

Evolutionary Deep Fusion Method and Its Application in Chemical Structure Recognition

Xinyan Liang^{ID}, *Student Member, IEEE*, Qian Guo^{ID}, *Student Member, IEEE*, Yuhua Qian^{ID}, *Member, IEEE*,
Weiping Ding^{ID}, *Senior Member, IEEE*, and Qingfu Zhang, *Fellow, IEEE*

Abstract—Feature extraction is a critical issue in many machine learning systems. A number of basic fusion operators have been proposed and studied. This article proposes an evolutionary algorithm, called evolutionary deep fusion method, for searching an optimal combination scheme of different basic fusion operators to fuse multiview features. We apply our proposed method to chemical structure recognition. Our proposed method can directly take images as inputs, and users do not need to transform images to other formats. The experimental results demonstrate that our proposed method can achieve a better performance than those designed by human experts on this real-life problem.

Index Terms—Deep learning, evolutionary algorithms (EAs), molecular structure recognition, multiview fusion.

Manuscript received October 2, 2020; revised January 1, 2021 and February 22, 2021; accepted March 5, 2021. Date of publication March 9, 2021; date of current version October 1, 2021. This work was supported in part by National Key Research and Development Program of China under Grant 2018YFB1004300; in part by the National Natural Science Fund of China under Grant 61672332, Grant 61432011, Grant 61976129, Grant 61976120, and Grant 61502289; in part by the Key Research and Development Program (International Science and Technology Cooperation Project) of Shanxi Province, China, under Grant 201903D421003; in part by the Program for the Young San Jin Scholars of Shanxi under Grant 2016769; in part by the Young Scientists Fund of the National Natural Science Foundation of China under Grant 61802238, Grant 61906115, Grant 61603228, Grant 62006146, and Grant 61906114; in part by the Shanxi Province Science Foundation for Youths under Grant 201901D211169, Grant 201901D211170, and Grant 201901D211171; in part by the Research Project Supported by Shanxi Scholarship Council of China under Grant HGKY2019001; and in part by the Scientific and Technological Innovation Programs of Higher Education Institutions in Shanxi under Grant 2020L0036. (Xinyan Liang and Qian Guo contributed equally to this work.) (Corresponding author: Yuhua Qian.)

Xinyan Liang and Qian Guo are with the Institute of Big Data Science and Industry, Shanxi University, Taiyuan 030006, China (e-mail: liangxinyan48@163.com; czguoqian@163.com).

Yuhua Qian is with the Institute of Big Data Science and Industry, Shanxi University, Taiyuan 030006, China, and also with the Key Laboratory of Computational Intelligence and Chinese Information Processing of Ministry of Education, Shanxi University, Taiyuan 030006, China (e-mail: jinzhengqyh@126.com).

Weiping Ding is with the School of Information Science and Technology, Nantong University, Nantong 226019, China, and also with the Centre for Artificial Intelligence, FEIT, University of Technology Sydney, Ultimo, NSW 2007, Australia (e-mail: dwp9988@163.com).

Qingfu Zhang is with the Department of Computer Science, City University of Hong Kong, Hong Kong, China (e-mail: qingfu.zhang@cityu.edu.hk).

Color versions of one or more figures in this article are available at <https://doi.org/10.1109/TEVC.2021.3064943>.

Digital Object Identifier 10.1109/TEVC.2021.3064943

I. INTRODUCTION

FEATURE extraction is a key in many machine learning systems. A number of deep neural networks (DNNs), such as Inception [1], ResNet [2], and DenseNet [3] have been used for this purpose. Different networks extract features from different views. It is natural to use several neural networks to extract multiview features and then do fuse them [4]. A number of basic fusion operators, such as concatenation [5], elementwise addition [6], elementwise multiplication [7], elementwise max [8], bilinear pooling [9], and tensor-based fusion [10] have been proposed and widely used in machine learning field. To the best of our knowledge, all the existing fusion methods use only one single basic fusion operator, and the features for fusion are manually selected by human experts. This article will investigate how to design an algorithm for searching an optimal combination scheme of different basic fusion operators to fuse multiview features. More specifically, We address the following two issues.

- 1) How to select view features for fusion?
- 2) How to select and use different basic fusion operators for fusing these selected features?

Inspired by the recent success of evolutionary algorithms (EAs) on neural architecture search (NAS), we code a fusion scheme as a chromosome vector which consists of selected features and basic fusion operators. The performance of each fusion scheme can be evaluated by its corresponding deep fusion network. We use an EA for finding an optimal fusion scheme. Our work represents a first attempt to automatically construct an optimal fusion scheme.

We apply our proposed method, named evolutionary deep fusion method (EDF), to chemical structure recognition. In this article, chemical structure recognition is defined as a task for identifying and/or verifying what compounds are in a chemical structure.¹ It can be used in many cheminformatics applications, such as patent search and drug search [11]. Many methods have been developed for this recognition problem (e.g., [12], [13]). Molecules can be naturally represented as graphs [14] and chemical structure recognition can be modeled as a graph search problem [15], [16]. VF2, a widely used molecular graph matching algorithm [17] for this recognition problem is of high time complexity. Some heuristics methods (e.g., [13]) have also been proposed for solving molecular graph search problems. However, all these methods require

¹Chemical structure recognition is to automatically convert images into some special formats in some research papers.

some special formats (e.g., SMILES [18], MOLfile [19]) designed by human experts for chemical structures. Designing these formats and collecting data requires a lot of human labor work, and these human designed formats often have some deficiencies. For example, markush structures cannot be represented in MOL files [20]. These deficiencies could deteriorate the performance of a chemical structure recognition system.

The most commonly used form of compound structures is images. Some Image2Structure tools (e.g., ChemGrapher [21], MolRec [22], Imago [23], OSRA [24], MolVec,² more see [25]) have been developed for automatically converting images into some special formats. However, their performances are not very satisfactory. For example, as reported in [25], the accuracy recognition rates of MolVec 0.9.7, Imago 2.0, and OSRA 2.1 are 66.67%, 40.00%, and 57.78%, respectively, on the JPO dataset. Thus, using these generated special formats, molecular graph matching algorithms may not work very well.

Over the last few years, artificial intelligence techniques, especially deep learning, have been extensively applied in chemistry, such as medical diagnosis [26] and chemical syntheses [27]. Several datasets are available for training neural networks and other machine learning systems. For example, ChEMBL [28] and PubChem [29] contain a large number of chemical structure images, ChEMBL is a large bioactivity dataset and PubChem's BioAssay is a small molecules dataset. For our research purpose, we have also collected a dataset, named ChemBook, which contains only natural compounds. These datasets make it feasible to train a DNN which can directly identify compounds from images of chemical structure.

Effective features are very important for chemical structure recognition [30]. Using different neural networks, we can easily obtain many features for chemical structures from different views and then transform chemical structure recognition into a multiview learning problem.

Our major contributions include:

- 1) we propose a simple yet efficient evolutionary EDF. It is a mix of deep learning, multiview learning and EA. EDF can not only automatically select proper DNNs to extract multiview features and select proper views from a candidate view set, but also find a suitable fusion scheme for different views from a candidate basic fusion operator set;
- 2) we have applied the proposed EDF to the chemical structure recognition problem. The experimental results have demonstrated the effectiveness of EDF. EDF has been successfully integrated into a patent data analysis platform at Shanxi University.

The remainder of this article is organized as follows. In Section II, we review the related work of multiview learning and neural architecture search (NAS). In Section III, we present the details of the EDF method. In Section IV, the performance of the EDF is evaluated on three chemical structure recognition datasets. Finally, we draw conclusions in Section V.

II. RELATED WORK

In this section, we give a review of multiview learning and neural architecture search for DNNs.

A. Multiview Learning

Multiview learning aims to build models that can process multiview data so that it can achieve a better classification performance and make the system more robust. It has successfully been applied to many fields, such as drug target prediction [31], concept approximation [32], among other [33]–[35]. Formally, let $\mathcal{X} = \mathbb{R}^{m_1} \times \mathbb{R}^{m_2} \times \dots \times \mathbb{R}^{m_{|V|}}$ denote the instance space (or feature space) of $|V|$ view representations, where $m_i (1 \leq i \leq |V|)$ denotes the feature dimension of i th view and $\mathcal{Y} = \{l_1, l_2, \dots, l_q\}$ denotes the label space with q class labels. Denote \mathcal{D} as an unknown distribution over $\mathcal{X} \times \mathcal{Y}$. A training set $\mathcal{D} = \{(\mathbf{x}_i^v, y_i) | 1 \leq v \leq |V|, 1 \leq i \leq n\} \in (\mathcal{X} \times \mathcal{Y})^n$ is drawn identically and independently from \mathcal{D} , where $\mathbf{x}_i^v = (x_{i1}^v, x_{i2}^v, \dots, x_{im_v}^v) \in \mathbb{R}^{m_v}$ is the v th view feature vector with dimension m_v and $y_i \in \mathcal{Y}$ is the known label associated with \mathbf{x}_i^v . The task of multiview recognition is to learn a predictive function $f: \mathcal{X} \mapsto \mathcal{Y}$ from \mathcal{D} which can assign a proper label $f(\mathbf{x}) \in \mathcal{Y}$ to an unseen instance \mathbf{x} .

A learner can be denoted as a two-tuple $\mathcal{L} = (h, \mathcal{F})$, where h is a learned decision function also called a classifier, and \mathcal{F} is a fusion function. Fusion plays a very important role in multiview learning and it has attracted much research effort [36].

1) *Basic Fusion Operators*: There are some simple yet efficient fusion operators, such as concatenation [5], elementwise addition [6], elementwise multiplication [7], elementwise max [8], and elementwise average.

Concatenation: The information from multiple views is fused as follows:

$$o(\mathbf{x}_i) = [x_i^1, x_i^2, \dots, x_i^{|V|}] \quad (1)$$

where $[\cdot, \cdot]$ is the concatenation operator.

Elementwise fusion operators require that the dimensions of input vectors are the same, hence different view features need to be mapped into the same dimension space by a linear function before fusion. This can be achieved using a fully connected (FC) layer without any activation function.

Addition: The information from $|V|$ views is fused as follows:

$$o(\mathbf{x}_i) = \text{FC}(\mathbf{x}_i^1) + \text{FC}(\mathbf{x}_i^2) + \dots + \text{FC}(\mathbf{x}_i^{|V|}). \quad (2)$$

Multiplication: The information from $|V|$ views is fused as follows:

$$o(\mathbf{x}_i) = \text{FC}(\mathbf{x}_i^1) \circ \text{FC}(\mathbf{x}_i^2) \circ \dots \circ \text{FC}(\mathbf{x}_i^{|V|}) \quad (3)$$

where \circ denotes Hadamard product, namely elementwise multiplication.

Max: The information from $|V|$ views is fused as follows:

$$o(\mathbf{x}_i) = \max(\text{FC}(\mathbf{x}_i^1), \text{FC}(\mathbf{x}_i^2), \dots, \text{FC}(\mathbf{x}_i^{|V|})) \quad (4)$$

where max is elementwise max, also called max-pooling.

²<https://github.com/ncats/molvec>

Average: The information from $|V|$ views is fused as follows:

$$o(x_i) = \frac{1}{|V|} \left(\text{FC}(x_i^1) + \text{FC}(x_i^2) + \dots + \text{FC}(x_i^{|V|}) \right) \quad (5)$$

where $+$ denotes elementwise addition, also called average-pooling.

2) *Advanced Fusion Methods:* Recently, two advanced fusion methods, namely, bilinear-based fusion [4], [9] and tensor-based fusion [10], [37], have been proposed.

Bilinear methods model all pairwise interactions among features from different views and provide a richer representation than linear methods. For example, multimodal low-rank bilinear (MLB) pooling approach [38] is to solve the dimension curse in feature fusion, it approximates the outer product by projecting first different view features into low-dimensional spaces and then performs elementwise multiplication on the projected features. The fusion process can be formalized as follows:

$$\begin{aligned} c &= \text{MLB}(v_1, v_2, \dots, v_{|V|}) \\ &= U^T \left(U_1^T v_1 \circ U_2^T v_2 \circ \dots \circ U_{|V|}^T v_{|V|} \right) + b \end{aligned} \quad (6)$$

where \circ denotes elementwise multiplication. $U_i \in \mathbb{R}^{M_i \times d}$ and $c \in \mathbb{R}^m$, where d and m are hyperparameters to determine the dimension of joint embeddings and the output dimension of low-rank bilinear models, respectively.

Noting that MLB could result in insufficient representation, [9] proposed a multimodal factorized bilinear (MFB) pooling. In MFB, the features from different views are first expanded to a high-dimensional space and then integrated the expanded vectors with Hadamard product. Then sum pooling followed by the normalization layers is conducted to squeeze the high-dimensional feature into the compact output feature. The fusion process can be formalized as follows:

$$\begin{aligned} c &= \text{MFB}(v_1, v_2, \dots, v_{|V|}) \\ &= \text{SumPool} \left(\hat{U}_1^T v_1 \circ \hat{U}_2^T v_2 \circ \dots \circ \hat{U}_{|V|}^T v_{|V|}, k \right) \end{aligned} \quad (7)$$

where the function $\text{SumPool}(x, k)$ uses a 1-D nonoverlap window with size k to do sum pooling over x .

Tensor-based methods model interactions among different view features by using a $|V|$ -fold Cartesian product from view embeddings. Recently, many efficient models have been proposed. For example, [10] developed a tensor fusion network (TFN) by introducing a tensor fusion layer. Given $|V|$ view vectors $\{v_i \in \mathbb{R}^{m_i}\}_{i=1}^{|V|}$, they are fused as follows:

$$c = \begin{bmatrix} v_1 \\ 1 \end{bmatrix} \otimes \begin{bmatrix} v_2 \\ 1 \end{bmatrix} \otimes \dots \otimes \begin{bmatrix} v_{|V|} \\ 1 \end{bmatrix} \quad (8)$$

where \otimes is the Kronecker product operator. It is worth noting that the output tensor $c \in \mathbb{R}^{(m_1+1) \times (m_2+1) \times \dots \times (m_{|V|}+1)}$ could be of high dimension and this could easily cause curse of dimensionality. Hence, it is only applicable on a very small number of views.

B. Network Architecture Search

DNN learning has successfully been applied to many areas, such as face recognition and speaker recognition. It is well known that network architectures play an critical role. Neural architecture search (NAS) is to search for an optimal network structure in an automatic manner. An NAS often consists of its search space definition, search strategy selection, and model evaluation.

The search space can be classified into macro and micro search spaces [39]. The macro search space is mainly for information of global structure [40], such as the number of layers, the operation types of each layer, and the hyper parameters of each operation. The micro search space is mainly for the change of repeated blocks or cells [1], [41].

The search strategy of network structure often uses reinforcement learning (RL), EA, and gradient-based method.

The RL-based search strategy gives an agent a reward as instructional feedback in an interactive way to find the optimal strategy in a finite-horizon environment. MetaQNN [42] models the neural architecture search as Markov decision process, and uses RL method to generate the convolutional neural network (CNN) architecture. Zoph and Le [43] used the recurrent neural network (RNN) as a controller to sample and generate the string description of a network structure. This structure is then trained and evaluated, and then the RL is used to learn the controller's parameters so that it can produce a more accurate network structure.

The EA-based search strategy uses the validation accuracy as instructional feedback to select the optimal model. In [40] and [44], some neural network structures with one input layer, one output layer and one global pooling layer are first initialized as initial individuals. In the process of evolution, new network structures are obtained using crossover and mutation, new parent population will be selected from the parent and offspring population. Compared with RL, EA can achieve similar accuracy, but it is faster and can produce smaller models.

Gradient-based search is much faster than RL-based and EA-based NAS methods. A gradient-based strategy using differentiable architecture sampling proposed in [39] needs only a few hours to obtain an optimal model. However, gradient-based search often requires much more computer memories.

III. PROPOSED METHOD

In this section, the details of the proposed EDF are presented. As shown in Fig. 1, the framework of EDF consists of two main stages, this first one is to extract multiview features (Section III-A) and the second one is to find a proper deep fusion network (Section III-B).

A. Extracting Multiview Features

We use some different DNNs as different view feature extractors, these DNNs will be trained on three datasets: 1) ChemBook-10k; 2) ChEMBL-10k; and 3) PubChem-10k. Next, similar to other works [45], chemical structure images will be successively fed into the trained models to extract the penultimate layer vector as data representation, i.e., a view. The pseudocode of this process is given in Algorithm 1.

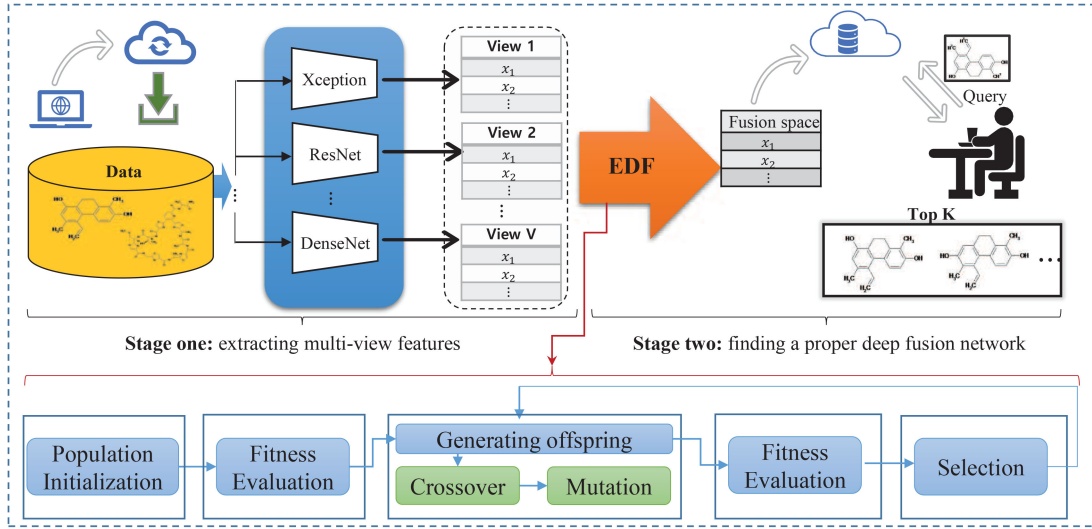


Fig. 1. Overall framework of EDF.

Algorithm 1 Pseudocode of Extracting Multiview Features

Input: A chemical structure recognition training dataset $D = (X, Y)$, test dataset $\hat{D} = (\hat{X}, \hat{Y})$, and multiple deep network set $NET = \{Net_i\}_{i=1}^{|NET|}$.

Output: A multiview training dataset $V = \{V_i\}_{i=1}^{|V|}$ and test dataset $\hat{V} = \{\hat{V}_i\}_{i=1}^{|V|}$.

```

1: for  $i = 1$  to  $|V|$  do
2:   Train  $Net_i$  on  $D$ ;
3:    $V_i \leftarrow Net_i(X)$ , take  $X$  as input and output  $V_i$ ;
4:    $\hat{V}_i \leftarrow Net_i(\hat{X})$ , take  $\hat{X}$  as input and output  $\hat{V}_i$ ;
5: end for
6: return  $V$  and  $\hat{V}$ .

```

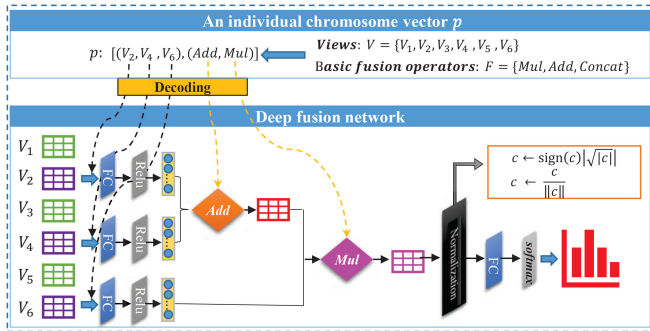


Fig. 2. Illustrative example of the decoding process from an individual chromosome vector to a deep fusion network.

B. Finding Proper Deep Fusion Network With EDF

In the following, we will introduce the encoding and decoding methods, and the framework of the proposed evolutionary multiview fusion method.

1) *Encoding and Decoding (Encoding):* We propose a variable-length encoding strategy for deep fusion networks.

Specifically, an individual is encoded as a list of two parts, the first part is for views, and the second is for the fusion scheme used.

Let V be a set of views and F be a set of basic fusion operators. Then an individual chromosome vector p can be represented as

$$p = [v, f]$$

where $v = (v_1, v_2, \dots, v_k)$ where each v_i is an element from V ; and $f = (f_1, f_2, \dots, f_{k-1})$ where each f_i is an element from F .

We should point out that each individual chromosome vector may have a different k value.

Decoding: We decode an individual chromosome vector $p = [v, f]$ to a deep fusion network as shown in Fig. 2. The corresponding network takes $v = (v_1, \dots, v_k)$ as its input and works as follows.

- 1) Transfer each v_i to u_i by a fully connected layer and then a Relu function.
- 2) Fuse u_1, \dots, u_k as follows:
 - a) $c = u_1$;
 - b) for $i = 1$ to $k-1$, $c \leftarrow$ the result of fusion operator f_i on c and u_{i+1} .
- 3) Normalize c

$$c \leftarrow \text{sign}(c) \cdot \sqrt{|c|} \quad (9)$$

$$c \leftarrow \frac{c}{\|c\|}. \quad (10)$$

- 4) Transfer c to a probability vector \hat{y} by a fully connected layer and a softmax function.

Our major reason for transferring v_i to u_i is to make sure that all the u_i 's are of the same dimension.

The parameters to learn in the deep fusion work include weights in these fully connected layers. This network can be used for classification.

2) *Framework of EDF:* In the following, we give the detailed steps of EDF including population initialization, fitness evaluation, offspring generation, and selection.

Population Initialization: We randomly generate an initial population of N individual chromosome vectors. Each individual chromosome p can have a different k value.

Fitness Evaluation: To evaluate the fitness of each individual chromosome vector p in the current population, we decode it to a deep fusion neural network, train it on a multiview training dataset and then evaluate its classification accuracy on a test dataset. The fitness of p is the classification accuracy.

Noting that at each generation, we need to evaluate the fitness of N individual chromosome vectors (when N is the population size). In our implementation, we train and evaluate their corresponding deep fusion networks in parallel. To further reduce the computational overhead, we record all the evaluated vectors and do not re-evaluate a individual chromosome vector p if it has already been evaluated.

Crossover: Given two chromosome vectors $p^1 = [v^1, f^1]$ and $p^2 = [v^2, f^2]$, we do the following crossover to generate two new chromosome vectors $p_o^1 = [v_o^1, f_o^1]$ and $p_o^2 = [v_o^2, f_o^2]$.

- 1) Do one-point crossover on v^1 and v^2 to produce v_o^1 and v_o^2 .
- 2) Do one-point crossover on f^1 and f^2 to produce f_o^1 and f_o^2 .
- 3) Set $p_o^1 = [v_o^1, f_o^1]$ and $p_o^2 = [v_o^2, f_o^2]$.
- 4) Repair each of p_o^1 and p_o^2 to make sure that it is feasible as follows.
 - a) If $|v_o| - 1 < |f_o|$, delete the $(|f_o| - |v_o| + 1)$ most left elements in f_o .
 - b) If $|v_o| - 1 > |f_o|$, delete the $(|v_o| - |f_o| - 1)$ most left elements in v_o .

Now we give an example of crossover. Let $p_1 = [(3, 1, 5, 4), (1, 3, 3)]$ and $p_2 = [(4, 3, 6), (2, 1)]$. Suppose that 1) gives $v_o^1 = (3, 3, 6)$ and $v_o^2 = (4, 1, 5, 4)$ and 2) produces $f_o^1 = (2, 3, 3)$ and $f_o^2 = (1, 1)$. Then 3) will gives $p_o^1 = [(3, 3, 6), (2, 3, 3)]$ and $p_o^2 = [(4, 1, 5, 4), (1, 1)]$. After repairing 4), $p_o^1 = [(3, 3, 6), (3, 3)]$ and $p_o^2 = [(1, 5, 4), (1, 1)]$.

Mutation: Given $p = [v, f]$, mutation alters some randomly selected elements in v and f .

Selection: We use binary tournament selection in our experiments [46].

The EDF is shown in Algorithm 2.

IV. EXPERIMENTAL STUDIES

A. Datasets

In our experiments, three chemical structure recognition datasets are used to study our proposed EDF. Each dataset includes 10 000 classes. These three datasets are ChemBook-10k, ChEMBL-10k, and PubChem-10k collected from the Chemical Book Website,³ Pubchem⁴ and ChEMBL,⁵ respectively. In the following, we take ChemBook-10k as an example to explain how these datasets are collected.

We first collect 10 000 chemical structure images of different compounds. Each image is classified as a different class. Then we perform the following nine operators on each image to generate another nine images to each class.

Flip: It flips along horizontal or vertical orientation. Fig. 3(2) gives an example of horizontal reflection. We randomly choose one from horizontal and vertical reflection.

³<https://www.chemicalbook.com/>

⁴<https://pubchem.ncbi.nlm.nih.gov/>

⁵<https://www.ebi.ac.uk/chembl/>

Algorithm 2 Evolutionary EDF

Input: N : population size;

T : maximal generation number;

$D = (X, Y)$: training dataset;

$\hat{D} = (\hat{X}, \hat{Y})$: test dataset;

F : a set of basic fusion operators;

NET : a set of DNNs.

Output: A deep fusion network.

- 1: Extract multi-view features using Algorithm 1 that takes D , \hat{D} and Net as inputs and outputs V and \hat{V} ;
- 2: Generate an initial population P_0 ;
- 3: Evaluate the fitness of each chromosome vector in P_0 ;
- 4: **for** $t = 1$ to T **do**
- 5: Generate offspring Q_t using the crossover operator;
- 6: Conduct mutation on each chromosome in Q_t ;
- 7: Evaluate the fitness of each chromosome in Q_t ;
- 8: Select next generation population P_{t+1} from $Q_t \cup P_t$ using a selection operator;
- 9: **end for**
- 10: $p_{best} \leftarrow$ Select the chromosome with the best fitness from P_T .
- 11: **return** the fusion network corresponding to p_{best} .

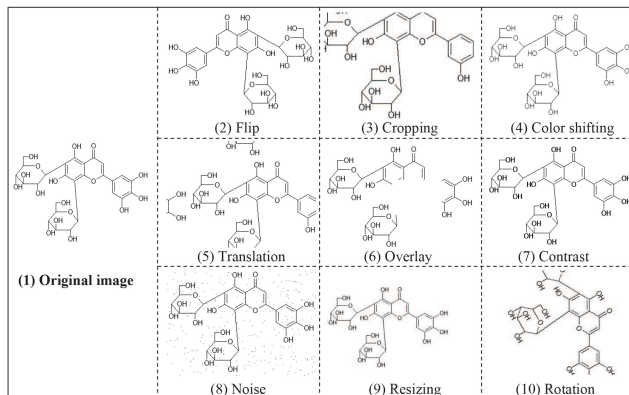


Fig. 3. Original image and new images generated by nine operators.

Cropping: It crops a rectangle region of any size on an image randomly, and then resizes it to the original size. Fig. 3(3) is an example of cropping.

Color Shifting: It generates a new image by adjusting the saturation, brightness, contrast, and sharpness of an image. Fig. 3(4) is an example of color shifting.

Translation: It translates the original image by random values along horizontal and vertical orientation. Fig. 3(5) is an example of translation.

Overlay: It takes a rectangle region of any size on the original image randomly. Fig. 3(6) is an example of overlay.

Contrast: It generates a new image by adjusting the contrast of the original image. Fig. 3(7) is an example of contrast.

Noising: It generates a new image by adding a Gaussian ($\sigma = 0.3$) noise to the original image. Fig. 3(8) is an example of noise.

Resizing: It resizes the original image to a small one, and then uses the background color of original image to fill the

gap between the new size and original size. Fig. 3(9) is an example of resizing.

Rotation: It rotates the original image at a random angle. The new region beyond original size is cropped, and then the gap between the original size and new size is filled by the background color of original image. Fig. 3(10) is an example of rotation.

Then each class has ten images. Then we randomly choose one image from each class to form the test set. The remaining images will be used as the training set.

To facilitate neural network training process, the following preprocessing operations are conducted.

- 1) Resize the size of each image to the same size 230×230 to ensure DNNs can take them as inputs.
- 2) The background of original images is white and contents are black, turn them from RGB into grayscale to reduce the size of the dataset.
- 3) Normalize each pixel value by

$$x = \frac{x}{127.5} - 1$$

where x denotes a pixel value.

B. Experimental Settings

In our experiments, all methods are implemented using Tensorflow⁶ (version: 2.0.3). Our computational environment is Ubuntu 16.04.4, 512-GB DDR4 RDIMM, 2X 40-Core Intel Xeon CPU E5-2698 v4 @ 2.20GH and NVIDIA Tesla P100 with 16-GB GPU memory.

1) Parameter Settings:

- a) *Training of DNNs:* All DNN models are trained using the Adam algorithm. The learning rate is 0.001, the exponential decay rate for the first moment estimates is 0.9, the exponential decay rate for the second moment estimates is 0.999. Every network is trained for 100 epochs. To avoid overfitting, training process will stop when a neural network model performance does not improve after 10 epochs.
- b) *EA:* To efficiently utilize the GPU resources, the population size is set to be a multiple of the number of GPUs. 7 NVIDIA Tesla P100 GPUs are used, and the population size is set to be 28. Following [47], the number of generations is set to be 20, the probabilities of crossover and mutation are set to 0.9 and 0.2, respectively.
- 2) *Chromosome Vector $p = [v, f]$:* We consider two versions: 1) *reused = False*, different elements in v are not allowed to be the same and 2) *reused = True*, there is no such constraint on v .

3) *Candidate Views and Fusion Operators:* In Algorithm 1 for extracting multiview features, two settings for *NET* are used in the experiments. One is *NET5* = {Resnet50, Densenet121, Xception, InceptionV3, MobileNetV2}, and the other is *NET10* = {Resnet50, Densenet121, Xception, InceptionV3, MobileNetV2, Resnet18, Resnet34, Densenet169, Densenet201, NASNetMobile}.

F , the set of basic fusion operators, is set to include elementwise addition (Add), elementwise multiplication (Mul),

concatenation (Concat), elementwise max (Max), and elementwise average (Avg). Note that the dimension of a fused feature obtained by the concatenation operator will be larger, we use a linear mapping to transfer it back to the feature space of the same dimension.

4) *Performance Metrics:* Top-1 accuracy and Top-5 accuracy are used to evaluate the performances of all the methods

$$\text{Top-1} = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(\text{in_top_k}(\hat{y}_i, y_i, 1)) \quad (11)$$

$$\text{Top-5} = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(\text{in_top_k}(\hat{y}_i, y_i, 5)) \quad (12)$$

where \hat{y}_i denotes the probability vector that a deep fusion network outputs, y_i denotes the ground truth class, and the function $\text{in_top_k}(\hat{y}_i, y_i, k)$ returns whether y_i is in a list that consists of these prediction classes corresponding to the first k highest probability values in \hat{y}_i . $\mathbb{I}(\cdot)$ is a indicator function

$$\mathbb{I}(\cdot) = \begin{cases} 1, & \text{if True} \\ 0, & \text{if False.} \end{cases}$$

The higher values of Top-1 and Top-5 are, the better the performance of the evaluated method is.

C. Compared Methods

1) *Five Single View Methods:* ResNet50 [2], DenseNet121 [41], Xception [48], InceptionV3 [41], and Mobilenetv2 [49].

2) *Five Multiview Baseline Methods:* Addition, average, max, multiplication, and concatenation.

3) *Two Ensemble Learning Methods:*

- a) *Simple Soft Voting (SSV)* [50]: It simply averages the outputs of the five single view methods.
- b) *Maximum Rule (MR)* [51]: It selects the highest confidence score among the outputs of the five single view methods.

4) *Three State-of-the-Art Multiview Methods:*

- a) *MLB* [38]: It has been explained in Section II-A. m is set to be 128 and d takes a value from {64, 128, 256, 512}.
- b) *MFB* [9]: It has been explained in Section II-A. m is set to be 128 and k takes a value from {1, 2, 3, 4}.
- c) *TFN* [10]: It has been explained in Section II-A. Batch normalization (BN) is used in order to avoid overfitting [52]. m is set to be 128 and m_i takes values from {5, 10, 15, 20}.

D. Experimental Results

In the experiments, we first extract five view features using Resnet50, Densenet121, Xception, InceptionV3, and MobileNetV2, respectively. Then, these extracted views of different dimensions are mapped into a dimension of $m = 128$ by a FC layer, so that elementwise fusion operators can be used.

The experimental results are summarized in Table I, where *#Paras.* is the number of parameters to learn, and *Time* is the computing time (in second) for training each neural network model. It is clear from Table I that:

⁶<https://github.com/tensorflow/tensorflow>

TABLE I
COMPARATIVE STUDY

Method	ChemBook-10k				ChEMBL-10k				PubChem-10k			
	# Paras.	Top-1	Top-5	Time (s)	# Paras.	Top-1	Top-5	Time (s)	# Paras.	Top-1	Top-5	Time (s)
ResNet50	44,018,320	64.39%	81.91%	62723.90	44,018,320	65.31%	84.96%	39923.59	44,018,320	69.13%	86.04%	33451.81
DenseNet121	17,197,584	72.53%	85.46%	20737.68	17,197,584	75.08%	90.10%	23896.40	17,197,584	81.94%	93.29%	37778.09
Mobilenetv2	15,033,296	69.98%	83.77%	47982.39	15,033,296	72.99%	86.86%	37790.09	15,033,296	74.57%	87.73%	48738.91
Xception	41,296,376	71.41%	85.81%	47652.38	41,296,376	74.11%	88.24%	54131.52	41,296,376	68.14%	82.45%	46188.60
InceptionV3	42,257,776	76.48%	89.07%	16804.38	42,257,776	75.29%	90.71%	16593.01	42,257,776	74.01%	84.84%	34125.32
SSV	-	81.01%	91.14%	85.65	-	84.78%	93.95%	85.03	-	85.84%	93.53%	86.22
MR	-	78.92%	90.75%	85.65	-	82.39%	93.50%	85.03	-	84.35%	93.31%	86.22
Addition	2,389,136	82.41%	92.71%	208.55	2,389,136	88.89%	98.03%	197.42	2,389,136	87.22%	96.45%	177.34
Average	2,389,136	82.19%	93.07%	193.75	2,389,136	88.82%	97.49%	201.10	2,389,136	87.50%	96.07%	197.06
Max	2,389,136	80.23%	91.96%	189.31	2,389,136	87.39%	97.02%	207.07	2,389,136	86.92%	96.14%	186.09
Multiplication	2,389,136	81.40%	92.38%	214.12	2,389,136	88.39%	97.90%	201.06	2,389,136	86.38%	95.94%	179.16
Concatenation	7,510,160	80.23%	89.80%	654.09	7,510,160	87.52%	95.19%	528.38	7,510,160	85.79%	92.74%	536.09
MLB	2,785,040	80.04%	91.84%	1227.43	2,785,040	86.38%	96.45%	395.87	2,785,040	85.73%	95.07%	696.21
MFB	2,719,120	84.14%	95.00%	774.68	2,719,120	91.48%	98.85%	764.33	2,636,560	90.77%	97.86%	900.22
TFN	533,335,550	78.11%	90.55%	18132.88	533,335,550	86.60%	96.38%	15097.85	533,335,550	84.53%	93.87%	14564.84
EDF (<i>reused</i> = <i>False</i>)	2,155,664	86.84%	96.66%	17482.56	2,389,136	93.33%	99.46%	19466.58	2,122,768	93.55%	99.16%	19246.73
EDF (<i>reused</i> = <i>True</i>)	2,522,384	87.49%	96.94%	47145.77	3,654,544	93.75%	99.38%	68329.77	3,954,320	93.85%	99.20%	52204.32

TABLE II
EXPERIMENTAL RESULTS OF MLB, MFB, AND TFN IN DIFFERENT SETTINGS ON CHEMBOOK-10K

MLB					MFB					TFN				
m	d	# $Paras.$	Top-1	Top-5	k	# $Paras.$	Top-1	Top-5	m_i	m_c	# $Paras.$	Top-1	Top-5	
64	64	1,794,448	69.34%	86.39%	1	1,790,160	78.93%	92.14%	5	7,776	2,265,641	73.61%	87.51%	
	128	1,839,824	74.25%	88.84%	2	1,831,440	81.62%	93.59%	10	161,051	12,385,016	76.40%	90.32%	
	256	1,930,576	78.75%	92.81%	3	1,872,720	81.51%	93.94%	15	1,048,576	70,964,891	77.78%	91.24%	
	512	2,112,080	79.25%	92.98%	4	1,914,000	82.16%	94.04%	20	4,084,101	271,312,766	77.91%	91.54%	
128	64	2,438,736	67.14%	85.62%	1	2,471,440	82.60%	94.39%	5	7,776	3,403,625	74.91%	87.64%	
	128	2,488,208	79.02%	89.74%	2	2,554,000	83.36%	94.81%	10	161,051	23,332,600	76.79%	89.46%	
	256	2,587,152	79.28%	90.94%	3	2,636,560	84.12%	95.29%	15	1,048,576	138,714,075	78.00%	90.86%	
	512	2,785,040	80.04%	91.84%	4	2,719,120	84.14%	95.00%	20	4,084,101	533,335,550	78.11%	90.55%	
256	64	3,727,312	74.81%	87.49%	1	3,834,000	84.02%	94.97%	5	7,776	5,679,593	74.58%	86.89%	
	128	3,784,976	77.30%	88.37%	2	3,999,120	84.45%	95.27%	10	161,051	45,227,768	76.41%	88.16%	
	256	3,900,304	79.16%	89.44%	3	4,164,240	84.97%	95.70%	15	1,048,576	274,212,443	77.34%	88.46%	
	512	4,130,960	79.49%	89.56%	4	4,329,360	85.03%	95.82%	20	-	-	-	-	
512	64	6,304,464	72.15%	86.79%	1	6,559,120	84.54%	95.18%	5	7,776	10,231,529	73.37%	86.38%	
	128	6,378,512	78.50%	89.29%	2	6,889,360	84.96%	95.59%	10	161,051	89,018,104	76.24%	87.50%	
	256	6,526,608	79.16%	89.58%	3	7,219,600	85.46%	95.91%	15	-	-	-	-	
	512	6,822,800	81.33%	90.75%	4	7,549,840	85.49%	95.77%	20	-	-	-	-	

- 1) multiview methods perform better than single view methods. This suggests that multiview fusion does have advantages. It also implies that the first stage of EDF is very useful for the performance improvement;
- 2) baseline fusion methods statistically work better than MLB and TFN. Noting that MLB and TFN have achieved the-state-of-art results on VQA task and multiview sentiment analysis task [10], [38], respectively. We can conclude that a fusion scheme of different views is very crucial;
- 3) using five simple fusion operators, EDF is 3.35%, 2.27%, 3.08% better on the Top-1 accuracy than the best one of all the compared methods, some of them were well designed for multiviews learning by human experts;
- 4) compared to the other multiview methods, EDF needs long training time. This is because EDF needs to train $N \times T$ deep fusion networks in the worst case ($N = 28$, $T = 20$ in our experiments). It is clear that the training time of EDF is about 84 times of that of Addition on ChemBook-10k. This indicates that parallel

implementation can reduce the clock time. Section IV-E. will further discuss this issue.

In summary, EDF is very competitive compared with other manually designed multiview algorithms. View selection in EDF can remove the redundancy view information, and the fusion scheme automatically obtained by EDF does work.

E. More Analysis

We further investigate the performance of EDF under different experimental settings on ChemBook-10k. The experimental results summarized in Table III show that:

- 1) in general, it can improve the fusion performance if the dimension of the fusion space and the size of candidate view set increase. It is also evident that it is better to allow elements of v in chromosome vector p duplicate. For example, Top-1 and Top-5 accuracy metrics improve from 85.31% to 90.06% and from 95.46% to 98.43% when the setting is changed from $m = 64$, *reused* = *False* and *NET* = *NET5* to $m = 512$, *reused* = *True* and *NET* = *NET10*, respectively;

TABLE III
EXPERIMENTAL RESULTS OF EDF IN DIFFERENT SETTINGS ON CHEMBOOK-10K

<i>NET</i>	<i>reuse</i>	<i>m</i>	Best chromosome	# <i>Paras.</i>	Top-1	Top-5
<i>NET5</i>	<i>False</i>	64	[(1, 4, 2, 0, 3), (0, 0, 4, 4)]	1,208,016	85.31%	95.46%
		128	[(0, 1, 4, 2), (2, 1, 1)]	2,155,664	86.84%	96.66%
		256	[(1, 3, 2, 4), (2, 2, 0)]	4,485,392	88.07%	97.31%
		512	[(2, 3, 4, 1), (3, 4, 2)]	8,947,472	88.46%	97.64%
	<i>True</i>	64	[(3, 0, 1, 4, 2, 1), (2, 4, 0, 4, 4)]	1,283,920	85.52%	96.07%
		128	[(3, 0, 1, 2, 1, 4), (1, 0, 4, 1, 1)]	2,522,384	87.49%	96.94%
		256	[(4, 0, 3, 2, 1, 1, 1, 1, 3, 2, 4, 3, 2, 4, 0, 2, 1, 1), (0, 4, 4, 0, 2, 0, 0, 4, 4, 0, 4, 0, 0, 4, 0, 2)]	9,970,960	88.50%	97.73%
		512	[(2, 1, 0, 4, 1, 3), (4, 0, 3, 2, 0)]	10,527,504	88.58%	97.73%
<i>NET10</i>	<i>False</i>	64	[(8, 9, 1, 4, 6, 7), (4, 4, 0, 0, 0)]	1,193,296	86.52%	96.41%
		128	[(3, 4, 8, 7, 1), (4, 4, 0, 4)]	2,422,416	88.41%	97.54%
		256	[(9, 0, 4, 8, 2, 5, 1, 7), (4, 3, 3, 4, 3, 4, 4)]	5,552,976	89.26%	97.67%
		512	[(0, 9, 3, 4, 1, 2, 7, 6, 8, 5), (3, 3, 0, 4, 0, 0, 0, 4, 2)]	12,914,512	89.89%	98.34%
	<i>True</i>	64	[(0, 1, 6, 5, 1, 6, 2, 4, 7), (4, 1, 1, 4, 1, 0, 0, 1)]	1,351,888	86.73%	96.48%
		128	[(8, 5, 2, 2, 6, 6, 7, 4, 9), (0, 3, 3, 0, 0, 2, 0, 0)]	2,726,224	88.51%	97.67%
		256	[(7, 2, 8, 9, 1, 4, 0, 5, 4, 0, 1, 8, 2, 7, 8, 7, 9, 2, 7), (3, 0, 0, 3, 2, 1, 3, 4, 1, 0, 0, 0, 0, 0, 0, 0)]	10,219,664	89.50%	97.98%
		512	[(0, 6, 6, 3, 2, 8, 9, 7, 6, 4, 7, 1, 3, 8, 4, 7, 1, 2, 5), (2, 2, 4, 2, 2, 3, 4, 4, 0, 4, 0, 4, 0, 4, 0, 2, 4)]	21,531,728	90.06%	98.43%

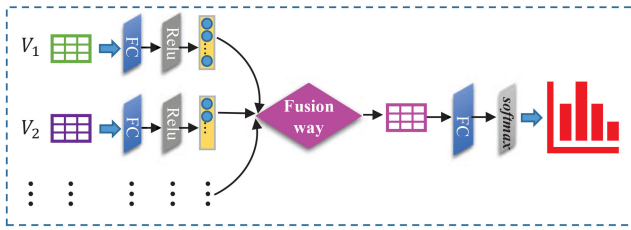


Fig. 4. Architecture of comparative methods.

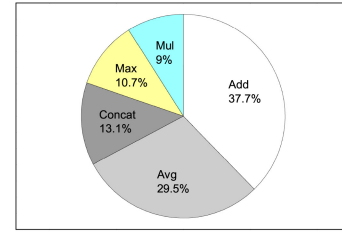


Fig. 6. Use frequency distribution of the basic fusion operators.

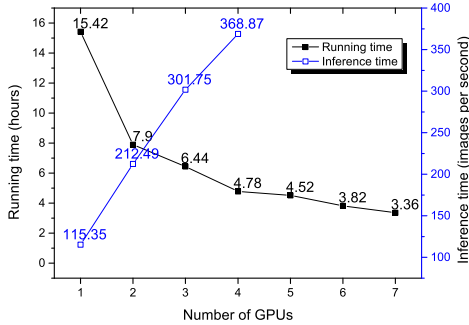


Fig. 5. Training and inference times change with the number of GPUs. The fusion model is obtained in the setup: *NET* = *NET5*, *reuse* = *False*, *m* = 512. Because four views are selected, the maximum number of GPUs is four in the inference process.

- EDF with *m* = 64 performs better than all other compared methods with *m* = 128. For example, EDF with *reuse* = *True* and *m* = 64 obtains the Top-1 accuracy metric of 85.52%, whereas it is 78.11% in TFN with *m* = 128. It indicates that EDF needs much fewer parameters than other methods;
- as shown in Table II, the number of the parameters used in TFN with tensor-based fusion is much larger than other compared methods. This is because the dimension of the fused vector increases exponentially as the dimension of embedding vectors increases. In comparison, EDF with basic fusion operators except Concat does not introduce extra parameters. Actually, the number of parameters can be reduced when some redundancy views are removed. For example, EDF (*m* = 128, *NET* =

NET5, *reuse* = *False*) does not use the view extracted by Xception, and leads to decrease of the number of parameters from 2 389 136 to 2 155 664.

In summary, compared with other methods, EDF with different settings works well and does not introduce more extra parameters. One of major drawbacks of EDF is that its training time is longer than other methods. One possible way for minimizing this drawback to use parallel computing environments. We can do EDP training and inference in parallel. Fig. 5 shows that the training time reduces and the number of inferring images per second increases as the number of GPUs increases.

We have also analyzed the use frequency distribution of the basic fusion operators in 16 final deep fusion networks obtained by EDF in Table III. The use frequency distribution is shown in Fig. 6. It can be observed that elementwise addition and elementwise average are more frequently used than other operators. Surprisingly, the elementwise multiplication used in recent bilinear fusion models is used the least. This observation suggests that more deep understanding on these basic operators is needed.

V. CONCLUSION

We have developed an evolutionary EDF, which can automatically build a good fusion model from given candidate views and basic fusion operator sets. The experimental studies have demonstrated that multiview fusion neural networks generated by EDF perform better than those manually designed by human experts.

This work is a first step toward use of NAS and EAs on multiview learning. Several issues are worthwhile investigating

TABLE A1
COMPARATIVE STUDY ON TINY IMAGENET

Method	# Paras.	Top-1	Top-5
ResNet50	58,539,720	42.43%	69.49%
DenseNet121	7,158,856	53.55%	76.84%
Mobilenetv2	2,480,072	47.55%	73.36%
Xception	21,216,752	48.67%	72.03%
InceptionV3	22,178,152	49.71%	73.43%
Addition	1,124,936	56.63%	78.24%
Average	1,124,936	56.69%	78.21%
Max	1,124,936	53.66%	76.42%
Multiplication	1,124,936	56.58%	78.51%
Concatenation	1,228,360	56.24%	78.26%
MLB	1,520,840	54.89%	76.24%
MFB	1,372,360	54.84%	76.89%
TFN	532,071,350	52.34%	73.90%
EDF (<i>reused</i> = <i>False</i>)	1,124,936	58.54%	78.61%
EDF (<i>reused</i> = <i>True</i>)	3,789,128	59.64%	79.79%

along this direction. For example, how can attention mechanisms or other methods be used to control the contribution of each view to the fusion model [53], how can the training cost of EDF be reduced by using expensive optimization techniques [54] and how can multiobjective techniques be used in EDF [55].

APPENDIX

A. Results on Tiny Imagenet

The Tiny ImageNet dataset⁷ is a subset of the ImageNet. It consists 200 classes while ImageNet has 1000 classes. Each class in Tiny ImageNet contains 500 training images and 50 validation images. The resolution of the images is 64×64 pixels, which makes it more difficult to extract information from it. To make sure that the DNNs used in our experiments can take these images as inputs, we have resized them to 230×230 pixels. In our experiment, we have not used image augmentation. The results are shown in Table A1. It is evident that EDF performs the best.

B. Application

In real-world applications, there exist hundreds of millions of molecules. A practical model has to be able to recognize complete unseen chemical images, i.e., recognition in open-set scenario.

Given a chemical structure image dataset $D = \{(x_i, y_i)\}_{i=1}^n$, where x_i denotes a chemical structure image and y_i is its name. In open-set scenario, EDF works as follows.

- 1) Obtain the deep fusion network with the best classification accuracy $EDFNet$ trained on $\hat{D} = \{(x_i, y_i)\}_{i=1}^m$ ($m \ll n$ in real world) that consists of m random molecules from D .
- 2) Construct a retrieve database $R = \{(c_i, y_i)\}_{i=1}^n$ as follows: each chemical structure image from D is successively fed into the trained model $EDFNet$ to extract the penultimate layer vector as data representation, i.e., $c_i \leftarrow EDFNet(x_i)$, take x_i as input and output c_i .
- 3) Given an unseen image list Q and the name of each molecule x from Q can be obtained as follows:

TABLE A2
EXPERIMENTAL RESULTS OF EDF ON OPEN-SET TASKS

Settings	Rank@1	Rank@5	Rank@10
i	84.84%	95.21%	97.52%
ii	81.57%	93.88%	96.24%

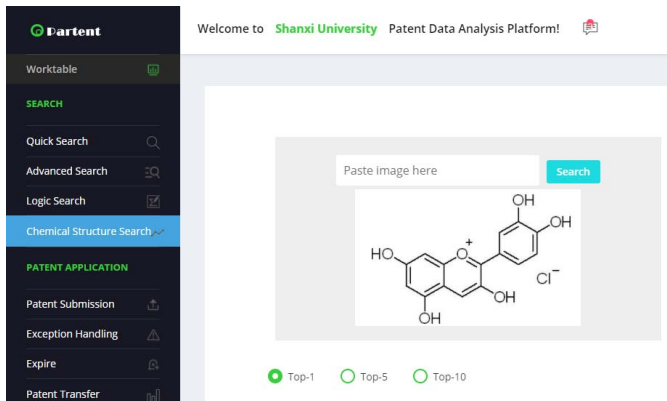


Fig. A1. Interface of image-based patent search.

- a) $c \leftarrow EDFNet(x)$;
- b) Calculate similarity s_i between c and each c_i from R by the Euclidean distance;
- c) Return k chemical structure images corresponding to the first k maximum values in $\{s_i\}_{i=1}^n$, and their name list $\{\hat{y}_i\}_i^k$.

In this way, EDF can generalize for all the molecules available in the real world. The EDF in open-set scenario can be evaluated by

$$\text{Rank}@k = \frac{1}{|Q|} \sum_{(x,y) \in Q} y \in \{\hat{y}_i\}_i^k$$

where y is the true name of the unseen image x from Q .

The results of EDF on open-set tasks are shown in Table A2. In our experiment, D consists of all images from PubChem-10k and ChEMBL-10k. We consider two settings for \hat{D} and Q : 1) \hat{D} consists of all images from PubChem-10k and Q consists of all images from the test set of ChEMBL-10k and 2) \hat{D} consists of all images from ChEMBL-10k and Q consists of all images from the test set of PubChem-10k. It is evident that EDF still works well.

Using EDF,⁸ we have developed an image-based patent search system in a patent data analysis platform at Shanxi University. As shown in Fig. A1, molecular structure search based on EDF has been used as one of the four search ways (other three are quick search, advanced search and logic search). Different from other three types of search ways based on text query, EDF based on image query may be more convenient and efficient for cheminformatics researchers in most cases. It is worth noting that there is little restriction for query image, such as size, format, resolution, which brings very good user experience.

⁷<http://cs231n.stanford.edu/tiny-imagenet-200.zip>

⁸The code is available at <https://github.com/xinyanliang/EDF>.

REFERENCES

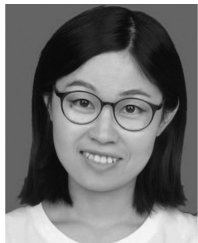
- [1] C. Szegedy *et al.*, "Going deeper with convolutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2015, pp. 1–9.
- [2] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.
- [3] G. Huang, Z. Liu, L. V. Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 2261–2269.
- [4] T. Lin, A. RoyChowdhury, and S. Maji, "Bilinear CNN models for fine-grained visual recognition," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, 2015, pp. 1449–1457.
- [5] L. Wang, W. Li, W. Li, and L. Van Gool, "Appearance-and-relation networks for video classification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 1430–1439.
- [6] Z. Wu, Y. Jiang, J. Wang, J. Pu, and X. Xue, "Exploring inter-feature and inter-class relationships with deep neural networks for video classification," in *Proc. 22nd ACM Int. Conf. Multimedia (MM)*, 2014, pp. 167–176.
- [7] P. Gao *et al.*, "Dynamic fusion with intra- and inter-modality attention flow for visual question answering," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 6639–6648.
- [8] C. T. Duong, R. Lebrete, and K. Aberer, "Multimodal classification for analysing social media," in *Proc. Eur. Conf. Mach. Learn. Principles Practice Knowl. Disc. Databases (ECML-PKDD)*, 2017, pp. 999–1014.
- [9] Z. Yu, J. Yu, C. Xiang, J. Fan, and D. Tao, "Beyond bilinear: Generalized multimodal factorized high-order pooling for visual question answering," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 29, no. 12, pp. 5947–5959, Dec. 2018.
- [10] A. Zadeh, M. Chen, S. Poria, E. Cambria, and L. Morency, "Tensor fusion network for multimodal sentiment analysis," in *Proc. Conf. Empirical Methods Nat. Lang. Process.*, 2017, pp. 1103–1114.
- [11] P. George *et al.*, "SureChEMBL: A large-scale, chemically annotated patent document database," *Nucl. Acids Res.*, vol. 44, no. D1, pp. D1220–D1228, 2016.
- [12] M. Kratochvil, J. Vondrášek, and J. Galgonek, "Sachem: A chemical cartridge for high-performance substructure search," *J. Cheminformatics*, vol. 10, no. 1, pp. 1–11, 2018.
- [13] H. Shang, Y. Tao, Y. Gao, C. Zhang, and X. Wang, "An improved invariant for matching molecular graphs based on VF2 algorithm," *IEEE Trans. Syst., Man, Cybern., Syst.*, vol. 45, no. 1, pp. 122–128, Jan. 2016.
- [14] E. Proschak, J. K. Wegner, A. Schüller, G. Schneider, and U. Fechner, "Molecular query language (MQL)—A context-free grammar for substructure matching," *J. Chem. Inf. Model.*, vol. 47, no. 2, pp. 295–301, 2007.
- [15] A. Juttner and P. Madarasi, "VF2++—An improved subgraph isomorphism algorithm," *Discr. Appl. Math.*, vol. 242, pp. 69–81, Jun. 2018.
- [16] V. Carletti, P. Foggia, A. Saggese, and M. Vento, "Introducing VF3: A new algorithm for subgraph isomorphism," in *Graph-Based Representations in Pattern Recognition*, P. Foggia, C.-L. Liu, and M. Vento, Eds. Cham, Switzerland: Springer, 2017, pp. 128–139.
- [17] L. P. Cordella, P. Foggia, C. Sansone, and M. Vento, "A (sub)graph isomorphism algorithm for matching large graphs," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 26, no. 10, pp. 1367–1372, Oct. 2004.
- [18] D. Weininger, "SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules," *J. Chem. Inf. Model.*, vol. 28, no. 1, pp. 31–36, 1988.
- [19] A. Dalby *et al.*, "Description of several chemical structure file formats used by computer programs developed at molecular design limited," *J. Chem. Inf. Comput. Sci.*, vol. 32, no. 3, pp. 244–255, 1992.
- [20] L. Mihai, M. Katja, K. Noriko, and J. T. Anthony, Eds., *Current Challenges in Patent Information Retrieval* (The Information Retrieval Series), vol. 37. Berlin, Germany: Springer-Verlag, 2017.
- [21] M. Oldenhof, A. Arany, Y. Moreau, and J. Simm, "Chemgrapher: Optical graph recognition of chemical compounds by deep learning," *J. Chem. Inf. Model.*, vol. 60, no. 10, pp. 4506–4517, 2020.
- [22] N. M. Sadawi, A. P. Sexton, and V. Sorge, "Chemical structure recognition: A rule-based approach," in *Proc. Recognit. Retrieval XIX*, 2012, Art. no. 82970E.
- [23] V. Smolov, F. Zentsev, and M. Rybalkin, "IMAGO: Open-source toolkit for 2D chemical structure image recognition," in *Proc. 12th Text Retrieval Conf.*, 2011, p. 675.
- [24] I. V. Filippov and M. C. Nicklaus, "Optical structure recognition software to recover chemical information: Osra, an open source solution," *J. Chem. Inf. Model.*, vol. 49, no. 3, pp. 740–743, 2009.
- [25] K. Rajan, H. O. Brinkhaus, A. Zieslesny, and C. Steinbeck, "A review of optical chemical structure recognition tools," *J. Cheminform.*, vol. 60, no. 10, pp. 1–13, 2020.
- [26] J. G. Richens, C. M. Lee, and S. Johri, "Improving the accuracy of medical diagnosis with causal machine learning," *Nat. Commun.*, vol. 11, p. 3923, Aug. 2020.
- [27] M. H. S. Segler, M. Preuss, and M. P. Waller, "Planning chemical syntheses with deep neural networks and symbolic Ai," *Nature*, vol. 555, no. 3923, pp. 604–610, 2018.
- [28] D. Mendez *et al.*, "ChEMBL: Towards direct deposition of bioassay data," *Nucl. Acids Res.*, vol. 47, no. D1, pp. D930–D940, 2018.
- [29] S. Kim *et al.*, "PubChem 2019 update: Improved access to chemical data," *Nucl. Acids Res.*, vol. 47, no. D1, pp. D1102–D1109, 2018.
- [30] S. Jaeger, S. Fulle, and S. Turk, "MOL2VEC: Unsupervised machine learning approach with chemical intuition," *J. Chem. Inf. Model.*, vol. 58, no. 1, pp. 27–35, 2017.
- [31] L. Li and M. Cai, "Drug target prediction by multi-view low rank embedding," *IEEE/ACM Trans. Comput. Biol. Bioinform.*, vol. 16, no. 5, pp. 1712–1721, Sep./Oct. 2019.
- [32] Y. Qian, X. Liang, G. Lin, Q. Guo, and J. Liang, "Local multigranulation decision-theoretic rough sets," *Int. J. Approx. Reason.*, vol. 82, pp. 119–137, Mar. 2017.
- [33] J. Wang, Y. Qian, F. Li, J. Liang, and W. Ding, "Fusing fuzzy monotonic decision trees," *IEEE Trans. Fuzzy Syst.*, vol. 28, no. 5, pp. 887–900, May 2020.
- [34] F. Li, Y. Qian, J. Wang, C. Dang, and L. Jing, "Clustering ensemble based on sample's stability," *Artif. Intell.*, vol. 273, pp. 37–55, Aug. 2019.
- [35] T. Yan, Z. Hu, Y. Qian, Z. Qiao, and L. Zhang, "3D shape reconstruction from multifocus image fusion using a multidirectional modified Laplacian operator," *Pattern Recognit.*, vol. 98, Feb. 2020, Art. no. 107065.
- [36] T. Baltrušaitis, C. Ahuja, and L.-P. Morency, "Multimodal machine learning: A survey and taxonomy," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 2, pp. 423–443, Feb. 2019.
- [37] Z. Liu, Y. Shen, V. B. Lakshminarasimhan, P. P. Liang, A. Zadeh, and L. Morency, "Efficient low-rank multimodal fusion with modality-specific factors," in *Proc. 56th Annu. Meeting Assoc. Comput. Linguist.*, 2018, pp. 2247–2256.
- [38] J.-H. Kim, K. W. On, W. Lim, J. Kim, J.-W. Ha, and B.-T. Zhang, "Hadamard product for low-rank bilinear pooling," in *Proc. 5th Int. Conf. Learn. Represent.*, 2017, pp. 2285–2294.
- [39] X. Dong and Y. Yang, "Searching for a robust neural architecture in four GPU hours," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 1761–1770.
- [40] E. Real *et al.*, "Large-scale evolution of image classifiers," in *Proc. Int. Conf. Mach. Learn.*, 2017, pp. 2902–2911.
- [41] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 2818–2826.
- [42] B. Baker, O. Gupta, N. Naik, and R. Raskar, "Designing neural network architectures using reinforcement learning," in *Proc. Int. Conf. Learn. Represent.*, 2017, pp. 1–6.
- [43] B. Zoph and Q. V. Le, "Neural architecture search with reinforcement learning," in *Proc. Int. Conf. Learn. Represent.*, 2017, pp. 1–5.
- [44] E. Real, A. Aggarwal, Y. Huang, and Q. V. Le, "Regularized evolution for image classifier architecture search," in *Proc. AAAI*, vol. 33, 2019, pp. 4780–4789.
- [45] A. Nagrani, J. S. Chung, W. Xie, and A. Zisserman, "Voxceleb: Large-scale speaker verification in the wild," *Comput. Speech Lang.*, vol. 60, pp. 1–15, Mar. 2020.
- [46] B. L. Miller and D. E. Goldberg, "Genetic algorithms, tournament selection, and the effects of noise," *Complex Syst.*, vol. 9, no. 3, pp. 193–212, 1995.
- [47] T. Bäck, *Evolutionary Algorithms in Theory and Practice. Evolution Strategies, Evolutionary Programming, Genetic Algorithms*. New York, NY, USA: Oxford Univ. Press, 1996.
- [48] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 1800–1807.
- [49] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L. Chen, "MobileNetV2: Inverted residuals and linear bottlenecks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 4510–4520.
- [50] Z.-H. Zhou, *Ensemble Methods: Foundations and Algorithms*, 1st ed. London, U.K.: Chapman & Hall, 2012.
- [51] Y. Peng *et al.*, "Multimodal ensemble fusion for disambiguation and retrieval," *IEEE MultiMedia*, vol. 23, no. 2, pp. 42–52, Apr.–Jun. 2016.

- [52] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *Proc. Int. Conf. Mach. Learn.*, 2015, pp. 448–456.
- [53] A. Vaswani *et al.*, "Attention is all you need," in *Proc. 31st Int. Conf. Neural Inf. Process. Syst.*, 2017, pp. 6000–6010.
- [54] T. Dick, M. Li, V. K. Pillutla, C. White, N. Balcan, and A. J. Smola, "Data driven resource allocation for distributed learning," in *Proc. 20th Int. Conf. Artif. Intell. Stat.*, vol. 54, 2017, pp. 662–671.
- [55] Q. Zhang and H. Li, "MOEA/D: A multiobjective evolutionary algorithm based on decomposition," *IEEE Trans. Evol. Comput.*, vol. 11, no. 6, pp. 712–731, Dec. 2007.



Xinyan Liang (Student Member, IEEE) received the B.Sc. degree from the School of Computer and Information Technology, Shanxi University, Taiyuan, China, in 2014, where he is currently pursuing the Ph.D. degree with the Institute of Big Data Science and Industry.

His main research interests include multimodal machine learning, multiview learning, granular computing, and their applications.



Qian Guo (Student Member, IEEE) received the B.Sc. degree from the School of Computer and Information Technology, Shanxi University, Taiyuan, China, in 2014, where she is currently pursuing the Ph.D. degree with the Institute of Big Data Science and Industry.

Her main research interests include logic learning and its applications, such as multi-image retrieval and abstract reasoning.



Yuhua Qian (Member, IEEE) received the M.S. and Ph.D. degrees in computers with applications from Shanxi University, Taiyuan, China, in 2005 and 2011, respectively.

He is currently a Professor with the Key Laboratory of Computational Intelligence and Chinese Information Processing, Ministry of Education, Shanxi University. He is best known for multigranulation rough sets in learning from categorical data and granular computing. He is involved in research on machine learning, pattern recognition, feature selection, granular computing, and artificial intelligence. He has authored over 100 articles on these topics in international journals.

Prof. Qian has served as the Program Chair or Special Issue Chair of the Conference on Rough Sets and Knowledge Technology, the Joint Rough Set Symposium, and the Conference on Industrial Instrumentation and Control. He served on the Editorial Board of the *International Journal of Knowledge-Based Organizations* and *Artificial Intelligence Research*. He is a PC member of many machine learning, data mining conferences.



Weiping Ding (Senior Member, IEEE) received the Ph.D. degree in computation application from the Nanjing University of Aeronautics and Astronautics, Nanjing, China, in 2013.

He was a Visiting Scholar with the University of Lethbridge, Lethbridge, AB, Canada, in 2011. From 2014 to 2015, he was a Postdoctoral Researcher with the Brain Research Center, National Chiao Tung University, Hsinchu, Taiwan. In 2016, he was a Visiting Scholar with the National University of Singapore, Singapore. From 2017 to 2018, he was a

Visiting Professor with the University of Technology Sydney, Ultimo, NSW, Australia.

Dr. Ding is currently the Chair of IEEE CIS Task Force on Granular Data Mining for Big Data. His main research directions involve data mining, granular computing, evolutionary computing, machine learning, and big data analytics. He has published more than 80 research peer-reviewed journal and conference papers, including *IEEE TRANSACTIONS ON FUZZY SYSTEMS*, *IEEE TRANSACTIONS ON NEURAL NETWORKS AND LEARNING SYSTEMS*, *IEEE TRANSACTIONS ON CYBERNETICS*, *IEEE TRANSACTIONS ON SYSTEMS, MAN, AND CYBERNETICS: SYSTEMS*, *IEEE TRANSACTIONS ON BIOMEDICAL ENGINEERING*, *IEEE TRANSACTIONS ON INDUSTRIAL INFORMATICS*, *IEEE TRANSACTIONS ON EMERGING TOPICS IN COMPUTATIONAL INTELLIGENCE*, and *IEEE TRANSACTIONS ON INTELLIGENT TRANSPORTATION SYSTEMS*. He served on the Editorial Advisory Board of *Knowledge-Based Systems*, and serves on the Editorial Board of *Information Fusion*, and *Applied Soft Computing*. He serves as an Associate Editor for *IEEE TRANSACTIONS ON FUZZY SYSTEMS*, *IEEE/CAA Journal of Automatica Sinica*, *Information Sciences*, *Neurocomputing*, *Swarm and Evolutionary Computation*, *IEEE ACCESS*, and *Journal of Intelligent & Fuzzy Systems*, and Co-Editor-in-Chief of *Journal of Artificial Intelligence and System*. He is the Leading Guest Editor of Special Issues in several prestigious journals, including *IEEE TRANSACTIONS ON EVOLUTIONARY COMPUTATION*, *IEEE TRANSACTIONS ON FUZZY SYSTEMS*, and *Information Fusion*. He is a member of IEEE-CIS, ACM, CCAI, and Senior CCF. He is a member of Technical Committee on Soft Computing of IEEE SMCS, on Granular Computing of IEEE SMCS, and on Data Mining and Big Data Analytics of IEEE CIS.



Qingfu Zhang (Fellow, IEEE) received the B.Sc. degree in mathematics from Shanxi University, Taiyuan, China, in 1984, the M.Sc. degree in applied mathematics and the Ph.D. degree in information engineering from Xidian University, Xi'an, China, in 1991 and 1994, respectively.

He is a Chair Professor of Computational Intelligence with the Department of Computer Science, City University of Hong Kong, Hong Kong. His main research interests include evolutionary computation, optimization, neural networks, data analysis, and their applications. He is an Associate Editor of the *IEEE TRANSACTIONS ON EVOLUTIONARY COMPUTATION* and the *IEEE TRANSACTIONS ON CYBERNETICS*. He has been a Web of Science highly cited researcher in Computer Science since 2016.