## RESEARCH ARTICLE

# A Research of Gas Open-Set Identification Based on Data Augmentation Algorithm

**YE ZHU** AND **JINGYA WANG**
School of Computer Science and Technology, Anhui University of Technology, Maanshan 243032, China
Corresponding author: Ye Zhu (zhuye@ahut.edu.cn)

**ABSTRACT** Significant progress has been made in convolutional neural networks (CNN) based gas recognition. However, existing electronic nose (e-Nose) algorithms all use the closed-set assumption that the test and training samples are in the same label space and can only detect objects of known classes. However, in realistic scenarios, collecting data and training for every possible gas would waste much resource. Open-set identification aims to actively reject samples from unknown classes by reducing the intra-class spacing and, thus, not misclassifying them as known classes. In this study, we propose a data preprocessing method to enhance the performance of closed-set recognition by augmenting the eigenvalues of each gas. We then implement the open-set recognition task for gases using an open-set recognition model. These methods contribute to improved recognition accuracy for gases and provide an effective means of handling unknown class samples. Experimental results show that our approach can identify unknown samples well while maintaining accuracy for available classes.

**INDEX TERMS** Electronic nose, open-set recognition, feature augmentation, machine learning.

## I. INTRODUCTION

The e-Nose [1] is composed of an array of multiple cross-sensitive sensors and an algorithm for recognition. It can detect different gas mixtures and identify complex samples with accuracy to simulation human senses. This technology has been applied in various fields, including food safety for detecting the quality of food [2], medicine for diagnosing lung diseases [3], and environmental detection for detecting land pollution [4]. In these applications, the technology is used to identify different substances through analysis of gases or gas mixtures. These techniques are typically based on spectroscopic analysis or other physical measurement methods, and utilize machine learning models to process the measurement results for identification of different substances. These techniques have played a significant role in applications such as detecting food quality, diagnosing lung diseases, and detecting land pollution.

CNNs are a type of deep learning algorithm known for their strong fitting abilities and ability to process various types of data, including 1D signals, 2D images, and 3D data. CNNs have been utilized in the classification of electronic nose data, and multilayer perceptrons (MLPs), a simple but effective classification method [5], are often employed. In recent years, CNNs have gained popularity in image recognition due to their strong feature extraction and generalization abilities. While deep learning algorithms may require more computing resources compared to traditional algorithms, the implementation of such methods has been made possible through the advancement of GPU technology.

Among them, LeNet-5, a classical convolutional neural network (CNN) algorithm, and its improved version have been applied for pattern recognition in electronic nose data, resulting in excellent performance [6]. This demonstrates the feasibility of utilizing relevant CNN algorithms for this application. However, the existing approach does not modify the features of electronic nose data, which are represented as one-dimensional sequences. To address this issue, we propose a method for preprocessing the data by converting each data

The associate editor coordinating the review of this manuscript and approving it for publication was Mehul S. Raval.

point into a two-dimensional sequence, and augmenting it with a practical algorithm that utilizes CNN structures commonly used in computer vision while preserving the original features (LeNet-5 [7], ResNet18 [8], VGG-16 [9], and InceptionNet [10]). This approach enables high accuracy in classification.

However, the number of gas species present in the real world far exceeds those present in the training dataset. In real-world scenarios, gas species that are not present in the training dataset may be encountered during the deployment of gas recognition instruments. These unknown gas samples represent classes that were not observed during the training phase. Since the models used are based on the closed-set assumption, the system will incorrectly assume that these gases belong to known classes, leading to labeling errors (as shown in FIGURE 1).

To address this issue, we attempt to introduce open-set Recognition (OSR) [11] as a solution, which can improve the system's usability in more realistic environments [12]. In open-set recognition, not all types of samples are included in the training dataset, meaning that there are samples in the test category that were not encountered during the training process. In open-set recognition systems, the goal is to maintain the classification accuracy of known classes while actively rejecting any unknown samples from the test set, as depicted in FIGURE 2.

To improve the robustness of Convolutional Neural Networks (CNNs) in open-set recognition (OSR) and maintain high accuracy in closed-set recognition (CSR), we applied the ARPL [13] and OLTR [14] algorithms for open-set recognition of gases. We also made relevant improvements to the OLTR algorithm to obtain the MOLTR model. Experimental results showed that our method was able to effectively identify finite classes in closed-set conditions. In addition, our approach was able to effectively recognize unknown gases under open-set conditions.

This research is significant in that it addresses the practical challenges encountered in recognizing gases in real-world environments. By utilizing preprocessing techniques and a limited number of gas features, the proposed method can improve the accuracy of gas identification, even in the presence of unknown gases. This has important implications for a range of industries, including those that rely on the identification of specific gases to ensure safety and compliance. This will further increase the practicality of gas identification and enhance the prospects for its applications.

In summary, the main contributions of this paper are as follows:

- In order to address the interference of unknown gases in practical operations on gas recognition, we have long been aware of the problem of open set recognition for gases.
- We propose an approach to improve the correct rate of gas identification by preprocessing sensor data under a limited number of gas features.

- We improved the accuracy of recognition by implementing an improvement to the existing open-set recognition model.
- Experiments on publicly available datasets show that our proposed related algorithm can effectively improve the accuracy of gas identification as well as identify unknown gases that are outside the sample range.

## II. RELATED WORK
We provide the details of the dataset utilized in our experiments, including the data preprocessing techniques applied, the comparison between the closed-set and open-set hypothesis CNN models, the evaluation criteria employed, and the specifications of the training devices.

### A. CLOSED-SET GAS RECOGNITION
Existing closed-set gas recognition methods can generally be divided into two categories: conventional and neural network methods. In this subsection, we will discuss the pros and cons of both approaches.

#### 1) CONVENTIONAL GAS RECOGNITION CLASS
Classic algorithms for traditional gas recognition involve extracting as many unique features as possible. Some examples of these algorithms include:

Support Vector Machines (SVMs) are linear classifiers used for binary classification in supervised learning. They aim to find the maximum margin hyperplane that separates the data points into two classes. The decision boundary is determined by solving for the learned samples. SVMs employ a hinge loss function to compute the empirical risk and augment a regularization term to the optimization system to minimize the structural risk. They are known for their sparsity and robustness. Least Squares Support Vector Machines (LSSVMs) have been used in [15] for gas mixture determination, attempting to optimize SVMs for different application contexts. The model is also employed for the prediction of the Water Quality Index (WQI) [16]. Calculation of WQI can be extremely complex and time-consuming, involving the calculation of sub-indices such as BOD and COD. However, by utilizing SVM and LS-SVM, WQI can be instantly predicted using the physical data directly measured by numerical methods with the same predictors, without the need for any sub-index calculations. The CMSVM [17] based on SVM can accurately predict the type of application flow through the Internet, including various anonymous networks.

Random Forests (RF) are ensembles of decision trees, where the output is determined by the majority vote of the individual trees [18]. RFs classify data using multiple classification trees, which can give importance scores for individual variables (or ''genes'') while evaluating the role played by each variable in the classification. RFs are known for their versatility and can handle a wide range of data types. In [19], the RF method was used to classify herbs with different characteristics. Additionally, the model can also achieve prediction of heart diseases [20]. In this study, the use of Random
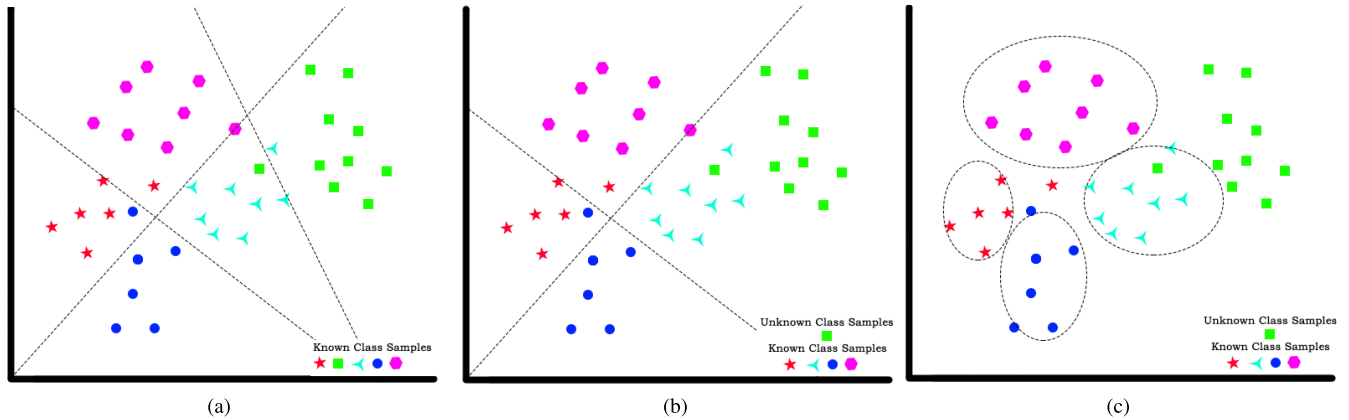
**FIGURE 1.** (a) the classical closed-set recognition model, (b) the use of the closed-set recognition model under the open-set assumption, which results in unknown samples being incorrectly classified to known class, and (c) the classical open-set recognition model, which achieves recognition of unknown samples by actively rejecting types that are not in the known samples.

Forest (RF) achieved an accuracy of 86.9% in predicting heart diseases.

K- Nearest Neighbor (KNN) is a simple and effective classification method that assigns a sample to a particular category based on the categories of its nearest neighbors in the feature space [21]. Specifically, if the majority of the K nearest neighboring samples belong to a particular category, the sample is classified as belonging to that category and is assumed to have similar characteristics to the samples in that category. This algorithm can be utilized to capture uncertainty and reflect the fluctuation range of electrical load [22], reducing computational cost and improving prediction efficiency and accuracy. Furthermore, the model can perform online public sentiment analysis [23] and has achieved a high level of accuracy.

### 2) NEURAL NETWORK GAS RECOGNITION CLASS

Neural networks are algorithms inspired by the functions of the human brain and consist of many highly interconnected processing units (neurons) working together to solve a specific problem. There are three common types of neural networks:

Recurrent Neural Networks (RNNs) are a type of neural network that are specifically designed to process sequential data and capture temporal and semantic information in data with sequential characteristics. This model effectively handles sequential data, as it can extract both temporal and semantic information from the data. Utilizing the capabilities of RNN, deep learning models have made breakthroughs in NLP fields such as speech recognition [24], language modeling [25], machine translation [26], and temporal analysis [27]. An RNN-based Long-Short-Term-Memory (LSTM) network was employed for estimation of different gas concentrations in the article [28].

Convolutional Neural Networks (CNNs) are a type of neural network architecture specifically designed for processing data with a grid-like structure. They use convolutional operations in at least one layer of the network rather than traditional

matrix multiplication. These networks are inspired by the biological mechanisms of visual perception and have several characteristics that make them useful for image processing tasks, including translation invariance, the use of convolution kernels, the ability to apply local information, the preservation of planar structural information, and the ability to perform both supervised and unsupervised learning. A series of studies have indicated that CNNs are highly effective in representing spatial patterns, as they can extract various vegetation attributes from remote sensing images [29], achieving high prediction accuracy and satisfying the growing demand for vegetation assessment and monitoring. In another study [6], a CNN-based LeNet-5 network was used for identifying CO, CH4, and their mixtures, with an accuracy of 98.67%.

Generative Adversarial Network (GAN), a generative model used in data generation, particularly for images. GAN consists of two parts: the generator and the discriminator. The generator tries to produce false targets that are indistinguishable from true targets, while the discriminator tries to differentiate between the two. Through an adversarial process, the generator and discriminator continuously improve their abilities, with the generator trying to produce better false targets and the discriminator trying to more accurately differentiate between true and false targets. While GAN has been primarily used for image generation [30], it has also been applied to text data, as in research [31], where it was used to address the imbalance problem in a sample dataset for dissolved gas analysis through data augmentation.

### 3) SUMMARIZATION AND COMPARISON

Traditionally, gas recognition models that achieve higher accuracy rates do so due to a smaller feature space and a reduced risk of overfitting when dealing with fewer species. However, as the number of species and the amount of data increases, traditional gas recognition models may encounter performance limitations. In contrast, neural network-based

gas recognition models often have the ability to utilize more parameters and achieve higher accuracy rates in these cases.

## B. OPEN-SET RECOGNITION METHODS IN COMPUTER VISION FIELD

In this subsection, we introduce two early methods of open-set recognition and two more advanced methods of open-set recognition. In the subsequent summary, we use these four algorithms to implement an application for gas open-set recognition.

SoftmaxThreshold (SoftMax-ST) method is a method for open-set identification that applies a threshold to the traditional closed-set identification method. This method allows for the recognition of unknown classes by setting a certain threshold. The SoftMax-ST method is a simple but effective approach for open-set recognition, as demonstrated in the paper [33] where it was used to achieve open-set iris recognition with an accuracy of 98.5%.Simultaneously, the method can achieve the prediction of regional vessel behavior [32] using historical AIS data, helping actively avoid collisions and enhance the maritime transportation system.

OpenMAX is a method for open-set identification that was proposed by Abhijit Bendale and Terrance E. Boult in the paper "Towards open-set Recognition" [34]. It involves the introduction of an OpenMAX layer that replaces several deep networks based on SoftMax probability thresholds in order to actively reject unknown classes. The method works by using Extreme Value Theory (EVT) to model the distance of activation vectors from the mean of each class and generate an updated penultimate vector that is used to identify unknown class test samples. A study has demonstrated the use of this technique to realize automatic target recognition (ATR) of synthetic aperture radar (SAR) images [35]. In this study, researchers employed this technique to classify open food powders and achieved an accuracy of 91.2% [36].

Open Long-Tailed Recognition (OLTR) Liu et al. [14] propose a comprehensive OLTR algorithm. By mapping images to feature spaces and using a learned metric that respects closed-world classification while acknowledging open-world novelty, the algorithm was able to handle the robustness of tail recognition through the use of dynamic meta-embeddings that are dynamically calibrated to visual memory. The embeddings were shown to be inversely proportional to their distance to the nearest center of mass, improving open-set recognition. The algorithm was tested on large-scale open-set recognition tasks and demonstrated improved performance compared to previous methods.

Adversarial Reciprocal Points Learning (ARPL) This method was proposed as a method for reducing empirical classification risk by Chen et al. [13]. This was achieved through the introduction of reciprocal points to model the potential open space of each known category in the feature space and the use of adversarial techniques between multiple known categories. Additionally, a new instantiated adversarial augmentation was introduced to estimate the unknown

distribution in the open space by generating a variety of confusion training samples from the known data and the adversarial swap points.

## C. DATASET

In this study, we utilize the Gas Sensor Array Drift Dataset[1] from the UCI Machine Learning Repository [37], which is composed of six gases: ammonia, ethylene, acetone, ethane, ethanol, and toluene. The dataset was collected over a period of 36 months, from January 2007 to February 2011, and features two types of information extracted from the response of Figaro metal oxide gas sensors, which are known to have a slower response to chemical compounds. These features include steady-state features (defined as the difference between the maximum resistance change and the baseline) and transient features (extracted using an exponential moving average (EMA) transformation). Then, by applying the EMA transformation to each of the 16 channels (sensors) within the pre-recorded time series, the sensor array response is mapped to a 128-dimensional feature vector. In order to alleviate the impact of sensor drift, the data is randomly shuffled and processed as individual sample space. To ensure consistent model training and testing, all models utilize the same validation and test sets. The detailed information of the dataset is shown in Table 1.

## D. DATA PRE-PROCESSING METHODS

In order to further analyze the dataset, we concatenated ten batches of data in the order of 1 to 10, forming a 2-dimensional matrix of $13910 \times 128$, then transformed it into a 3-dimensional matrix of $13910 \times 16x8$. We then used the equation in equation 1 to transform it into a matrix of $13910 \times 16x16$, as depicted in Figure 2. In this process, we introduced the concept of feature augmentation degree, which is calculated as:

$$A = \frac{S_{Exp}}{S_{Matrix}} \quad (1)$$

where $S_{Exp}$ denotes the area expanded using the algorithm and $S_{Matrix}$ denotes the total area after expansion.

In this study, we transformed the gas drift dataset, which is originally a 2D matrix of $13910 \times 128$, into a three-dimensional matrix of $13910 \times 16x8$, and then further transformed it into a matrix of $13910 \times 16x16$ using a specific algorithm Eq2. This process can be visualized as converting the dataset into $13910 \times 16x16$ images, which can then be used for classification with a modified deep learning model. In the later sections of this paper, we will demonstrate that this feature augmentation method can significantly improve sample recognition accuracy under closed-set conditions.

$$Y_{i,j,k+8} = \begin{cases} \dfrac{\sum\limits_{n=1}^{8} X_{i,j,n}}{8} & (k = 1) \\ \dfrac{\sum\limits_{n=k-1}^{k} X_{i,j,n}}{2} & (k \in [2, 8]) \end{cases} \quad (2)$$

[1] https://archive.ics.uci.edu/ml/datasets/gas+sensor+array+drift+dataset

**TABLE 1.** Sample details of batches.

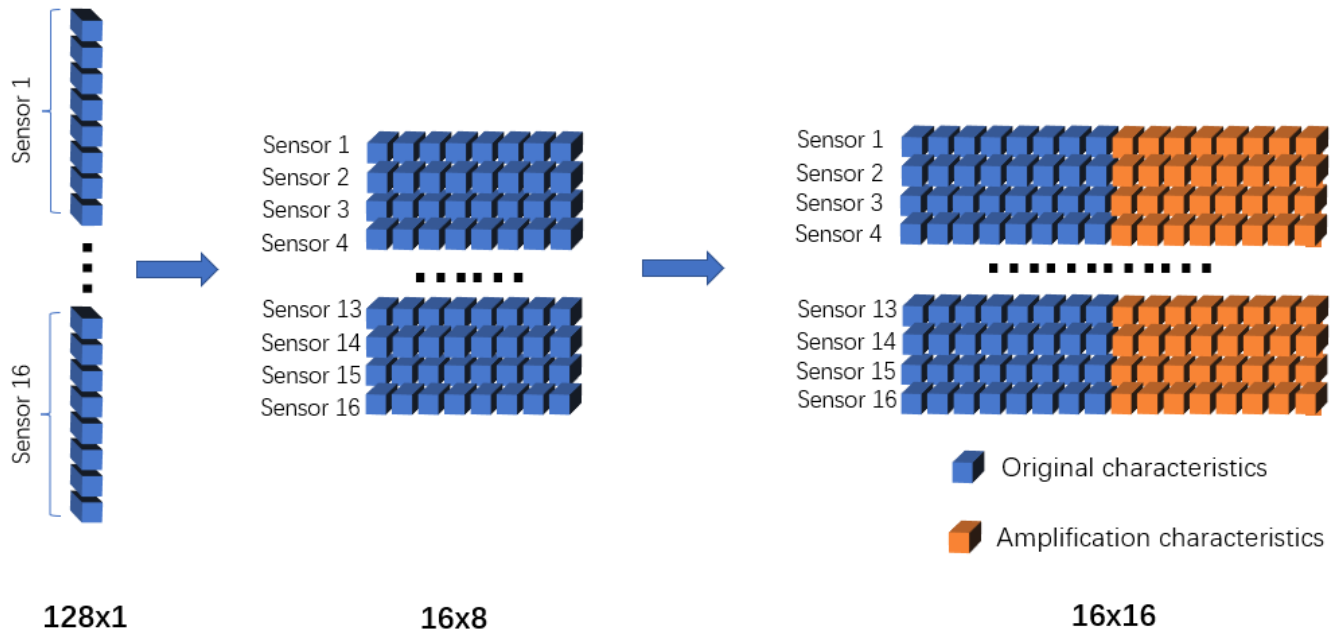| Batch ID | Month | Ammonia | Acetaldehyde | Acetone | Ethylene | Ethanol | Toluene | Total |
|----------|-------|---------|--------------|---------|----------|---------|---------|-------|
| Batch 1 | 1,2 | 83 | 30 | 70 | 98 | 90 | 74 | 445 |
| Batch 2 | 3,4,8,9,10 | 100 | 109 | 532 | 334 | 164 | 5 | 1244 |
| Batch 3 | 11,12,13 | 216 | 240 | 275 | 490 | 365 | 0 | 1586 |
| Batch 4 | 14,15 | 12 | 30 | 12 | 43 | 64 | 0 | 161 |
| Batch 5 | 16 | 20 | 46 | 63 | 40 | 28 | 0 | 197 |
| Batch 6 | 17,18,19,20 | 110 | 29 | 606 | 574 | 514 | 467 | 2300 |
| Batch 7 | 21 | 360 | 774 | 630 | 662 | 649 | 568 | 3613 |
| Batch 8 | 22,23 | 40 | 33 | 143 | 30 | 30 | 18 | 294 |
| Batch 9 | 24,30 | 100 | 75 | 78 | 55 | 61 | 101 | 470 |
| Batch 10 | 36 | 600 | 600 | 600 | 600 | 600 | 600 | 3600 |



**FIGURE 2.** Data pre-processing.

where $i$, $j$, $k$ denotes the three dimensions of the matrix, $X$ marks the value of the current position, $Y$ marks the value to be calculated.

## III. ARCHITECTURE OVERVIEW

In this subsection, we detail the similarities and differences between open-set recognition and closed-set recognition. The similarities and differences between their training and testing are shown in FIGURE 3.

### A. CLOSED-SET IDENTIFICATION MODEL

In this section, we study the closed-set recognition assumption, in which it is assumed that all classes of test samples are present in the training set. We compare the performance of several classical deep learning models, including AlexNet, LeNet-5, ResNet18, VGG-16, and InceptionNet, in this task. These models contain convolutional and pooling layers denoted as CBA (Conv2D, BatchNormalization, Activation) and PD (MaxPool, Dropout), respectively. To improve the accuracy of the results and adapt to the detection of
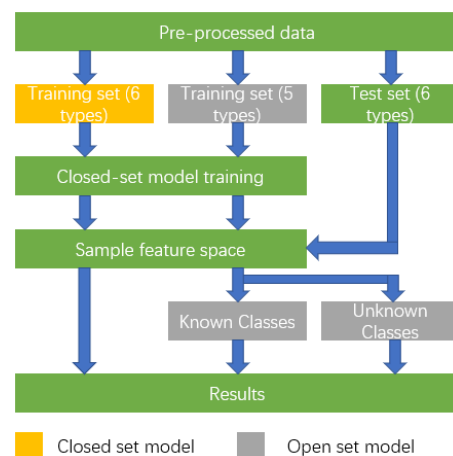


**FIGURE 3.** Open-set and closed-set models overview.

smaller input samples, we made various modifications to the models, including replacing the Sigmoid activation function with ReLu, which has been shown to have better

**TABLE 2.** Model hyperparameters.

| | Conv2d | | Max_pool2d |
|---|---|---|---|
| Stride | Padding | Kernel_size | Kernel_size |
| 1 | 2 | 5 | 2 |
| Activation | Optimizer | Learning Rate | Batch Size |
| PReLU | SGD | 0.01 | 128 |

**TABLE 3.** Configuration of computer.

| Devices | Model |
|---|---|
| System | Microsoft Windows 11 Pro 21H2 |
| CPU | Intel(R) Core(TM) i7-9750H CPU @ 2.60GHz |
| GPU | NVIDIA GeForce RTX 2060 6GB |
| Memory | 32.0 GB |

performance. The final fully connected layer of all four models in closed-set identification uses the SoftMax activation function (Eq3) to directly obtain the probability distribution of multiple prediction categories. However, this activation function is used for probability value generation and is forced to choose a type, which can lead to normalization issues and is a inherent property of closed-set recognition.

$$\sigma(z_i) = \frac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}} \quad for \ i = 1, 2, \ldots, K \quad (3)$$

where $\sigma(z_i)$: softmax function output for class i. $e^{z_i}$: exponential of $z_i$. $\sum_{j=1}^{K} e^{z_j}$: sum of exponential of $z_j$ for all classes j. $i = 1, 2, \ldots, K$: for all classes i in total of K classes.

### B. OPEN-SET IDENTIFICATION MODEL

In this study, we examine the performance of open-set recognition in the context of classifying test samples where not all classes are present in the training set. To simulate open-set identification, we reduce the number of categories in the training set and evaluate the performance of various open-set models using the same dataset. The openness of the dataset can be expressed using Eq4 [12]. To address the complexity of real-world situations and reduce the limitations of the model, we propose adding an unknown class to the closed-set model. That is, there are $N$ categories in the sample, but the total number of types is represented as $N + 1$, assuming the existence of unknown classes in the test set.

$$openness = 1 - \sqrt{\frac{2 \times |N_{train}|}{|N_{test}| + |N_{target}|}}. \quad (4)$$

where $N_{train}$ denotes the kind of training set, $N_{test}$ denotes the type of test set, $N_{target}$ denotes the type of target set.

We evaluated four algorithms for open-set recognition of gas: SoftMax-ST, OpenMAX, MOLTR, and ARPL. The results showed that each algorithm had its own strengths and weaknesses, and we will discuss these in more detail in the subsequent summary. The SoftMax activation function is not appropriate for open-set recognition tasks due to its inability to handle the distinct constraints between classes required in this type of recognition. Therefore, we modified the OLTR algorithm by incorporating pooling layers, resulting in the MOLTR model. In TABLE 2, we present the relevant hyperparameters of the model. These four algorithms were tested on a gas dataset with an openness of 0.5 and were used to achieve open-set recognition in this study.

### C. EVALUATION CRITERIA

To evaluate the closed-set model, we set the validation set and the training set in the same sample space and calculate the accuracy of the closed-set model $\mathcal{A}$ using Eq5.

$$\mathcal{A} = \frac{\sum_{i=1}^{C} (TP_i + TN_i)}{\sum_{i=1}^{C} (TP_i + TN_i + FP_i + FN_i)}. \quad (5)$$

where $C$ is the number of classes. $TP_i$ is true positives for class i. $TN_i$ is true negatives for class i. $FP_i$ is false positives for class i. $FN_i$ is false negatives for class i.

For the open-set recognition task, we consider known classes in the test set as positive and unknown classes as negative. Samples below a certain threshold are considered open-sets. In addition, for the open-set recognition model, the optimal case is to achieve the accuracy of closed-set recognition for known label classes. To measure the performance of our classifier, we utilize the area under the receiver operating characteristic curve (AUROC) [38]. A value of AUROC close to 1 indicates that the classifier can effectively classify positive and negative samples. Additionally, in order to accurately measure the performance of our model, we also calculated the proportion of known classes correctly classified in the test set. We refer to this metric as Accuracy(open-set)(ACC(os)) in the subsequent figures and tables in this section.

### D. TRAINING DEVICE INFORMATION

The models in this paper are trained, tested, and evaluated on a laptop. To improve the reproducibility of the article, we decided to disclose the device's leading software and hardware information, as shown in TABLE 3.
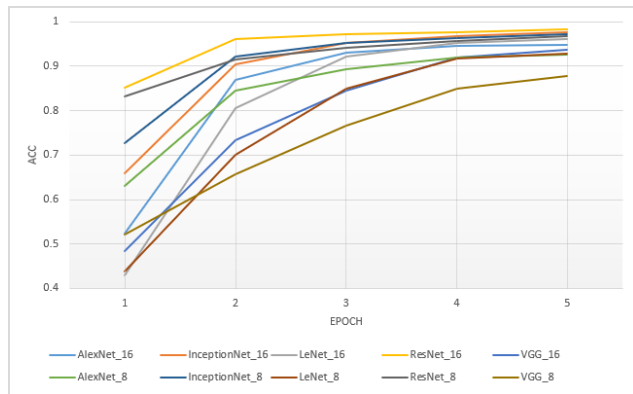
### IV. RESULTS AND DISCUSSION

This section presents the results of our experiments, along with a detailed analysis and discussion of the obtained outcomes.
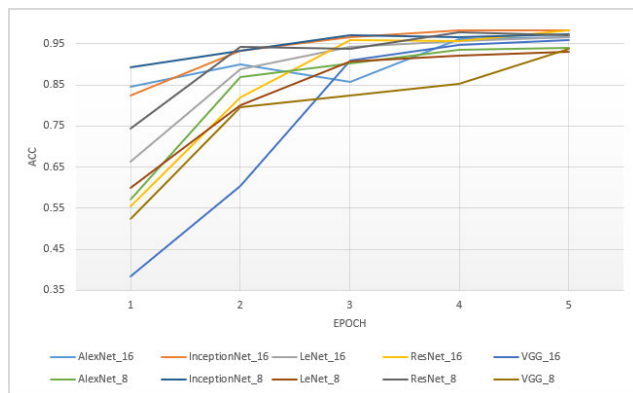
### A. CLOSED-SET EXPERIMENTS

In our closed-set experiments, we compared the effects of different data preprocessing methods on the experimental results. In order to preserve the original characteristics of the data, we transformed the data into a 16 × 8 matrix and randomly selected 80% of the samples (11128) to comprise the training set. The remaining 20% of the samples (2782) were utilized to create a validation set, which was used in conjunction with the training set from all batches during the training process. We also expanded the matrix into a 16 × 16 matrix and conducted a comparison test using both

**TABLE 4.** Closed-set experimental results.

| Matrix | AlexNet | | Inception | | LeNet | | ResNet | | VGG | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Train | Val | Train | Val | Train | Val | Train | Val | Train | Val |
| 16x16 | 0.9486 | 0.9651 | 0.9768 | 0.9832 | 0.9612 | 0.9672 | 0.9833 | 0.9834 | 0.9359 | 0.9585 |
| 16x8 | 0.9249 | 0.9405 | 0.9721 | 0.9721 | 0.9278 | 0.9312 | 0.9674 | 0.9701 | 0.8782 | 0.9384 |



(a) Training set



(b) Validation set

**FIGURE 4.** Accuracy changes during training.

**TABLE 5.** Open-set experimental results.

| | SoftMax | SoftMax-ST | OpenMax | MOLTR | OLTR | ARPL |
|---|---|---|---|---|---|---|
| ACC(os) | 0.58 | 0.808 | 0.84 | 0.795 | 0.774 | 0.995 |
| AUROC | 0.979 | 0.979 | 0.957 | 0.99 | 0.988 | 0.943 |



(a) ACC(os)



(b) AUROC

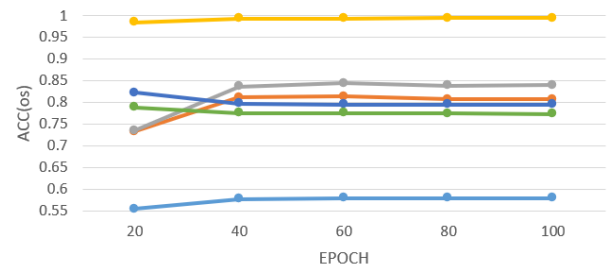**FIGURE 5.** ACC(os) and AUROC changes during training.

matrix samples. The change in accuracy of the validation set during the training process is shown in FIGURE 4, while the final training and validation set accuracies are shown in TABLE 4.
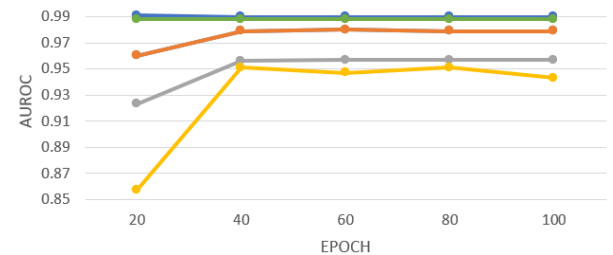
The experimental results indicate that the proposed data preprocessing algorithm effectively expands the feature space of the samples without changing the structure of CNNs, thereby improving the accuracy of gas classification. The degree of improvement in accuracy varies with the underlying recognition model, with a maximum improvement of approximately 3% for the LeNet model and a minimum improvement of approximately 1% for the Inception network. However, there is still room for optimization for the feature enhancement algorithm, and different data enhancement methods should be considered for different models. Future work may aim to determine a more appropriate feature enhancement algorithm and feature enhancement degree for specific gas classification models.

## B. OPEN-SET EXPERIMENTS

In the open-set experiments, we used data sets that have been amplified by data preprocessing methods and divided them into training and validation sets with a ratio of 4:1. During the experiments, we recorded and analyzed the changes in the ACC and AUROC values while using the SoftMax-ST, Open-MAX, MOLTR, ARPL, and SoftMax algorithms for open-set recognition, as shown in FIGURE 5. The final accuracies of the five models after 100 epochs are shown in TABLE 5.

The performance of four algorithms, SoftMax-ST, Open-MAX, MOLTR, ARPL, and the classic SoftMax algorithm, was compared in open-set recognition experiments. The results in Figure 5 indicate that the SoftMax algorithm is not suitable for open-set recognition. Additionally, the SoftMax-ST, OpenMAX, and MOLTR algorithms exhibit overfitting

**TABLE 6.** Ablation study results.

| Methods | OLTR+16x8 | MOLTR+16x8 | OLTR+16x16 | MOLTR+16x16 |
|---------|-----------|------------|------------|-------------|
| ACC(os) | 0.761 | 0.793 | 0.774 | 0.795 |
| AUROC | 0.985 | 0.985 | 0.988 | 0.99 |

phenomena. The ARPL algorithm had the highest accuracy of 99.5%, but the lowest AUROC value. In contrast, the MOLTR algorithm had a higher AUROC value of 0.985. Therefore, the MOLTR model proposed in this paper is capable of completing the gas open-set recognition task with relatively good performance.

### C. ABLATION STUDY

In this study, we compared the efficacy of our two proposed methods for improving the accuracy of open-set identification for gases through ablation experiments. The results, presented in TABLE 6, demonstrate that both methods are effective in improving the accuracy of the model.

## V. CONCLUSION

In this study, we propose a novel data augmentation method and demonstrate its effectiveness in improving the accuracy of open-set recognition for gases. Additionally, we apply image open-set recognition algorithms to the gas dataset for the task of open-set recognition. This approach significantly improves the performance of e-Nose in identifying unknown gases and classifying known gases. Our data augmentation algorithm, using several classical CNN models (AlexNet, Inception, LeNet, ResNet, VGG), effectively increases the accuracy of gas recognition without altering the original data features. We also demonstrate the recognition of unknown gases using various open-set recognition algorithms (SoftMax-ST, OpenMAX, MOLTR, ARPL), with MOLTR performing best in unknown gas recognition and ARPL performing best in classification of known gas types. While we have implemented data augmentation for the case of limited features, the results achieved with different CNN models vary. Therefore, the appropriate data augmentation algorithms for different models also differ. In this work, we provide a preliminary investigation of this approach and further targeted research is necessary for different targets. Additionally, we have implemented open-set identification for gases in this study, but this is only a preliminary examination. Due to the black-box nature of deep learning, it is challenging to optimize the model itself, however, new modules can be inserted to reduce inter-class variance while increasing intra-class variance. Future research should aim to enhance the recognition accuracy of open set recognition of unknown classes while maintaining the accuracy of known classes.

## REFERENCES

[1] T. Kuchmenko, R. Umarkhanov, and L. Lvova, "E-nose for the monitoring of plastics catalytic degradation through the released volatile organic compounds (VOCs) detection," *Sens. Actuators B, Chem.*, vol. 322, Nov. 2020, Art. no. 128585.

[2] S. Faal, M. Loghavi, and S. Kamgar, "Physicochemical properties of Iranian ziziphus honey and emerging approach for predicting them using electronic nose," *Measurement*, vol. 148, Dec. 2019, Art. no. 106936.

[3] W. Li, Z. Jia, D. Xie, K. Chen, J. Cui, and H. Liu, "Recognizing lung cancer using a homemade e-nose: A comprehensive study," *Comput. Biol. Med.*, vol. 120, May 2020, Art. no. 103706.

[4] A. Bieganowski, G. Józefaciuk, L. Bandura, Ł. Guz, G. Łagód, and W. Franus, "Evaluation of hydrocarbon soil pollution using e-nose," *Sensors*, vol. 18, no. 8, p. 2463, Jul. 2018.

[5] J. C. R. Gamboa, A. J. da Silva, I. C. S. Araujo, E. S. Albarracin, and C. M. Duran, "Validation of the rapid detection approach for enhancing the electronic nose systems performance, using different deep learning models and support vector machines," *Sens. Actuators B, Chem.*, vol. 327, Jan. 2021, Art. no. 128921.

[6] G. Wei, G. Li, J. Zhao, and A. He, "Development of a LeNet-5 gas identification CNN structure for electronic noses," *Sensors*, vol. 19, no. 1, p. 217, 2019.

[7] R. Liu, Y. Liu, Z. Wang, and H. Tian, "Research on face recognition technology based on an improved LeNet-5 system," in *Proc. Int. Seminar Comput. Sci. Eng. Technol. (SCSET)*, Jan. 2022, pp. 121–123.

[8] S. H. Karaddi and L. D. Sharma, "Automated multi-class classification of lung diseases from CXR-images using pre-trained convolutional neural networks," *Exp. Syst. Appl.*, vol. 211, Jan. 2023, Art. no. 118650.

[9] H. Nie, "Face expression classification using squeeze-excitation based VGG16 network," in *Proc. 2nd Int. Conf. Consum. Electron. Comput. Eng. (ICCECE)*, Jan. 2022, pp. 482–485.

[10] Z. Xuan-yu and D. Zuo-jie, "An improved method of clothing image classification based on CNN," *Int. J. Adv. Netw. Appl.*, vol. 12, no. 6, pp. 4742–4745, 2021.

[11] C. Liu, C. Yang, H.-B. Qin, X. Zhu, C.-L. Liu, and X.-C. Yin, "Towards open-set text recognition via label-to-prototype learning," *Pattern Recognit.*, vol. 134, Feb. 2023, Art. no. 109109.

[12] W. J. Scheirer, A. D. R. Rocha, A. Sapkota, and T. E. Boult, "Toward open set recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 7, pp. 1757–1772, Jul. 2012.

[13] G. Chen, P. Peng, X. Wang, and Y. Tian, "Adversarial reciprocal points learning for open set recognition," 2021, *arXiv:2103.00953*.

[14] Z. Liu, Z. Miao, X. Zhan, J. Wang, B. Gong, and S. X. Yu, "Large-scale long-tailed recognition in an open world," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 2537–2546.

[15] Ł. Lentka, J. M. Smulko, R. Ionescu, C. G. Granqvist, and L. B. Kish, "Determination of gas mixture components using fluctuation enhanced sensing and the LS-SVM regression algorithm," *Metrol. Meas. Syst.*, vol. 22, no. 3, pp. 341–350, 2015.

[16] W. C. Leong, A. Bahadori, J. Zhang, and Z. Ahmad, "Prediction of water quality index (WQI) using support vector machine (SVM) and least square-support vector machine (LS-SVM)," *Int. J. River Basin Manage.*, vol. 19, no. 2, pp. 149–156, Apr. 2021.

[17] S. Dong, "Multi class SVM algorithm with active learning for network traffic classification," *Exp. Syst. Appl.*, vol. 176, Aug. 2021, Art. no. 114885.

[18] L. Breiman, "Random forests," *Mach. Learn.*, vol. 45, no. 1, pp. 5–32, 2001.

[19] P. P. Kaur and S. Singh, "Random forest classifier used for modelling and classification of herbal plants considering different features using machine learning," in *Mobile Radio Communications and 5G Networks*. Singapore: Springer, 2022, pp. 83–94.

[20] M. Pal and S. Parija, "Prediction of heart diseases using random forest," *J. Phys., Conf.*, vol. 1817, no. 1, Mar. 2021, Art. no. 012009.

[21] N. S. Altman, "An introduction to kernel and nearest-neighbor nonparametric regression," *Amer. Statistician*, vol. 46, pp. 175–185, Aug. 1992.

[22] Y. Dong, X. Ma, and T. Fu, "Electrical load forecasting: A deep learning approach based on K-nearest neighbors," *Appl. Soft Comput.*, vol. 99, Feb. 2021, Art. no. 106900.

[23] A. R. Isnain, J. Supriyanto, and M. P. Kharisma, "Implementation of K-nearest neighbor (K-NN) algorithm for public sentiment analysis of online learning," *Indonesian J. Comput. Cybern. Systems*, vol. 15, no. 2, pp. 121–130, Apr. 2021.

[24] D. Baby, P. D'Alterio, and V. Mendelev, "Incremental learning for RNN-transducer based speech recognition models," in *Proc. Interspeech*, Sep. 2022, pp. 1–5.

[25] H. Zhang, Y.-C. Cheng, S. Kumar, W. R. Huang, M. Chen, and R. Mathews, "Capitalization normalization for language modeling with an accurate and efficient hierarchical RNN model," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP)*, May 2022, pp. 6097–6101.

[26] Y. Dong, "RNN neural network model for Chinese-Korean translation learning," *Secur. Commun. Netw.*, vol. 2022, pp. 1–13, May 2022.

[27] J. Wang, X. Li, J. Li, Q. Sun, and H. Wang, "NGCU: A new RNN model for time-series data prediction," *Big Data Res.*, vol. 27, Feb. 2022, Art. no. 100296.

[28] H. Bakiler and S. Güney, "Estimation of concentration values of different gases based on long short-term memory by using electronic nose," *Biomed. Signal Process. Control*, vol. 69, Aug. 2021, Art. no. 102908.

[29] T. Kattenborn, J. Leitloff, F. Schiefe, and S. Hinz, "Review on convolutional neural networks (CNN) in vegetation remote sensing," *ISPRS J. Photogramm. Remote Sens.*, vol. 173, Mar. 2021, pp. 24–49.

[30] Z. Ren, X. Y. Stella, and D. Whitney, "Controllable medical image generation via GAN," *J. Perceptual Imag.*, vol. 5, pp. 1–15, Mar. 2022.

[31] Y. Li, Y. Xu, X. Li, R. Li, J. Lin, and G. Zhang, "Addressing imbalance of sample datasets in dissolved gas analysis by data augmentation: Generative adversarial networks," *IET Gener., Transmiss. Distribution*, vol. 16, no. 22, pp. 4494–4504, Nov. 2022.

[32] B. Murray and L. P. Perera, "An AIS-based deep learning framework for regional ship behavior prediction," *Rel. Eng. Syst. Saf.*, vol. 215, Nov. 2021, Art. no. 107819.

[33] J. Sun, S. Zhao, S. Miao, X. Wang, and Y. Yu, "Open-set iris recognition based on deep learning," *IET Image Process.*, vol. 16, no. 9, pp. 2361–2372, Jul. 2022.

[34] A. Bendale and T. E. Boult, "Towards open set deep networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 1563–1572.

[35] E. Giusti, S. Ghio, A. H. Oveis, and M. Martorella, "Open set recognition in synthetic aperture radar using the openmax classifier," in *Proc. IEEE Radar Conf. (RadarConf)*, Mar. 2022, pp. 1–6.

[36] H. You, H. Kim, D.-K. Joo, S. M. Lee, J. Kim, and S. Choi, "Classification of food powders with open set using portable VIS-NIR spectrometer," in *Proc. Int. Conf. Artif. Intell. Inf. Commun. (ICAIIC)*, Feb. 2019, pp. 423–426.

[37] A. Vergara, S. Vembu, T. Ayhan, M. A. Ryan, M. L. Homer, and R. Huerta, "Chemical gas sensor drift compensation using classifier ensembles," *Sens. Actuators B, Chem.*, vol. 166, pp. 320–329, May 2012.

[38] H. Zhang, A. Li, J. Guo, and Y. Guo, "Hybrid models for open set recognition," in *Proc. Eur. Conf. Comput. Vis.* Cham, Switzerland: Springer, 2020, pp. 102–117.

**YE ZHU** is currently pursuing the B.Eng. degree with the School of Computer Science and Technology, Anhui University of Technology, Maanshan, China. His research interests include machine learning and the Internet of Things.

**JINGYA WANG** is currently pursuing the B.Eng. degree with the School of Computer Science and Technology, Anhui University of Technology, Maanshan, China. Her research interests include image processing and machine learning.

● ● ●