shorten the development time for material innovation. To achieve this, a strong understanding of the relevant and most useful material descriptors influencing the polymer properties is required.

Describing the polymer in its most basic form typically constitutes identifying the atoms that make up the material as well as their arrangement with respect to each other, determined by the type of bonding that holds the material together. For inorganic crystals, where long range order is necessary, the material can be described as a lattice and a basis, where the lattice reflects the symmetry of the crystal structure, while the basis is the repeating unit. Polymers, on the other hand, are largely amorphous. This is because they do not exhibit such long range order as they do not conform to any lattice unless specifically phase controlled. However, certain polymeric properties are intrinsic to the monomer unit and translate well to actual applications. This is especially true for ground state properties such as the dielectric constant and the polymer bandgap. The dielectric constant is a reflection of the electronic polarizability of the polymer and is a consequence of the bonding nature between the constituent atoms. Typically, the bandgap depends upon the strength of the bonding. Therefore, one would expect that knowing the atoms that constitute the monomer as well as how they bond with each other - resulting in a fingerprint morphological character - is sufficient information to predict these properties. Such information is uniquely contained in a polymer’s Crystallographic Information File (CIF), which can then be converted to a two-dimensional (2D) graph and used as an input into a convolutional neural network to predict the above-mentioned properties.

In addition to morphological character intrinsic to the polymer, typically, environmental considerations are also important. Their physio-chemical environment is a complex multi-variable parameter space; intra and inter chain effects due to multi-valency and chain coiling affects functional properties, but is expected to be less influential for ground state properties. More complex descriptors that have been used in complementary approaches involve a combination of detailed atomic and morphological information along with environmental interactions.[7] In our work, we predict the dielectric constant and the bandgap of a large class of polymer compounds using the elegant graph convolutional neural network (GCNN) and find that the predictions are better than those where complex descriptors have been used in the recent past as illustrated in Figure 1. Specifically, by comparing to traditional descriptors, our mean absolute errors are lower than those obtained by other machine learning techniques. This clarifies that to achieve speed and accuracy in predictions, we only need to consider the morphological character of polymers.

## 2 Dataset

To investigate the accuracy of GCNN on predicting physical properties of polymers, we use the publicly available dataset at the Polymergenome Project.[8, 9, 10, 11, 7] The main purpose of this dataset is to develop polymers for energy storage and electronics applications. The dataset covers 1073 polymers that can be classified as organic or organometallic which can be further divided into three subsets according to their sources. The first subset has 34 common polymers that have been synthesized in experiments, like polyethylene, polyureas, polythioureas, polyesters, *etc.*. The second subset adopted from the Crystallography Open Database (COD) contains another 253 organic and

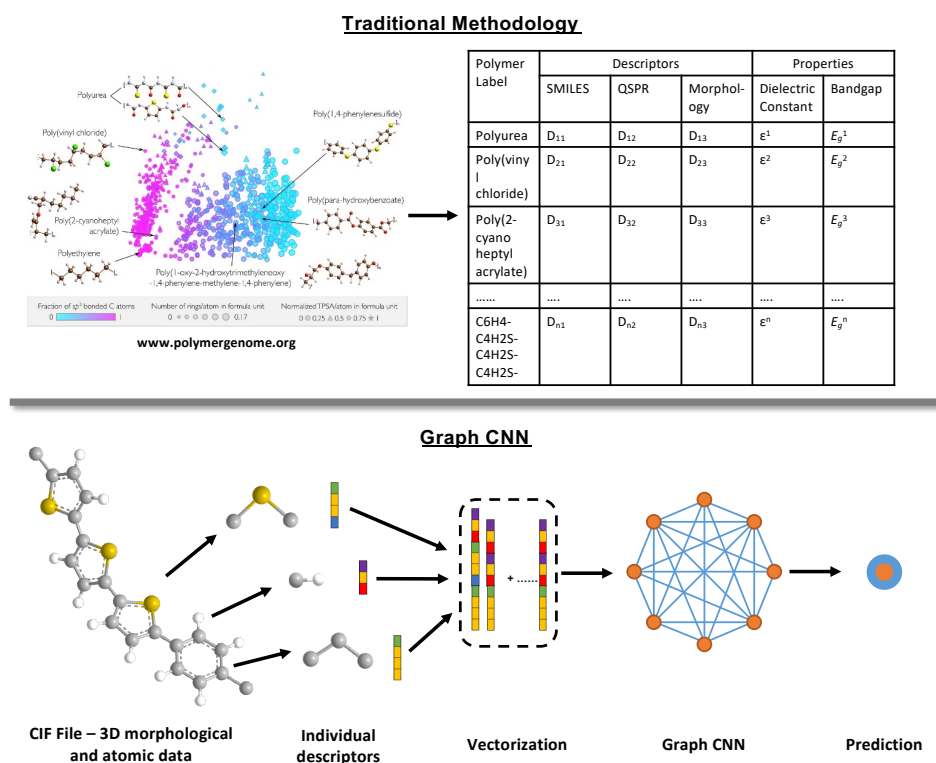


Figure 1: Schematic comparison between GCNN and other machine learning algorithms. The traditional methodology relies on hand engineered features crafted from SMILE string, quantitative structure-property relationship and morphology of polymers. One remarkable advantage of graph CNN is its ability to automatically learn the chemical environment of polymers and map polymer structure to abundant feature vectors for a fast and accurate prediction on polymer properties.

organometallic polymers. The last subset generated from computational methods includes 314 organic polymers and 472 organometallic polymers. The polymer building blocks in this subset include organic  $-\text{CH}_2-$ ,  $-\text{NH}-$ ,  $-\text{CO}-$ ,  $-\text{O}-$ ,  $-\text{CS}-$ ,  $-\text{C}_6\text{H}_4-$ , and  $-\text{C}_4\text{H}_2\text{S}-$ , as well as metal-containing inorganic blocks, such as  $-\text{COO}-\text{Sn}(\text{CH}_3)_2-\text{OCC}-$ ,  $-\text{SnF}_2-$ , and  $-\text{SnCl}_2-$ . The polymer building block repeats along the 3-dimensional crystal axis, and the repeat unit with symmetry information is fed into the GCNN.

### 3 Experimental Setup

#### 3.1 Input

To build a regression model, the input of GCNN includes CIF files recording the structure of the polymers, the target properties for each polymer and a JSON file that stores the initialization vector for each atom.

#### 3.2 Model Architecture

A polymer graph can be represented by nodes and edges using the atomic feature vector  $\mathbf{v}_i$  and the bonding feature vector  $\mathbf{u}_{(i,j)_k}$ , respectively.[4, 12, 13, 14] These vectors are obtained by one hot encoding. With the help of a non-linear convolution function, the multiple convolutional layers automatically learn the atomic features after iterating through surrounding chemical features (atoms/bonds) and different convolutional layers. A pooling layer of normalized summation is then

Table 1: Hyperparameters for different algorithms

Method	Hyperparameters
GCNN for $\varepsilon$ ( $E_g$ )	Learning rate: 0.001 (0.01), momentum: 0.9 (0.8), hidden-feature-length: 256 (64), weight decay: 1e-8 (1e-5), atomic-feature-length: 32 (32), number of CNN layers: 2 (5), number of hidden layer: 2 (5)
Kernel Regression	alpha: 1e-05, gamma: 1.25e-09, kernel: rbf
Random Forest	min_samples_leaf: 10, n_estimators: 150, oob_score: True
Gradient Boosting	alpha: 0.7, learning_rate: 0.1, max_depth: 5
Neural Network	alpha: 0.001, hidden_layer_sizes: 100, momentum: 0.7

used to generate an overall feature vector  $\mathbf{v}_c$ . To improve the performance of the GCNN, two fully connected hidden layers are added to capture the complicated structure-property relationship. Finally, an output layer connected to the top hidden layer is used to predict the physical property of polymers.

### 3.3 Hyperparameter Optimization

The database is randomly divided into training, validation and test sets with the ratio of 6:2:2. The network parameters are optimized via Stochastic Gradient Descent (SGD), and the optimum hyperparameters are determined by the lowest mean absolute error (MAE) in the validation set using the DFT result as the ground truth.

## 4 Predictive Performance

We ran two experiments to demonstrate that GCNN can obtain comparable accuracy to DFT for polymers. Here we focus on the bandgap ( $E_g$ ) and the dielectric constant ( $\varepsilon$ ). These two properties are important in order to screen polymer materials regardless of specific applications. The DFT results of  $E_g$  and  $\varepsilon$  are used as the ground truth. They are obtained with hybrid electron exchange-correlation functionals and density functional perturbation theory, respectively.[8] The optimized hyperparameters of GCNN for  $E_g$  and  $\varepsilon$  are listed in Table 1. Figures 2(a,b) compare the predictive performance of GCNN versus the DFT ground truth. Impressive agreement with DFT is found in predicting the dielectric constant. A MAE value of 0.24 is achieved on the test set, which is lower than the published work using a similar dataset and Gaussian process regression.[7] Given the error of DFT calculations compared with experiments and the small MAE obtained in our study, GCNN may achieve accuracy in predicting polymer properties as compared to experiments. Comparatively, a higher MAE (0.41 eV) is found using GCNN to predict the energy bandgap of the polymer dataset. We illustrate the data distribution for  $E_g$  and  $\varepsilon$  in Figure 2(c). The dielectric constant of the dataset is mainly concentrated in the region of  $\varepsilon = 3$ , with the mean and variation of 3.47 and 0.92, respectively. Comparatively, the  $E_g$  is much more dispersed (mean = 4.417 eV, variation = 2.75 eV) with an obvious tail weighted in the high bandgap range. Therefore, the MAE is high since lesser data is available for the polymer in this energy bandgap tail region. The inset of Figure 2(b) shows the MAE for  $E_g$  as the function of dataset size. We find a systematic decrease in MAE with increasing the size of dataset from 128, 256, 512 to 1024. This indicates the prediction of  $E_g$  by GCNN can be improved if more polymer data are provided.[7].

## 5 Comparison with other Machine Learning methods

Besides GCNN, where we only use the cif and the atomic .json files as inputs, we also use other Machine Learning (ML) methods to establish a regression model for the polymer structure-property relationship. The feature extraction of the polymer dataset is implemented with the Matminer package.[15] Table 3 lists the feature extraction modules that we apply, which produces 158 individual feature descriptors. These hand-crafted features reflect some key physical and chemical properties of polymer systems, such as structural heterogeneity and chemical ordering. Based on these physically relevant descriptors, we apply a series of machine learning algorithms using the Scikit-learn package

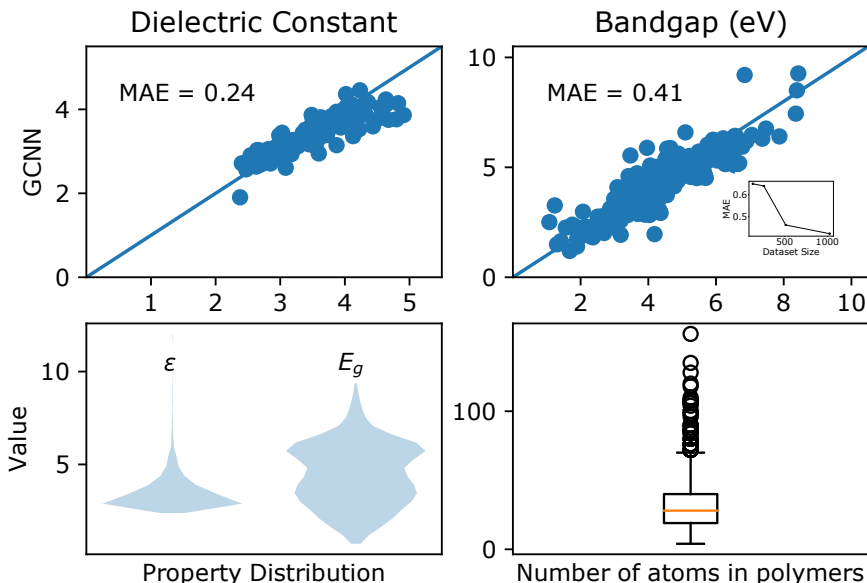


Figure 2: (a,b) The predicted dielectric constant and bandgap by GCNN compared with the DFT results, showing that the GCNN is an accurate predictive tool. The inset of (b) shows a systematic decrease in MAE with increasing the size of dataset from 128, 256, 512 to 1024. (c) Value distribution showing a narrow range for the dielectric constant and a wider range for the bandgap. (d) Statistics showing the large range in the number of atoms in each polymer; the GCNN does not discriminate against large polymers and is still able to achieve good prediction.

Table 2: Summary of the prediction performance of three different properties by different ML methods

Property	MAE				
	GCNN	Kernel Regression	Random Forest	Gradient Boosting	Neural Network
$\epsilon$	0.24	0.425	0.355	0.359	0.59
$E_g(eV)$	0.41	0.652	0.505	0.446	0.509

to study the structure-property relationship of the polymer dataset. The ML models have internal 4-fold cross-validation to minimise over-fitting and ensure model generality. The models are trained using SGD with the ADAM optimizer; and the GridSearchCV method in Scikit-Learn is applied to tune hyperparameters. The optimized hyperparameters are listed in Table 1; and the MAE values of these ML algorithms are listed in Table 2. It can be seen that GCNN outperforms all these ML algorithms in predicting the  $E_g$  and  $\epsilon$  within this polymer dataset. This indicates that the feature vectors generated from the GCNN function better in predicting the physical properties of polymers. This may come from the fact that GCNN feature vectors take into account global spatial geometry, as well as accurate local atomic configuration, as listed in Table 3. As shown by the statistic analysis of the polymer dataset in Figure 2(d), GCNN works well even on long polymer chains, regardless of a simplified graphical representation. Our results suggest that GCNN has the potential to serve as an excellent forward predictive model for polymers.

## 6 Discussion and Future Work

Given the strength of the GCNN predictions, it is good to note that the present approach has two limitations. The first is that the ground truth is based on simulations rather than experimental data, which is a result of the difficulty of performing high-throughput experiments and hence sparse experimental data. The second is that the present technique only allows for the prediction of bulk polymeric properties as opposed to polymer-solvent interactions and composites. While the first limitation is due to lack of data availability, the second could be addressed by investigating

Table 3: Comparison between the GCNN generated feature classes and the hand-crafted feature classes obtained via the Matminer package. The Kernel Regression, Random Forest, Gradient Boosting and Neural Network algorithms use these listed hand-crafted feature classes.

Hand-crafted feature classes	GCNN generated feature classes
Structural Heterogeneity	Group number
Chemical Ordering	Period number
Maximum Packing Efficiency	Electronegativity
Stoichiometry	Covalent radius
Element Property	Valence electrons
Valence Orbital	First ionization energy
Ionic Property	Electron affinity
	(s,p,d,f) Block
	Atomic volume
	Atomic distance

feature importance using the GCNN to allow for interpretability. We seek inspiration from a recent work on developing a methodology of vectorizing individual molecular descriptors via multi-dimensional correlation to account for such complexities, thereby enabling the discovery of novel functional molecules.[16] This will then allow us to build an accurate yet fundamental design-property understanding which in turn could predict experimental polymer behaviour. Moreover, this sort of understanding could also extend the capability beyond bulk polymers to polymer composites and polymers in solution. Most importantly this strong understanding of polymers will allow us to predict the required chemical and structural design of a polymer based on its physical property requirement - the problem of inverse design - which would lead to huge impact in both industry and academia.

## 7 Conclusion

We applied several machine learning algorithms to predict dielectric constant and bandgap from a large dataset of crystallography data of polymers. Our models included graph convolution neural network, random forest, kernel regression, gradient boosting and conventional neural network, where the lowest mean absolute error for both properties were reported for GCNN. These results indicate that GCNN offers an effective approach for fast and accurate prediction of polymer properties starting from the atomic and morphological character of polymer. This conceptual advance allows us to move into the realm of faster predictions, relying on a smaller amount of metadata, and paves the way for inverse design, where polymers can be designed with a final property in mind.

## References

- [1] Andrew Gregory and Martina H Stenzel. Complex polymer architectures via raft polymerization: From fundamental process to extending the scope using click chemistry and nature’s building blocks. *Progress in Polymer Science*, 37(1):38, 2012.
- [2] Charles E Carraher Jr and RB Seymour. *Structure—Property Relationships in Polymers*. Springer Science & Business Media, 2012.
- [3] Johannes Hachmann, Roberto Olivares-Amaya, Sule Atahan-Evrenk, Carlos Amador-Bedolla, Roel S Sánchez-Carrera, Aryeh Gold-Parker, Leslie Vogt, Anna M Brockway, and Alán Aspuru-Guzik. The harvard clean energy project: large-scale computational screening and design of organic photovoltaics on the world community grid. *The Journal of Physical Chemistry Letters*, 2(17):2241, 2011.
- [4] David K Duvenaud, Dougal Maclaurin, Jorge Iparraguirre, Rafael Bombarell, Timothy Hirzel, Alán Aspuru-Guzik, and Ryan P Adams. Convolutional networks on graphs for learning molecular fingerprints. In *Advances in neural information processing systems*, page 2224, 2015.
- [5] Evgeny Putin, Arip Asadulaev, Yan Ivanenkov, Vladimir Aladinskiy, Benjamin Sánchez-Lengeling, Alán Aspuru-Guzik, and Alex Zhavoronkov. Reinforced adversarial neural computer for de novo molecular design. *Journal of chemical information and modeling*, 2018.

- [6] Arun Mannodi-Kanakkithodi, Ghanshyam Pilania, Tran Doan Huan, Turab Lookman, and Rampi Ramprasad. Machine learning strategy for accelerated design of polymer dielectrics. *Scientific reports*, 6:20952, 2016.
- [7] Chiho Kim, Anand Chandrasekaran, Tran Doan Huan, Deya Das, and Rampi Ramprasad. Polymer genome: A data-powered polymer informatics platform for property predictions. *The Journal of Physical Chemistry C*, 2018.
- [8] Tran Doan Huan, Arun Mannodi-Kanakkithodi, Chiho Kim, Vinit Sharma, Ghanshyam Pilania, and Rampi Ramprasad. A polymer dataset for accelerated property prediction and design. *Scientific data*, 3:160012, 2016.
- [9] Arun Mannodi-Kanakkithodi, Tran Doan Huan, and Rampi Ramprasad. Mining Materials Design Rules from Data: The Example of Polymer Dielectrics. *CHEMISTRY OF MATERIALS*, 29(21):9001, 2017.
- [10] Arun Mannodi-Kanakkithodi, Ghanshyam Pilania, and Rampi Ramprasad. Critical assessment of regression-based machine learning methods for polymer dielectrics. *COMPUTATIONAL MATERIALS SCIENCE*, 125:123, 2016.
- [11] Arun Mannodi-Kanakkithodi, Anand Chandrasekaran, Chiho Kim, Tran Doan Huan, Ghanshyam Pilania, Venkatesh Botu, and Rampi Ramprasad. Scoping the polymer genome: A roadmap for rational polymer dielectrics design and beyond. *Materials Today*, 2017.
- [12] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. *arXiv preprint arXiv:1704.01212*, 2017.
- [13] Tian Xie and Jeffrey C Grossman. Hierarchical visualization of materials space with graph convolutional neural networks. *arXiv preprint arXiv:1807.03404*, 2018.
- [14] Tian Xie and Jeffrey C. Grossman. Crystal graph convolutional neural networks for an accurate and interpretable prediction of material properties. *Phys. Rev. Lett.*, 120:145301, Apr 2018.
- [15] Logan Ward, Alexander Dunn, Alireza Faghaninia, Nils ER Zimmermann, Saurabh Bajaj, Qi Wang, Joseph Montoya, Jiming Chen, Kyle Bystrom, Maxwell Dylla, et al. Matminer: An open source toolkit for materials data mining. *Computational Materials Science*, 152:60, 2018.
- [16] Rafael Gómez-Bombarelli, Jennifer N Wei, David Duvenaud, José Miguel Hernández-Lobato, Benjamín Sánchez-Lengeling, Dennis Sheberla, Jorge Aguilera-Iparraguirre, Timothy D Hirzel, Ryan P Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven continuous representation of molecules. *ACS central science*, 4(2):268, 2018.