

Deeply learning molecular structure-property relationships using graph attention neural network

Seongok Ryu,¹ Jaechang Lim,¹ and Woo Youn Kim^{1,2*}

¹Department of Chemistry, KAIST, 291 Daehak-ro, Yuseong-gu, Daejeon 34141, Republic of Korea

²KI for Artificial Intelligence, KAIST, 291 Daehak-ro, Yuseong-gu, Daejeon 34141, Republic of Korea

Correspondence and requests for materials should be addressed to W.Y.K. (email : wooyoun@kaist.ac.kr).

Abstract

Molecular structure-property relationships are the key to molecular design for materials and drug discovery. The rise of deep learning offers a new viable solution to elucidate the structure-property relationships directly from chemical data. Here we show that graph attention networks can greatly improve performance of the deep learning for chemistry. The attention mechanism enables to distinguish atoms in different environments and thus to extract important structural features determining target properties. We demonstrated that our model can detect appropriate features for molecular polarity, solubility, and energy. Interestingly, it identified two distinct parts of molecules as essential structural features for high photovoltaic efficiency, each of which coincided with the area of donor and acceptor orbitals in charge-transfer excitations, respectively. As a result, it could accurately predict molecular properties. Moreover, the resultant latent space was well-organized such that molecules with similar properties were closely located, which is critical for successful molecular design.

Introduction

Elucidating molecular structure-property relationships plays a pivotal role for rational molecular design. However, the structure-property relationships are usually unknown, so molecular engineering based on physical principles and heuristic rules is inevitable with many trials and errors. As a more systematic approach, a computer-aided molecular design has attracted great attention especially in drug and materials discovery^{1,2}. Promising molecules with desired properties are first selected through high-throughput virtual screening of a large set of molecules using computational chemistry before experiments³⁻⁵. In this procedure, expensive computational costs are required for reliable calculations. Practical approximations adopted for efficient screening provoke undesirable errors, giving rise to less reliable results.

The rise of deep learning (DL) techniques⁶ is expected to open a new paradigm for efficient molecular design. Unlike traditional computational methods based on physical principles, the DL can find out structure-property relationships directly from chemical data and apply them in a novel way to molecular design. For example, supervised learning methods have been widely used to learn and predict molecular energetics⁷⁻¹⁰, toxicity¹¹⁻¹⁴, and drug efficacy¹⁵⁻¹⁷. A class of unsupervised learning in particular with novel generative models has been utilized for *de novo* molecular design¹⁸⁻³¹. Reinforcement learning techniques have also been applied to planning synthetic routes and designing drug molecules³⁰⁻³⁴. As such, the DL, although at an early stage, is rapidly spreading to various chemical fields.

The key to the success of the DL in chemistry is how accurately elucidating the structure-property relationships from existing data. That is equivalent to make a DL model to best approximate a function F in $y = F(x)$ where y and x denote molecular properties and structures, respectively. In principle, it is possible if a vast amount of high quality data are available because the DL is known as a universal approximation kernel^{35,36}. However, in practice, the lack of chemical data limits its capability. Thus, it is essential to develop a high performance DL model specialized for chemical problems, as can be seen from the great success of convolutional neural networks (CNN)^{37,38} and recurrent neural networks (RNN)³⁹⁻⁴¹ in the vision recognition and the natural language processing, respectively. Here the high performance DL model means that it can extract important structural features determining a target property from the limited data. In chemistry, both the CNN and RNN have been used by representing molecules with SMILES¹⁸⁻²³ and molecular fingerprints¹¹⁻¹³, because those models are readily available. However, the SMILES and fingerprints are too simplified to deliver the

topological information of molecular structures, resulting in relatively low learning accuracies. Schütt *et al* proposed namely the deep tensor neural network model and achieved a high accuracy for molecular energetics by exploiting 3D molecular structures⁹. Unfortunately, most chemical data provides simplified molecular structures such as the SMILES, fingerprints, and molecular graphs. Moreover, calculations of 3D molecular structures from them are very demanding.

In this aspect, the molecular graph representation would be the best compromise; it describes atoms and bonds in a molecule as vertices and edges, respectively. Molecular graphs intuitively and concisely express molecules with 2D topological information. Hence, it is widely adopted in chemical education as well as chemical informatics. Indeed, there have been efforts to develop DL models based on the molecular graphs. The graph convolutional network (GCN) as an extension of the CNN was proposed to deal with graph structures^{42,43}. The GCN benefits from the advantage of the CNN architecture; a high accuracy with a relatively low computational cost due to less parameters compared to a fully connected multi-layer perceptron model. It can identify important atom features determining molecular properties by analyzing relations between neighboring atoms. Weave model as a variant of the GCN considers not only atom features but also bond features¹². For further generalization, message passing neural network divides this procedure into two steps; it first extracts graph structural features and then relates them with target properties⁸.

Despite the aforementioned advantages, we doubt that those models are still missing an important factor for better structure-property relationships. Molecules are not just a simple collection of atoms. Same atoms can play different roles in determining molecular properties depending on their local chemical environments. For instance, carbon atoms of aromatic rings, aliphatic chains, and carbonyl groups have different characters. Chemists can identify functional groups related to molecular properties. Polar and nonpolar groups are such examples for molecular polarity and solubility. Therefore, it is critical to correctly identify molecular substructures determining a target property to obtain better structure-property relationships. However, the previous graph neural network models apply identical convolution weights to all atoms and bonds. In other words, they treat all atoms and bonds with equal importance regardless of their chemical environments.

To improve performance of the GCN and its variants, an obvious strategy is to apply the graph convolution to individual atoms with different weights depending on their chemical environments. To this end, one may add attention coefficients to atoms which are adaptive to the chemical environments. The resultant neural network optimizes both the convolution weights and the attention coefficients to give better structure-

property relationships. The so-called graph attention network (GAT) has been originally developed for a network analysis in computer science⁴⁴. In chemistry, Shang et al. first adopted the attention mechanism for prediction of molecular properties⁴⁵. It should be noted that Shang et al. described molecules as an assembly of chemical bonds. Thus, they applied the attention mechanism to the chemical bonds and then used the same bond features across all molecules. For instance, all C=O bonds share the same bond features. However, the C=O of a carboxyl group is chemically different from that of an ester group (e.g., bond length and strength), meaning that the bond features must depend on chemical environments as well. Therefore, the chemical bonds would not be appropriate building blocks for molecules.

In this regard, we propose to apply the attention mechanism to atoms instead of chemical bonds. The attention mechanism should differentiate atoms being in different chemical environments by considering pairwise interactions between neighboring atoms. The pairwise interactions in this process seem similar to the atom pair concept in Shang's model. The main difference from Shang's model is that our model may have different interactions even for identical atom pairs if they have different atom features due to different chemical environments. Therefore, it is more flexible to elucidate the structure-property relationships. As a cost for the flexibility, attention coefficients for the pairwise interactions should be trained for all atom pairs independently. Here we focus on the role of the attention mechanism to elucidate the structure-property relationships. Especially, we show that the GAT is able to identify important functional groups which are directly related to target properties. In addition, it helps chemically rationalize the important structural features by mapping them on molecular graphs. This is important because scientific interpretation of a result is often more valuable than the result itself. We show that not so surprisingly the GAT can recognize polar and nonpolar functional groups as important features for molecular polarity and solubility. In addition, we will show that as critical structural features for highly efficient photovoltaic molecules the GAT can precisely identify molecular regions to which donor and acceptor orbitals in charge-transfer excitations are distributed, respectively. Apparently, it is not a trivial task even to experts without any information on the electronic structure of molecules.

Experimental Setup

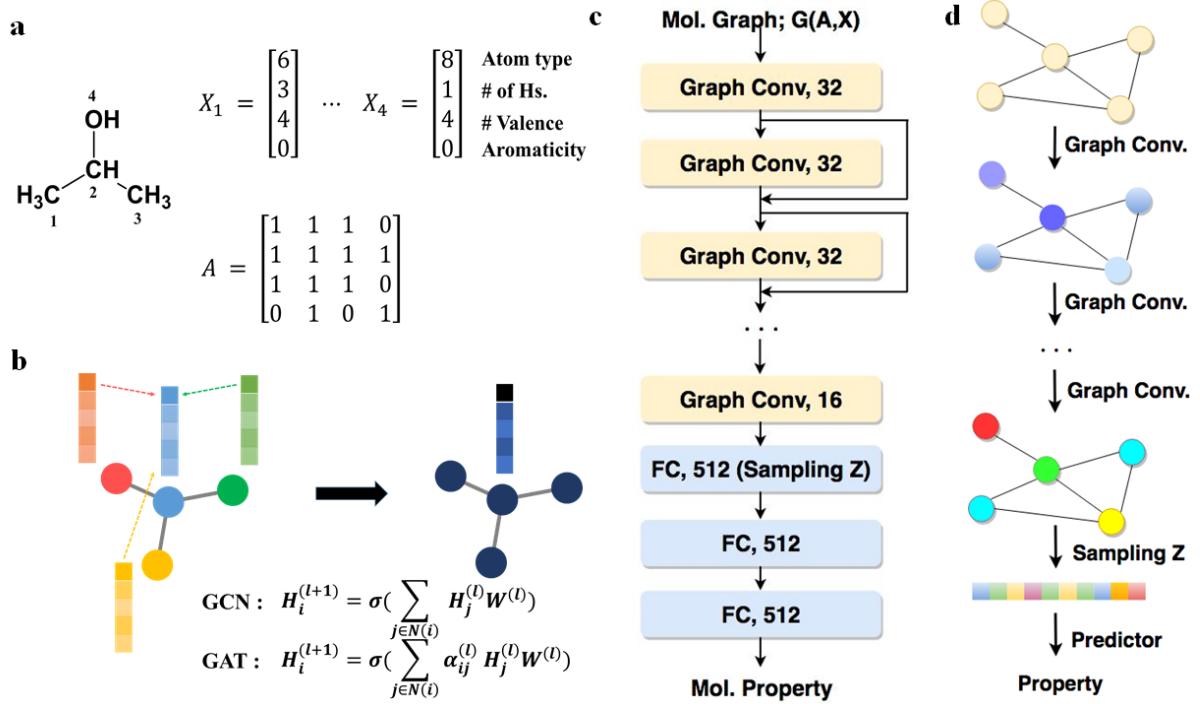


Figure 1. (a) Graph representation, $G(A, X)$, of a 2-propanol. The i -th atom vector X_i contains the initial atom features (atom type, the number of hydrogens attached, the number of valence electrons, and aromaticity) and the matrix A represents the connectivity between atoms including itself. (b) Schematic description of updating a hidden node state, $H_i^{(l)}$, in a graph convolution layer with/without the attention mechanism (c) Architecture of our GAT and GCN in this work. It is composed of six convolution layers and three fully connected layers. (d) Overall procedure of updating atom features and obtaining a target property processed by the neural network presented in c.

Graph convolution network and attention mechanism

Figure 1 shows the GCN and GAT implemented in this work. The GCN updates each state of the $l+1^{\text{th}}$ layer, $H^{(l+1)}$, from those of the l -th layer, $H^{(l)}$, as follows^{42,43}.

$$H^{(l+1)} = \sigma(AH^{(l)}W^{(l)}) \quad (1)$$

where σ , A , and $W^{(l)}$ denote an activation function, an adjacency matrix combined with an identity matrix

(Figure 1a), and the convolution weights of the l -th layer. Suppose that as shown in Figure 1a, the atom 2 has three adjacent atoms 1, 3 and 4. In this case, equation (1) can be written as

$$H_2^{(l+1)} = \sigma(H_2^{(l)}W^{(l)} + H_1^{(l)}W^{(l)} + H_3^{(l)}W^{(l)} + H_4^{(l)}W^{(l)}) \quad (2)$$

which means that all atom features of the current layer are first summed up with the convolution weights and then subjected to the activation function to update the corresponding hidden state in the next layer. A single graph convolution updates atom features only from its adjacent atoms as depicted in Figure 1b. n -times consecutive operations as shown in Figure 1c can update the atom features from neighboring atoms with the n -th graph distance. It is worth to note that the convolution weights in a given layer are identical for all atoms. Therefore, the GCN has a limitation to reflect different local chemical environments of individual atoms in the convolution process.

The attention mechanism enables us to incorporate such a local chemical environment in the graph convolution. The hidden state in equation (2) is updated with different extents of attention to the neighboring atoms, resulting in

$$H_2^{(l+1)} = \sigma(\alpha_{22}H_2^{(l)}W^{(l)} + \alpha_{21}H_1^{(l)}W^{(l)} + \alpha_{23}H_3^{(l)}W^{(l)} + \alpha_{24}H_4^{(l)}W^{(l)}) \quad (3)$$

where $\alpha_{ij}^{(l)}$ denotes an attention coefficient between the i -th and j -th atoms which is adaptive to local environments. The attention coefficient should be determined by considering the interaction strength between the two adjacent atoms i and j to best predict a target property. Thus, we modified the attention coefficient from the original GAT to incorporate the interaction between neighboring atoms as follows:

$$\alpha_{ij}^{(l)} = \sigma\left((H_i^{(l)}W^{(l)})C^{(l)}(H_j^{(l)}W^{(l)})^T\right), \quad (4)$$

where $C^{(l)}$ is a coupling matrix. Note that the coupling matrix is analogous to the dictionaries containing pairwise interactions in Shang’s model. However, the coupling matrix is determined for every individual atom pair forming chemical bonds in a given molecule and thus can be different even in the same atom pair depending on its environments. We used $ReLU(x) = \max(0, x)$ as the activation function in equation (2) and (3) and $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ for equation (4).

Neural network architecture

For comparison, the same architecture shown in Figure 1c is used for both the GCN and GAT. The former updates atom features with the vanilla graph convolution, while the latter adopts a graph convolution with the attention mechanism. We also use a skip-connection, $y = F(x) + x$, at every graph convolution/attention layer to avoid the vanishing gradient problem. Therefore, we are able to repeat the graph convolution many times. In the present work, we use six graph convolution/attention layers, and three fully

connected layers are followed to obtain a target property. Atom feature vectors and latent vectors for further analysis are sampled from the last convolution layer and the first fully-connected layer, respectively. Figure 1d visualizes the overall procedure of updating atom features and obtaining a target property processed by the graph neural network presented in Figure 1c.

Results and Discussion

Prediction results of various molecular properties

Model	logP	TPSA	Atomization energy (kcal/mol)	PVE (%)
GAT	0.019	0.088	4.12	0.63
GCN	0.073	0.75	6.09	0.89
Previous works	Graph - 0.05 ¹⁸ SMILES - 0.13 ¹⁸	-	-	Graph -1.43 ¹¹

Table 1. Mean absolute error of the prediction results of molecular properties. The best ones are shown in bold.

To evaluate the relative performance of the GCN and GAT, we trained each model with various molecular properties. Table 1 shows the prediction results of the partition coefficient (logP), the topological polar surface area (TPSA), the atomization energy, and the photovoltaic efficiency (PVE) for each test set in terms of mean absolute error. The GAT outperformed the GCN for all the cases. In particular, it was about 3.8 and 8.5 times more accurate for the logP and TPSA, respectively. This is because the logP and TPSA are explicitly related to molecular structures such as polar functional groups. In contrast, the atomization energy and photovoltaic efficiency are determined by electronic structures which can be obtained from quantum chemical calculations, but we did not use any information on the electronic structures. Therefore, the improvement of the GAT over the GCN is relatively small for the two electronic properties. In what follows, we explain the reason for the outperformance of the GAT by investigating the role of the attention mechanism in elucidating the structure-property relationships.

Atom features extracted by GAT and GCN

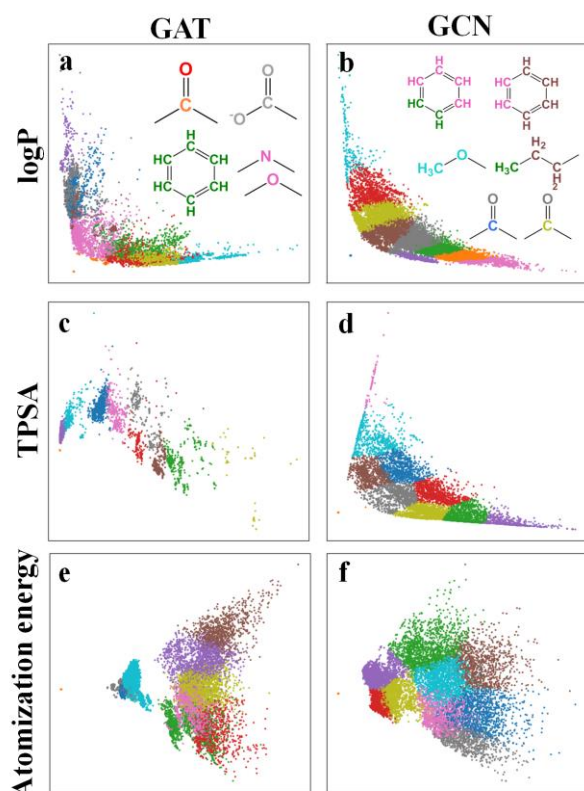


Figure 2. 2D plots generated by the principal component analysis of the atom feature vectors obtained from the GAT and GCN for various properties. Each color represents one of 10 groups classified by the k-means clustering of the latent vectors. The same color notation is applied to each atom in molecules as shown in a and b.

The most important role of DL models is to extract important features from data for correct mapping between input and output data (e.g., molecular structures and properties, respectively, in our work). In this process, the DL models may result in a well-structured latent space, meaning that inputs with similar output values are closely located in the latent space. In the natural language processing, for example, similar words are distributed in the same area of the latent space made of word embedding vectors³⁹. In chemistry, similar sub-molecular structures (e.g., functional groups) would contribute to a target property with similar extents and hence have similar feature vectors. Thus, recognizing important substructures associated with the target property is critical to achieve a high accuracy.

In this regard, we examined atom features extracted by the GAT and GCN for various target properties. To construct a relevant latent space, we used atom feature vectors sampled from the last graph convolution layer for each property as depicted in Figure 1c and 1d. Figure 2 shows the 2D plots of the high-

dimensional latent space generated by the principal component analysis of the atom feature vectors. Each color represents one of 10 groups classified by the k-means clustering⁴⁶ of the latent space. In other words, all atoms have been classified into the 10 groups depending on their feature vectors associated with a given target property. The same color notation is applied to each atom of molecules as shown in Figure 2a and 2b. We expected that atoms in a similar chemical environment have the same color because they would have similar feature vectors. Indeed, it was true for the case of the GAT. We noted that atoms with the same local environment have the same color across all molecules as shown in Figure 2a. For instance, all carbon atoms of aromatic rings are green, while those of carbonyl groups are orange. Similarly, ether oxygens and carbonyl oxygens are pink and red, respectively. However, the GCN assigned different colors to atoms in the same chemical environment or same colors to atoms in different chemical environments as depicted in Figure 2b; for example, carbon atoms of even the same aromatic ring have different colors or carbon atoms of a methyl group and an aromatic group have the same green color. A similar tendency was also found for the other properties.

	GAT	GCN
<div>logP</div> <div>TPSA</div> <div>Atomization energy (Hartree)</div>	<div> A (6.09) </div> <div>(6.09)</div>	<div> A (6.09) </div> <div>(6.09)</div>
	<div> B (-3.18) </div> <div>(-3.16)</div>	<div> B (-3.19) </div> <div>(-3.19)</div>
	<div> C (0.0) </div> <div>(-0.01)</div>	<div> C (0.2) </div> <div>(0.2)</div>
	<div> D (141.5) </div> <div>(141.7)</div>	<div> D (141.1) </div> <div>(141.1)</div>
	<div> E (-3.651) </div> <div>(-3.654)</div>	<div> E (-3.659) </div> <div>(-3.659)</div>
	<div> F (-1.899) </div> <div>(-1.901)</div>	<div> F (-1.93) </div> <div>(-1.93)</div>

Figure 3. Representative molecules for each property. Each atom is colored with the color notation for each property obtained in Figure 2. The numbers in parentheses denote true values (leftmost column) and predicted values (below each molecule) by the GAT and GCN, respectively.

To further investigate the dependence of atom features on local chemical environments, we analyzed a few representative molecules for the logP, TPSA, and atomization energy. We colored each atom of those molecules with the color notation obtained for each property in Figure 2. For the molecule A having a high logP value, the GAT assigned a consistent color to atoms in similar chemical environments. The oxygen atoms of the two ether groups are pink. All aromatic ring carbons are green, while the other carbon atoms have different colors depending their surroundings. The two carbonyl carbons share the orange color. Interestingly, the GAT was also able to differentiate carbon atoms of the terminal alkyl groups according to the neighboring atom type; alkyl carbons attached to another carbon are purple, while those adjacent to oxygen are cyan. For the same molecule, however, the GCN shows an inconsistent color mapping to atoms in similar environments. For the molecule B having a low logP value, the GAT shows the same color mapping with that of A. For instance, the carbon atoms of the aromatic rings in A and B have the same color as green, though the two molecules have completely different logP values. The reason for why the GAT assigned different colors to atoms according to those functional groups is because the logP strongly depends on them, indicating that the attention mechanism was able to well characterize the structure-property relationship.

As for the TPSA, we can analyze atom features in a similar fashion. In this case, more polar functional groups may induce larger TPSA values. Therefore, we expect that the GAT is able to identify polar and nonpolar functional groups as important structural features. For the molecule C having the zero TPSA value, the GAT assigned the same purple color to all atoms, whereas the GCN gave three different colors. The GAT result is more reasonable because the molecule C does not involve polar functional groups and so all carbon atoms are in a uniform environment in terms of polarity. On the contrary, the molecule D has a very large TPSA value. The GAT still assigned the same purple color to aromatic carbon atoms because they are nonpolar as was in the molecule C. On the other hand, the carbon atoms of carboxyl and carbonyl groups in D have not only different colors with that of aromatic carbons but also different from one another because the relative polarity of the two polar groups is different. In contrast, the GCN gave rise to an inconsistent color mapping. It is also noticeable that the GAT gave different colors to carbon atoms attached to polar functional groups (gray carbons connected to the nitro and amide groups).

In the case of the atomization energy, a well-trained neural network should differentiate atoms depending on their relative chemical bond strengths. The chemical bond strengths depend on local chemical environments as well as atom types. For example, the strength of the C-C bond is different in all cases of a

strained ring, a unstrained aliphatic chain, and an aromatic ring. Again, the GAT and GCN produced considerably different atom features for the example molecules E and F shown in Figure 3. It should be noted that we included hydrogen atoms explicitly in the atomization energy learning. Therefore, each hydrogen atom has its own color, while the previous two examples used implicit hydrogen atoms indicated by the number of hydrogen atoms attached to each atom as input data. The GAT identified hydrogen atoms of $\text{-CH}_2\text{-}$, -CH_3 , and acidic hydrogen (NH_2^+) differently. In contrast, the GCN assigned the same color to those hydrogen atoms. Unlike the GAT, the GCN even assigned the same color to the oxygen atom of the carboxylic group in F with that of the hydrogen atoms probably because they are all terminal atoms. Overall, the GAT distinguished atoms according to their chemical environments more delicately.

The three properties are directly related to molecular structures. Thus, it would be easy for the GAT to identify specific functional groups relevant to the target properties. Undergraduate students majoring chemistry may also recognize those functional groups readily. However, identifying key molecular substructures related to the PVE of a molecule would be challenging even to experts without any information on electronic structures. As the last example, we examined atom features of photovoltaic molecules obtained from the GAT and GCN. In this case, we used only two colors (red and blue) to simplify the visual analysis. Figure 4 shows the results of two representative molecules; the molecule G has a very high PVE value, while the molecule H has a very low PVE value. Interestingly, both the GAT and GCN categorized atoms of G into two groups with red and blue colors, respectively, while they assigned only a single color to all atoms of H. Especially, we noted that the GAT divided the molecule G into two parts denoted by the two colors. To check out if these color mappings are closely related to key structural factors for the corresponding PVE values or not, we plotted the highest occupied molecular orbital (HOMO) and the lowest unoccupied molecular orbital (LUMO) obtained from density functional theory calculations for G and H in Figure 4. Surprisingly, the red and blue regions of G coincided with the areas to which the HOMO and the LUMO are distributed, respectively. In the case of H, all atoms had the same color, and indeed both the HOMO and the LUMO of the molecule are delocalized over the entire region. The results can be rationalized as follows. To be a good photovoltaic molecule, an excited electron and the corresponding hole should be separated spatially to prevent from an easy recombination as well as to readily separate the electron-hole pair for energy harvesting. In this aspect, it is remarkable that the GAT was able to find out the important structural features determining the PVE just from molecular graphs.

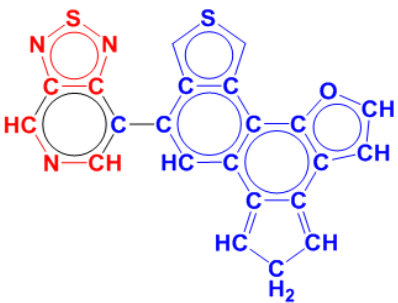
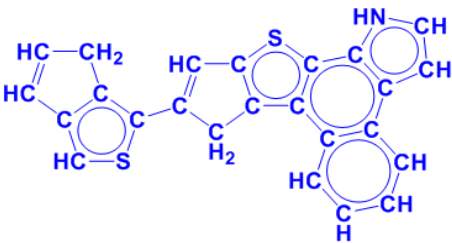
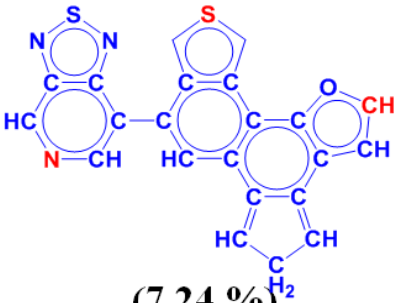
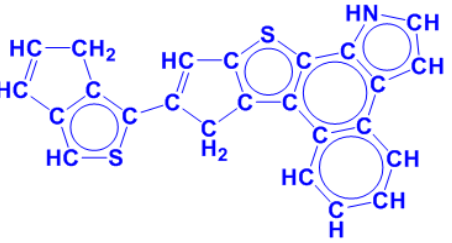
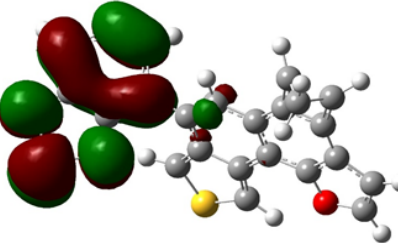
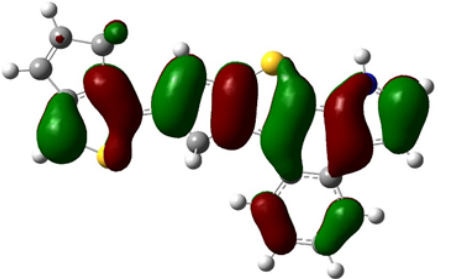
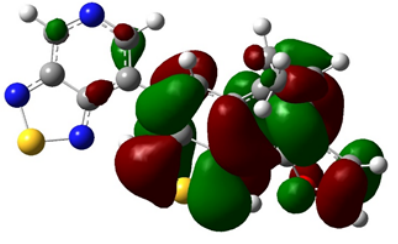
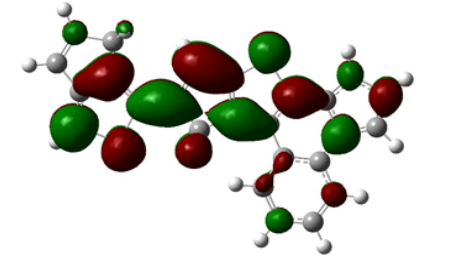
	G (PVE = 7.27 %)	H (PVE = 1.30 %)
GAT	 (8.12 %)	 (1.29 %)
GCN	 (7.24 %)	 (1.82 %)
HOMO		
LUMO		

Figure 4. Two example molecules G and H with high and low PVE values, respectively. Each atom is colored with the binary color notation obtained from the atom feature vector analysis as was done in Figure 2. The numbers in parentheses denote the true PVE values of the molecules.

Latent space implying the structure-property relationships

The above examples for the four molecular properties manifest that the attention mechanism greatly helps the graph neural networks capture key structural features closely related to target properties. It becomes

possible because of using adaptive attention coefficients for individual atoms. As a result, the graph convolution with the attention mechanism may lead to a desirable mapping between molecular structures and properties. Deeply learning the structure-property relationships is the key to successful molecular design. Previous works have shown how DL models utilize such relationships for *de novo* molecular design. A common idea is to use a latent space, which is supposed to contain the structure-property mapping, obtained from the DL models. In a well-trained latent space, molecules with similar properties are closely located. Thus, new molecules with desired properties can be generated by exploring the latent space. Along the same strategy, we examined the latent space of the GAT for the four properties. We used the latent vectors obtained from the first fully connected layer for each property as shown in Figure 1c and 1d.

As the first example, we examined the three closest and the three farthest molecules from a given reference molecule in the PVE latent space obtained from the GAT. The distance between two molecules was measured by the L2-norm distance, $d_{ij} = |z_i - z_j|$, between the corresponding latent vectors z_i and z_j . Figure 5 shows the reference molecule with the largest PVE value in the data set and the resulting six molecules with their distances from the reference and PVE values. In the case of the closest molecules, all of them have very high PVE values. In addition, they have two distinct atoms spatially well-separated by the red and blue colors. On the contrary, the three farthest molecules have very low PVE values and only blue color atoms except for a single red atom. These results indicate that the latent space of the GAT implies an accurate structure-property mapping. Furthermore, one can use the latent space to design new molecules with a high PVE value for example using generative models proposed in the literature^{18,20,21,23,29,31}. To our best knowledge, however, such a generative model with the attention mechanism is not available at present.

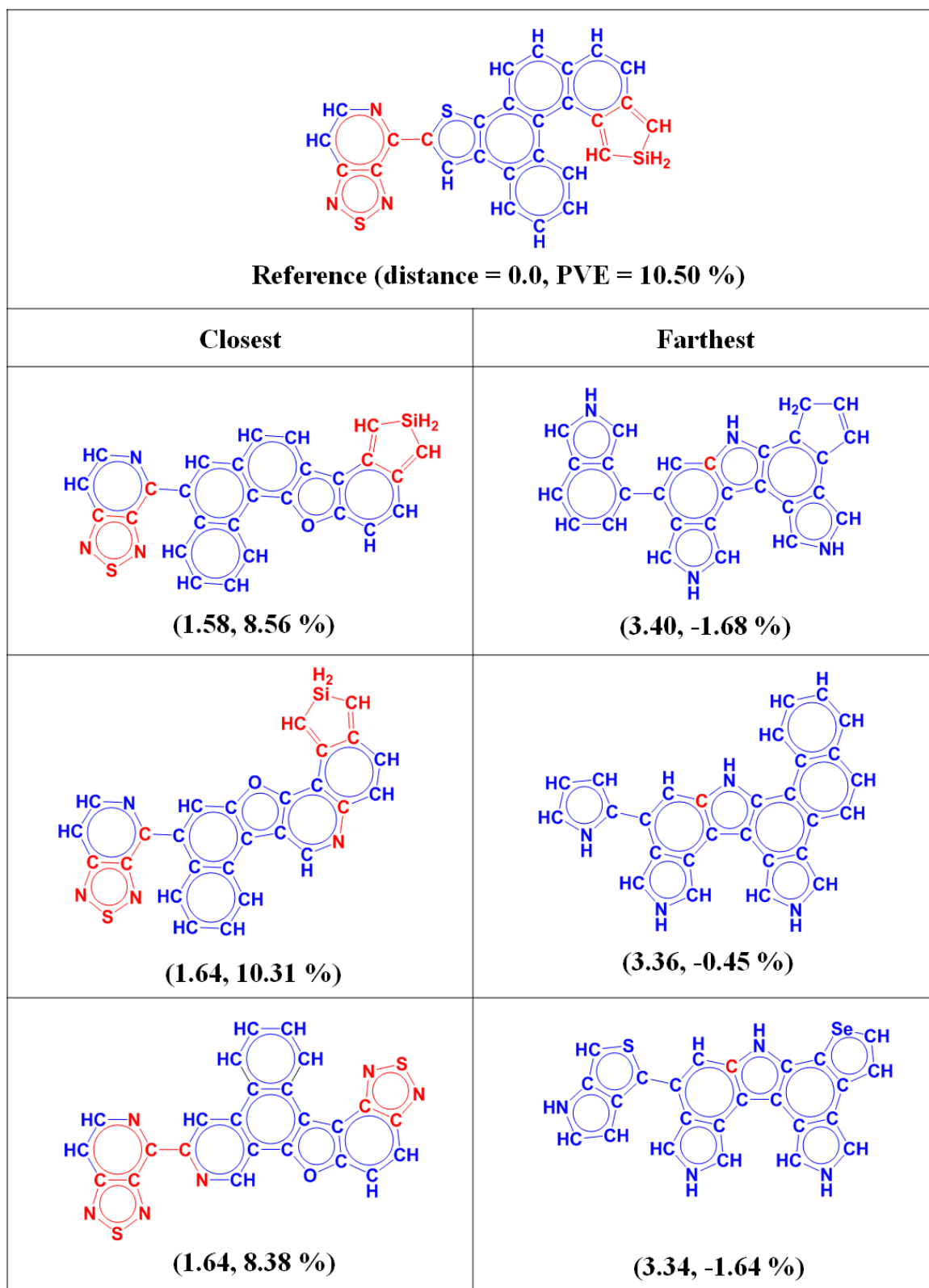


Figure 5. One reference molecule and the three closest and the three farthest molecules from the reference in the PVE latent space obtained from the GAT. Each atom is colored with the binary color notation obtained from the atom feature vector analysis as was done in Figure 2. The numbers in parentheses denote the L2-norm distances of the molecules from the reference and their true PVE values, respectively.

We expect similar results for the other properties. To statistically analyze the results, we sorted 5,000 molecules randomly chosen from each test set in the ascending order of their properties and measured the L2-norm distance between them. Figure 6 shows the 2D plot for the distance mapping of each case. The distance between two molecules labeled by the row and column indices is denoted by the scale bar next to each figure. The dark blue in the diagonal elements means the distance of a molecule from itself, while the dark red denotes the farthest distance. The GAT shows a gradual color change from the dark blue to the dark red, resulting in unique patterns. Especially, the color change of the TPSA seems very smooth compared to the other two cases. It is consistent to the relatively high accuracy for the TPSA prediction (Table 1). Such a gradual color change means that molecules are well clustered around a molecule with a similar property in the latent space. The GCN also shows a gradual color change at some regions. However, it shows an abrupt change from the dark blue to the white in particular for the logP. This result is a strong evidence indicating that the GAT learns the structure-property relationship better than the GCN does.

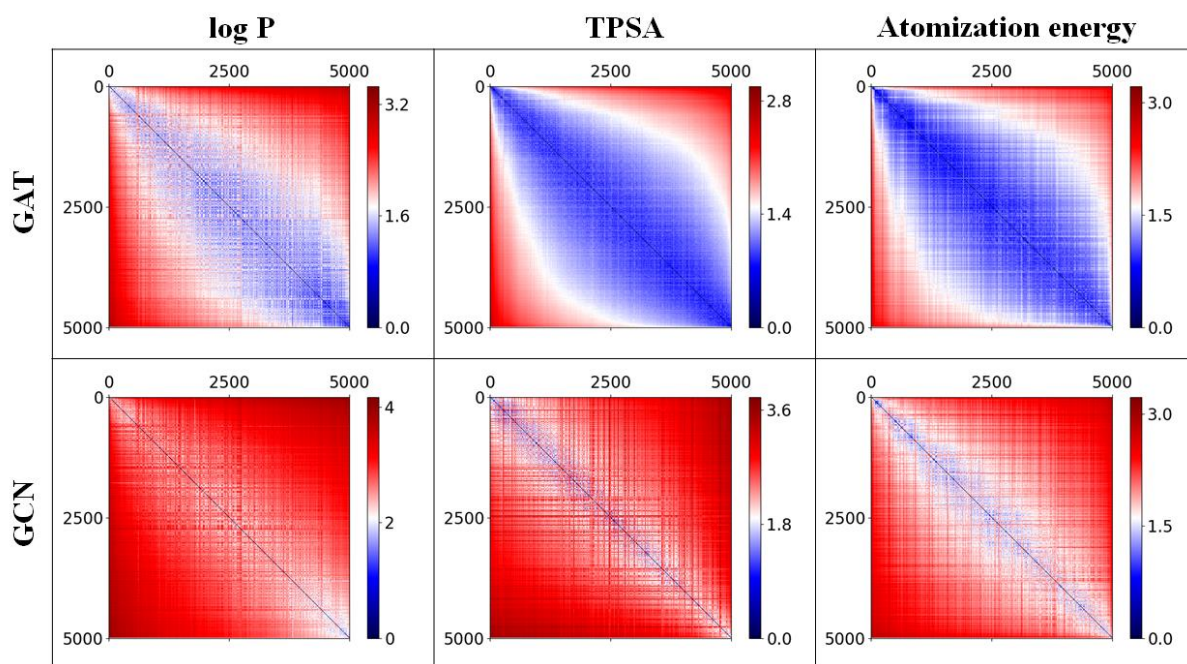


Figure 6. 2D plot of the L2-norm distances between molecules in the latent space associated with each target property. The 5,000 molecules randomly chosen from each test set are listed along the row and column in the ascending order of their properties. The distance between two molecules labeled by the row and column indices is denoted by the scale bar next to each figure.

Conclusions

We have shown that the graph attention network (GAT) proposed here outperforms the graph convolutional network (GCN) in the supervised learning of various molecular properties. The attention mechanism can distinguish atoms in different chemical environments by considering the neighborhood of atoms. Thus, the GAT can extract important structural features determining a target property. We demonstrated by analyzing atom features that the GAT was able to identify polar and nonpolar functional groups of molecules as key structural features for the logP and the TPSA. It could also delicately differentiate atoms in different environments for the atomization energy learning. More interestingly, the GAT identified two distinct parts of molecules as important structural features for high photovoltaic efficiency, which turned out that each of them corresponds to the areas around which donor and acceptor orbitals of the molecules in charge-transfer excitations reside, respectively. Evidently the GAT elucidates the structure-property relationships from chemical data better than the GCN does. As a result, the GAT produced well-organized latent spaces such that molecules with similar properties were closely located.

Such a high performance of the GAT offers various application possibilities. In particular, accurate learning of the structure-property relationship is essential to design new molecules with desired properties using molecular generative models. For instance, Gómez-Bombarelli and coworkers showed that a variational autoencoder can generate new molecules with a target property through the gradient-based optimization process in a latent space¹⁸. We also demonstrated that a conditional variational autoencoder can design molecules with simultaneous control of multiple target molecular properties for drug discovery by embedding them directly in latent vectors²³. Segler et al. and Gupta et al. designed molecules with specific biological activities using a natural language processing model combined with transfer learning^{20,21}. Jaques et al., Olivecrona et al., and Guimaraes et al. proposed methods which finely tune a pretrained generative model using reinforcement learning to generate molecules with certain desirable properties^{22,30,31}.

However, these models employed the SMILES representation of molecular structures, so they may have a low rate of valid molecules and no topological information. Recently, there is a growing interest in developing graph-based generative models²⁵⁻²⁹. The early stage works showed promising results with high rates of valid and novel molecules, which is important to explore an extended chemical space. We expect that the GAT will substantially improve such generative models for *de novo* molecular design via a well-trained structure-property

relationship. However, generative models with the attention mechanism have yet to be developed. In addition, the attention mechanism can be applied to other graph neural networks for chemical applications. Consequently, we believe that the graph attention network has a broad impact on molecular design for materials and drug discovery.

Methods

Training, validation and test conditions

Data sets and number of molecules for training, test and validation of each molecular property are indicated in Table 2. We use learning rate of 0.001 with decay rate 0.95 for all cases.

	logP	TPSA	Atomization energy	PVE
Dataset	ZINC ⁴⁷		QM9 ⁴⁸	Harvard Clean Energy Project ¹¹
Training	400,000		96,000	19,200
Validation	50,000		24,000	4,800
Test	50,000		13,885	5,978

Implementation detail

The input node states contain raw atom features including atom type, number of hydrogens attached, implicit valence number, and aromaticity. They were represented with one-hot encoding. Hydrogen atoms were represented explicitly for atomization energy learning and implicitly for the other cases. We obtained the raw atom features and the adjacency matrices of molecules from the open-source python toolkit, RDKit⁴⁹. We implemented both GAT and GCN using TensorFlow. We uploaded our code, supplementary figures, and the pre-trained models on the github (<https://github.com/SeongokRyu/Molecular-GAT>).

Electronic structure calculation

We obtained the molecular orbitals in Figure 4 from density functional calculations as implemented in GAUSSIAN09⁵⁰. All calculations were performed with B3LYP⁵¹ and 6-31G(d) basis-set.

References

1. McCammon, J. A. Computer-aided molecular design. *Science* **238**, 486-491 (1987).
2. Todd, M. H. Computer-aided organic synthesis. *Chem. Soc. Rev.* **34**, 247-266 (2005).
3. Shoichet, B. K. Virtual screening of chemical libraries. *Nature* **432**, 862-865 (2004).
4. Curtarolo, S. *et al.* The high-throughput highway to computational materials design. *Nat. Mater.* **12**, 191–201 (2013)
5. Pyzer-Knapp, E. O., Suh, C., Gómez-Bombarelli, R., Aguilera-Iparraguirre, J., & Aspuru-Guzik, A. What is high-throughput virtual screening? A perspective from organic materials discovery. *Annu. Rev. Mater. Res.* **45**, 195-216 (2015).
6. LeCun, Y., Bengio, Y., & Hinton, G. Deep learning. *Nature*, **521**, 436-444 (2015).
7. Rupp, M., Tkatchenko, A., Müller, K. R., & Von Lilienfeld, O. A. Fast and accurate modeling of molecular atomization energies with machine learning. *Phys. Rev. Lett.* **108**, 058301-058306 (2012).
8. Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., & Dahl, G. E. Neural message passing for quantum chemistry. *arXiv:1704.01212*. (2017).
9. Schütt, K. T., Arbabzadah, F., Chmiela, S., Müller, K. R., & Tkatchenko, A. Quantum-chemical insights from deep tensor neural networks. *Nat. Commun.* **8**, 13890 (2017).
10. Faber, F. A. *et al.* Prediction errors of molecular machine learning models lower than hybrid DFT error. *J. Chem. Theo. Comput.* **13**, 5255-5264 (2017).
11. Duvenaud, D. K., Maclaurin, D., Iparraguirre, J., Bombarelli, R., Hirzel, T., Aspuru-Guzik, A., & Adams, R. P. Convolutional networks on graphs for learning molecular fingerprints. *In Advances in Neural Information Processing Systems*. 2224-2232 (2015).
12. Kearnes, S., McCloskey, K., Berndl, M., Pande, V., & Riley, P. Molecular graph convolutions: moving beyond fingerprints. *J. Comput.-Aided Mol. Des.* **30**, 595-608 (2016).
13. Wu, Z. *et al.* MoleculeNet: a benchmark for molecular machine learning. *Chem. Sci.* **9**, 513-530 (2018).

14. Mayr, A., Klambauer, G., Unterthiner, T., & Hochreiter, S. DeepTox: toxicity prediction using deep learning. *Front. Environ. Sci.* **3**, 80 (2016).
15. Gomes, J., Ramsundar, B., Feinberg, E. N., & Pande, V. S. Atomic convolutional networks for predicting protein-ligand binding affinity. *arXiv:1703.10603*. (2017).
16. Öztürk, H., Ozkirimli, E., & Özgür, A. DeepDTA: Deep Drug-Target Binding Affinity Prediction. *arXiv:1801.10193*. (2018).
17. Jiménez Luna, J., Skalic, M., Martinez-Rosell, G., & De Fabritiis, G. K-DEEP: Protein-ligand absolute binding affinity prediction via 3D-convolutional neural networks. *J. Chem. Info. Model.* **58**, 287-296. (2018).
18. Gómez-Bombarelli, R. *et al.* A. Automatic chemical design using a data-driven continuous representation of molecules. *ACS Central Science.* **4**, 268-276 (2016).
19. Kusner, M. J., Paige, B., & Hernández-Lobato, J. M. Grammar variational autoencoder. *arXiv:1703.01925*. (2017).
20. Segler, M. H., Kogej, T., Tyrchan, C., & Waller, M. P. Generating focused molecule libraries for drug discovery with recurrent neural networks. *ACS Central Science.* **4**, 120-131 (2017).
21. Gupta, A., Müller, A. T., Huisman, B. J., Fuchs, J. A., Schneider, P., & Schneider, G. Generative recurrent networks for de novo drug design. *Molecular informatics* **37**, 1700111 (2018).
22. Jaques, N., Gu, S., Bahdanau, D., Hernández-Lobato, J. M., Turner, R. E., & Eck, D. Sequence tutor: Conservative fine-tuning of sequence generation models with KL-control. *arXiv:1611.02796*. (2016).
23. Lim, J., Ryu, S., Kim, J., & Kim, W. Y. Molecular generative model based on conditional variational autoencoder for de novo molecular design. *under revision*. (2018)
24. Müller, A. T., Hiss, J. A., & Schneider, G. Recurrent Neural Network Model for Constructive Peptide Design. *J. Chem. Info. Model.* **58**, 472-479 (2018).
25. Li, Y., Vinyals, O., Dyer, C., Pascanu, R., & Battaglia, P. Learning deep generative models of graphs. *arXiv:1803.03324*. (2018).
26. You, J., Ying, R., Ren, X., Hamilton, W. L., & Leskovec, J. GraphRNN: A Deep Generative Model for Graphs. *arXiv:1802.08773*. (2018).
27. Jin, W., Barzilay, R., & Jaakkola, T. Junction Tree Variational Autoencoder for Molecular Graph Generation. *arXiv:1802.04364*. (2018).

28. Simonovsky, M., & Komodakis, N. GraphVAE: Towards Generation of Small Graphs Using Variational Autoencoders. *arXiv:1802.03480*. (2018).
29. Li, Y., Zhang, L., & Liu, Z. Multi-Objective De Novo Drug Design with Conditional Graph Generative Model. *arXiv:1801.07299*. (2018).
30. Olivecrona, M., Blaschke, T., Engkvist, O., & Chen, H. Molecular de-novo design through deep reinforcement learning. *Journal of Cheminformatics* **9**, 48 (2017).
31. Guimaraes, G. L., Sanchez-Lengeling, B., Farias, P. L. C., & Aspuru-Guzik, A. Objective-reinforced generative adversarial networks (ORGAN) for sequence generation models. *arXiv:1705.10843*. (2017).
32. Zhou, Z., Li, X., & Zare, R. N. Optimizing Chemical Reactions with Deep Reinforcement Learning. *ACS Central Science* **3**, 1337-1344 (2017).
33. Segler, M. H., Preuss, M., & Waller, M. P. Planning chemical syntheses with deep neural networks and symbolic AI. *Nature* **555**, 604-610 (2018).
34. Wei, J. N., Duvenaud, D., & Aspuru-Guzik, A. Neural networks for the prediction of organic chemistry reactions. *ACS Central Science* **2**, 725-732 (2016).
35. Young, M. The Stone-Weierstrass Theorem. (2006)
36. Lin, H. W., Tegmark, M., & Rolnick, D. Why does deep and cheap learning work so well?. *J. Stat. Phys.* **168**, 1223-1247 (2017).
37. Krizhevsky, A., Sutskever, I., & Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*. 1097-1105 (2012).
38. Kim, Y. Convolutional neural networks for sentence classification. *arXiv:1408.5882*. (2014).
39. Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.
40. Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *arXiv:1406.1078*. (2014).

41. Bowman, S. R., Vilnis, L., Vinyals, O., Dai, A. M., Jozefowicz, R., & Bengio, S. Generating sentences from a continuous space. *arXiv:1511.06349*. (2015).
42. Kipf, T. N., & Welling, M. Semi-supervised classification with graph convolutional networks. *arXiv:1609.02907*. (2016).
43. Defferrard, M., Bresson, X., & Vandergheynst, P. Convolutional neural networks on graphs with fast localized spectral filtering. In *Advances in Neural Information Processing Systems*. 3844-3852 (2016).
44. Veličković, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., & Bengio, Y. Graph Attention Networks. *arXiv:1710.10903*. (2017).
45. Shang, C., Liu, Q., Chen, K. S., Sun, J., Lu, J., Yi, J., & Bi, J. Edge attention-based multi-relational graph convolutional networks. *arXiv:1802.04944*. (2018).
46. Arthur, D., & Vassilvitskii, S. k-means++: The advantages of careful seeding. In *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms* 1027-1035. Society for Industrial and Applied Mathematics. (2007)
47. Irwin, J. J., & Shoichet, B. K. ZINC— a free database of commercially available compounds for virtual screening. *J. Chem. Info. Model.* **45**, 177-182 (2005).
48. Ramakrishnan, R., Dral, P. O., Rupp, M., & Von Lilienfeld, O. A. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific Data* **1**, 140022. (2014).
49. Landrum, G. RDKit: Open-source cheminformatics. (2006).
50. Frisch, M. (2013). Gaussian09. <http://www.gaussian.com/>.
51. Stephens, P. J., Devlin, F. J., Chabalowski, C. F. N., & Frisch, M. J. Ab initio calculation of vibrational absorption and circular dichroism spectra using density functional force fields. *J. Phys. Chem.* **98**, 11623-11627 (1994).

Acknowledgements

This work was supported by Basic Science Research Programs through the National Research Foundation of Korea (NRF) funded by the Ministry of Science, ICT and Future Planning (NRF-2017R1E1A1A01078109).

Author contributions

S.R. and W.Y.K. conceived the idea, S.R. did the implementation and run the simulation. All the authors analyzed the results and wrote the manuscript together.

Additional information

Competing financial interests: The authors declare no competing financial interests.