Volume 3, Issue 2, 2025

Print ISSN: 2959-9865 Online ISSN: 2959-9873

# WORLD JOURNAL OF ENGINEERING RESEARCH



**Copyright© Upubscience Publisher** 

# World Journal of Engineering Research

Volume 3, Issue 2, 2025



**Published by Upubscience Publisher** 

#### Copyright<sup>©</sup> The Authors

Upubscience Publisher adheres to the principles of Creative Commons, meaning that we do not claim copyright of the work we publish. We only ask people using one of our publications to respect the integrity of the work and to refer to the original location, title and author(s).

Copyright on any article is retained by the author(s) under the Creative Commons Attribution license, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Authors grant us a license to publish the article and identify us as the original publisher.

Authors also grant any third party the right to use, distribute and reproduce the article in any medium, provided the original work is properly cited.

World Journal of Engineering Research Print ISSN: 2959-9865 Online ISSN: 2959-9873 Email: info@upubscience.com Website: http://www.upubscience.com/

# **Table of Content**

REPRODUCE LONG-TAIL CURRENT HYSTERESIS IN PEROVSKITE SOLAR CELLS BASED ON AN IONIC CAPACITOR MODELLED BY A VOLTAGE-RELATED INITIAL CURRENT TIMES THE LINEAR COMBINATION OF THREE EXPONENTIAL DECAY TERMS Biao Peng, RongXin Wu, YueWen Chen, MuYun Li, YingFeng Li*	1-7
DESIGN OF WATER PIPELINE MONITORING SYSTEM BASED ON MULTI-SOURCE INFORMATION FUSION LiTong Ma, Bo Ma*	8-14
OPTIMIZATION OF COAL MINE ROCKBURST EARLY WARNING SYSTEM JiaQi Wu*, YunMin Tian, TianLe Xiong, JunYao Hou, YunFeng Luo, Hao Chen	15-20
DESIGN OF DUST MONITORING SYSTEM FOR PRODUCTION WORKSHOP MingMao Gong*, SiZhe Zheng, SiYan Xu	21-26
ELUCIDATING THE DRIVING FACTORS IN SULFUR TRIOXIDE FORMATION UNDER SIMULATED ACTUAL ULTRA-LOW EMISSION PROCESS ZePeng Li, Yasser M. A. Mohamed*, YingHui Han*	27-31
DIAGNOSIS FOR WHEELSET OUT-OF-ROUNDNESS OF METRO VEHICLE USING VMD COMBINED WITH OPTIMIZED MCKD XiChun Luo, HaoRan Hu*	32-40
<b>SWM: AN OPTIMIZED DIFFERENTIAL EQUATION MODEL FOR STAIR WEAR</b> GuanYu Xu, SongHe Wang, ChaoJing Zhang, GaoHua Kong*	41-49
INTELLIGENT FAULT DIAGNOSIS OF ROLLING BEARINGS BASED ON VMD-CNN-TRANSFORMER JinYuan Hu	50-56
GRAPH AUTOENCODERS: A SURVEY LiNing Yuan	57-62
IMPROVING SMALL FIRE TARGET DETECTION IN UAV IMAGERY: AN ENHANCED RT-DETR WITH MULTI-SCALE FUSION AND EXPERT ROUTING ZhiCheng Zhang	63-74

## REPRODUCE LONG-TAIL CURRENT HYSTERESIS IN PEROVSKITE SOLAR CELLS BASED ON AN IONIC CAPACITOR MODELLED BY A VOLTAGE-RELATED INITIAL CURRENT TIMES THE LINEAR COMBINATION OF THREE EXPONENTIAL DECAY TERMS

Biao Peng, RongXin Wu, YueWen Chen, MuYun Li, YingFeng Li\*

State Key Laboratory of Alternate Electrical Power System with Renewable Energy Sources, School of Renewable Energy, North China Electric Power University, Beijing 102206, China. Corresponding Author: YingFeng Li, Email: livingfeng@ncepu.edu.cn

Abstract: It has been reported that the current-voltage (J-V) hysteresis loop of perovskite solar cells (PSCs) could be reproduced by incorporating extra resistances or capacitors in the equivalent circuits of PSCs. However, the exponential decay long-tail current in the hysteresis phenomenon, lasting about 2-5 seconds, remains inadequately modeled, yet it is crucial for maximum power point tracking in PSCs. We propose an ionic capacitor model to describe the impact of ion migration on the current of PSCs, which is composed of a voltage-related initial current multiplies by the linear combination of three exponential decay terms over time. The initial current term is formulated as the product of the voltage step size and a voltage-related conductance. The three exponential terms, each associated with a specific time constants  $\tau$ , correspond to the migration of electron, iodine ions, and other ions with lower mobilities, respectively. Based on this model, an equivalent circuit for PSCs is constructed, and corresponding parameters were numerical fitted based on available *J-V* data and current-time response curves. Numerical simulations demonstrate that the proposed model accurately reproduces both the *J-V* hysteresis loop and the exponential decay long-tail current. This work lays the foundation for the development of MPPT tracking algorithms tailored for PSCs.

Keywords: Perovskite solar cells; Hysteresis; Long-tail current; Ionic capacitor

#### **1 INTRODUTION**

Recently, perovskite solar cells (PSCs) have garnered significant attention in the research community[1-4]. These cells can be manufactured using solution-based methods, offering low production and material costs, which gives them advantages over traditional crystalline silicon solar cells. To date, the photoelectric conversion efficiencies of PSCs have reached 27.0%[5]. With the industrialization and practical applications of PSCs in power generation, maximum power point tracking (MPPT) has become indispensable. Common MPPT algorithms include the perturb and observe (P&O) method and the incremental conductance method (INC)[6-7]. Both algorithms require continuous adjustments to the battery's output voltage to locate the maximum power point, and their determination of the maximum power point's location is based on the battery's J-V characteristic curve[8].

Actually, adjusting the output voltage during the MPPT process in solar cell is equivalent to altering the electric field across the cell. In addition to conventional electrons, various ions in perovskite solar cells (PSCs) can also migrate in response to an external electric field. These migrated ions under the electric field will accumulate or dissipate at the electron transport layer (ETL)/perovskite interface and the hole transport layer (HTL)/perovskite interface. These ion accumulation and dissipation processes in response to the changing electric field are equivalent to introducing a series capacitor, characterized by a long charge-discharge time constant, within the PSC. This complex ion migration behavior is just the underlying mechanism behind the well-known hysteresis effect[10-11]. The hysteresis effect of PSCs shows two key characteristics. The first one is the non-overlapping current-voltage (J-V) curves during the forward and reverse scan of the PSCs, usually called the hysteresis loop, as illustrated in Figure 1(a). The other characteristic is the time-dependent decay of the dynamic non-steady-state photocurrent during stepwise voltage scanning, with a response time ranging from 2-5 seconds[9], which is named "long-tail current" in this work, as shown in Figure 1(b) and 1(c). This hysteresis effect renders conventional MPPT techniques unsuitable for PSCs[8].



Figure 1 (a) Hysteresis Loops of PSC under Different Scan Rates; (b) Time-dependent Photocurrent Response under Stepwise Reverse Scan with 100 mV Step Size and 5 s Step Time; (c) Decay of Dynamic Non-Steady-State Photocurrent with Time during the Stepwise Reverse Scan, the Inset Figure is the Normalized Exponential Decay Current Curves; (d) The Initial Current Extracted from (b) Changes with Voltage. Reprinted (adapted) with Permission from ref 18. Copyright © 2015, American Chemical Society

A critical step in constructing new MPPT algorithms for PSCs is to develop a model that can quantitatively describe the impact of ion migration on the hysteresis effect. In previous reports, the ion migration behaviors were modeled using a conventional electronic capacitor in parallel with an appropriate resistance[12-13]. Seki et al. have integrated such a parallel resistor-capacitor (RC) module (fixed resistor in parallel with fixed capacitor) in series into an equivalent circuit model of PSCs[14], and reproduced the hysteresis loop in both forward and reverse *J-V* sweep curves. However, such a simplified model effectively abstracts various carriers, including electrons and ions, in PSCs into a kind of "pseudo carrier" with lower carrier mobility. As a consequence, the long-tail current in PSCs, characterized by multiple segments corresponding to the migration of different ions, cannot be accurately reproduced[15].

Aware of this defect, J. B et al. have studied the long-tail current response of PSCs in response to a small voltage perturbation of 10 mV over a steady-state of 1 V[16]. They have divided the long-tail current into three distinct segments and modeled each segment using different RC/RL circuits. Based on impedance spectroscopy measurements, they fitted the parameters in the RC/RL circuits. Such constructed model reproduced the long-tail current curve under a given voltage very well. However, this model is also difficult to directly apply in MPPT algorithm design. On one hand, the parameters extraction rely on additional impedance spectroscopy measurements, on the other hand, it lacks to establish an analytical relationship between the model parameters and the voltage.

#### 2 IONIC CAPACITOR MODEL

For crystalline silicon solar cells, changing the voltage will disrupt the equilibrium of the built-in field, driving charge carriers (electrons and holes) to flow through the external circuit, as illustrated in Figure 2a, thereby generating current. In contrast, PSCs exhibit a different mechanism. Under the voltage changing not only electrons but also various ions migrate toward the respective transport layers. Due to the selective permeability of the ETL/HTL, these ions cannot fully penetrate the interfaces and instead accumulate at the surface, as illustrated in Figure 2b. These accumulated ions can act as an equivalent resistance which affects the electron transport. However, since the migration processes take time, it can be resembled as the charging and discharging behavior of a capacitor. Therefore, we use an ionic capacitor to describe the impacts of ions migration behavior on the electron transport.



Figure 2 The Electric Field and Carrier Dynamics in (a) Crystalline Silicon Solar Cells and (b) PSCs

To accurately reflect the long-tail current effect, referring to the model proposed by J.B. et al. [16], we have modified the conventional electronic capacitor into a linear combination of three different capacitors. The modeling process is as follows. According to the Kirchhoff laws, the charging or discharging current of a capacitor can be expressed by equation 1,

$$I_c(t) = (U - \frac{Q(t)}{C})/R \tag{1}$$

where U is the apply voltage, Q(t) is the charge on the conductor at time t, C denotes the capacitance, R is the external resistance. By substituting into equation (1) and solving the resulting first-order differential equation, we obtained

$$Q(t) = CU - \xi C e^{\frac{1}{RC}}$$
<sup>(2)</sup>

The integration constant can be determined by setting t=0, After substituting into equation (2), the charging or discharging current of the capacitor can be derived by differentiating Q(t) with respect to t,

$$I_{c}(t) = \frac{dQ(t)}{dt} = \frac{Q(0)/C - U}{R}e^{-\frac{t}{RC}} = \frac{U_{0} - U}{R}e^{-\frac{t}{RC}}$$
(3)

where represents the voltage on the capacitor before adjustment. can be represented as  $\Delta U$ , which represents the voltage adjustment step size. The reciprocal of the external resistance *R* can be expressed as the conductivity  $\sigma$ . Meanwhile, the charging or discharging time constant of the capacitor with external resistance *R*, given by *RC*, can be denoted by  $\tau$ . The time response equation of a given capacitor can be written as

$$I_c(t) = \Delta U \sigma e^{-t/\tau} = I_0 e^{-t/\tau} \tag{4}$$

The first term,  $I_0 = \Delta U \sigma$ , describes the initial current value on the capacitor after a voltage adjustment; and the second term,  $e^{-t/\tau}$ , corresponds to decay speed of the capacitor current over time.

As discussed above, the capacitor current in PSCs is not solely due to electron migration but is also influenced by the migration of various ions, including fast-moving species such as I and slower-moving ions like  $MA^+$ . Therefore, we try to interpret the capacitor current in PSCs as the joint contribution of three distinct carriers,

$$I_{c}(t) = I_{0}^{1} e^{-t/\tau_{1}} + I_{0}^{2} e^{-t/\tau_{2}} + I_{0}^{3} e^{-t/\tau_{3}}$$
(5)

In this equation, the three terms represent the contributions of electron migration, fast-moving ion migration, and slow-moving ion migration, respectively.  $I_0^1$ ,  $I_0^2$ , and denote the initial currents contributed by the three types of charge carriers, while  $\tau_1$ ,  $\tau_2$ , and represent the time constants of the three sub-capacitors. However, in practical measurements, directly obtaining the initial currents associated with individual charge carriers is unfeasible. Instead, an overall initial current of the capacitor in PSCs can be conveniently derived from time-dependent photocurrent measurements. Considering this fact, we reformulate equation 5 by replacing the individual initial currents with an overall initial current,  $I_0$ , while representing the contributions of different carriers through weighting factors  $c_1$ ,  $c_2$ , and  $c_3$ ,

$$I_{c}(t) = I_{0}(c_{1}e^{-t/\tau_{1}} + c_{2}e^{-t/\tau_{2}} + c_{3}e^{-t/\tau_{3}})$$
(6)

The values of  $I_0$  under various voltages were extracted from Figure 1(b) and plotted in Figure 1(d). It can be observed that  $I_0$  increases first and then decreases with the increase of the voltage. This is because  $I_0 = \Delta U/R$ , while *R*, which characterizes the influence of ion migrations on charge trapping and defect-induced recombination<sup>17</sup>, should be voltage-dependent. As the reciprocal of *R*, the conductivity should be also voltage-dependent. However, due to the complexity of the origins of *R*, it is challenging to establish a physically meaningful equation to describe the relationship between and *U*. Here, we adopt a mathematically feasible approach, namely a polynomial function, to express this relationship,

$$I_0 = \Delta U\sigma(U) = \Delta U(aU^4 + bU^3 + cU^2 + dU + e)$$
<sup>(7)</sup>

#### **3** NUMERICAL FITTING AND MODEL CONSTRUCTION

To determine the parameters in equation 7, we extracted data for  $I_0$  at various voltages (25°C, 1000W/m<sup>2</sup>) from the

Volume 3, Issue 2, Pp 1-7, 2025

measured data of Bo Chen et al.[18], as shown in Figure 1(d) and performed regression fitting to determine the parameters in equation 7. The obtained expression of  $I_0$  is shown in equation 7, with *R*-square =0.964.

$$I_0 = 0.1 * (-320.9U^4 + 242.9U^3 + 18.25U^2 + 10720U + 2.428)$$
(8)

The three exponential current decay terms in equation 6 represent the contributions of electron migration, fast-moving ion migration, and slow-moving ion migration, respectively. The first term is primarily responsible for the rapid decay in the initial stage of the long-tail response in PSCs[16, 19-20] within small than one second. It has been widely confirmed that iodide ions are the majority of migrating ions in PSCs, and the local electric field established by their migration results in a current decay over a longer time scale[17]. The second term takes account for the influence of the local electric field formed by iodide migration on the charge capture behaviors in PSCs. There are also some other slow-moving ions, like MA<sup>+</sup> or bound iodide ions, in PSCs[20]. To capture the effects of these slow-moving ions, we introduced the third term.

Then, regarding the long-tail exponential current decay term, the normalized curves in Figure 1c show that the differences in decay curves under varying voltages are not particularly significant and do not critically impact the manifestation of the hysteresis effect. For the convenience of fitting, we selected the representative current response curve corresponding to the reverse voltage scan from 0.4 V to 0.3 V as the experimental data for regression fitting. The obtained decay of capacitor current over time, with *R-square* =0.999, is shown in equation 9. The exponential items with different time constants can be seen as three sub-capacitors, and the current-time curves of them were plotted in Figure 4(a). The complete charge and discharge times of the three sub-capacitances can be calculated by  $5 \times \tau$ , i.e., 0.1810 s, 1.3455 s, and 4.6350 s, respectively. In equation 9, the value of  $c_2$  is much larger than  $c_3$ , indicating that the long-tail current should be mainly contributed by the migration of iodine ions. It can be also observed from Figure 4(a) that the actual output current of a PSC is strongly determined by the step interval time  $\Delta t$ , which represents the duration after the voltage adjustment before the cell's current is measured. In the scanning of solar cell J-V curve, the value of step voltage  $\Delta V (0.05 \text{ V} \text{ here})$  is usually fixed. Therefore,  $\Delta t$  is inversely proportional to the scanning speed  $\Delta V/\Delta t$ . Consequently, hysteresis in the *J-V* curves of PSCs will strongly depend on the scanning speed, as illustrated in Figure 1(a).

$$I_c(t) = I_0(0.7898 * e^{-\frac{1}{0.0362}} + 0.1467 * e^{-\frac{1}{0.2691}} + 0.0637 * e^{-\frac{1}{0.9270}})$$
(9)

According to equations 6 and 7 established in this paper, and the ionic capacitor model, the output current of PSC can be written as equation 10. This current model can be divided into two parts: the traditional equivalent circuit part, which can reflect the traditional steady photogenic current, and an extra capacitor part, which is proposed here to reflect the effect of the ion migration. For the carriers that respond rapidly to voltage changes and can quickly reach equilibrium after the application of a voltage change, we refer to the current output from the traditional equivalent circuit part,  $I_{trad}$ , as the steady-state current. In contrast, due to the slower response of ions to voltage changes, the current output from the extra capacitor part,  $I_{c_2}$  is referred to as the non-steady-state photocurrent.

$$I = I_{trad} - I_c = I_{ph} - I_s \left( e^{\frac{q(U+I_{trad}R_s)}{nK_BT}} - 1 \right) - \frac{U+I_{trad}R_s}{R_{sh}} - I_c$$
(10)

where  $I_{ph}$  is the photogenerated current,  $R_s$  is the series resistance,  $R_{sh}$  is the shunt resistance, the reverse saturation current can be written as  $I_s = (I_{ph} - V_{oc}/R_{sh}) / exp (qV_{oc}/(nk_BT))$ , and the ideal factor is fixed to be n = 1.1.



Figure 3 (a) Structure Diagram of PSC and the Positions of Ions at Different Time Points; (b) The Equivalent Circuit Model of a PSC with Hysteresis Effect; (c) PSC Matlab/Simulink Electrical Model

After the current model of PSC had been constructed, regression fittings were carried out to determine the parameters in

equation 10. The parameters that need to be determined are  $V_{oc}$ ,  $I_{ph}$ ,  $R_s$ , and  $R_{sh}$ , respectively. These parameters can be obtained by performing regression fitting on equation 10 using the stable *J-V* curve of a PSC, i.e., the current at a given voltage is measured after an adequate period. We derive the stable *J-V* curve of a PSC by averaging the currents in the forward and reverse scanning curves at every specific voltage, with a scanning speed of 200 mV/s, in Figure 1(a). Here, it should be noted that the *J-V* curve of PSC has a higher curvature than that of crystalline silicon cells. This is because, in PSCs,  $R_{sh}$  should be dependent on *U* as changes in operating voltage may cause extra leakage current related to ion migrations. Therefore, we tried to construct the relation between  $R_{sh}$  and *U* using a simple linear function  $R_{sh}=gU+h$ . After the regression fitting, the parameters obtained are  $V_{oc}=1.0118$ ,  $I_{ph}=20.0775$ ,  $R_{sh}=-0.151U+0.2284$ , and  $R_s=0.0054$ , respectively, with *R-square* =0.994.

Based on equations 10 and the structure of PSCs shown in Figure 3(a) we drew the equivalent circuit diagram of the PSCs, as shown in Figure 3(b), and we developed the electrical model of PSC in Matlab/Simulink, as shown in Figure 3(c), to simulate its *J-V* curves and long-tail current phenomena. The cyan-shaded area represents the  $I_{trad}$  component, while the yellow-shaded area represents the  $I_c$  component. For the  $I_{trad}$  component, we used a traditional crystalline silicon solar cell model to describe its dependence on scanning voltages. Then  $I_{trad}$  is used as an input parameter for the  $I_c$  component.  $I_c$  is calculated based on equations 6 and 7. At each scanning step, the scanning direction will be judged by comparing the values of the present input voltage U and the previous voltage  $U_{pre}$  (output a variable 'towards'); and if a voltage change has occurred ( $U \neq U_{pre}$ ), a time reset signal will be generated which will cause the variable *t* in equation 6 to be reset to zero.

#### 4 RESULT AND DISCUSSION

To verify the reliability and accuracy of the equivalent circuit model established in this work, at first, the hysteresis loop of the PSC J-V curve is simulated by Matlab/Simulink, using the same scanning speed as that of the experimental measurement. Corresponding results were shown in Figure 4(b), and for comparison, the measured results at scanning speeds of 200 mV/s and 1000 mV/s were given. It can be observed that the simulation results can perfectly reproduce the measured hysteresis loops at various scanning speeds. Then, the long-tail current in the hysteresis effect of PSCs was reproduced, as shown in Figure 4(c), which was obtained in reverse scanning with a stepwise voltage of 0.1 V and a duration of 5 seconds. The reproductions of both the hysteresis loops and long-tail currents together confirm the reliability and accuracy of the proposed equivalent circuit model with introduced exponential decay items.



Figure 4 (a) Current Decay Curve of the Capacitor and Three Sub-Capacitances with Different Time Constants in PSC,  $\tau_1=0.0362$ ,  $\tau_2=0.2691$ ,  $\tau_3=0.927$ ; (b) Hysteresis Loops of PSC at Different Scan Rates Simulated by Matlab/Simulink and the Experiment Data; (c) Simulation and Experiment Date of the Time-Dependent Photocurrent Response under Stepwise Reverse Scan with 100 mV Step Size and 5 s Step Time

By adjusting the parameters in the proposed model, some other typical hysteresis curves of different PSCs can also be reproduced[21]. For example, by adjusting  $I_0$  at U = 0 V from 0 to 1.5, a special hysteresis curve was reproduced, where the reverse scan current consistently exceeded the forward scan current, as shown in Figure 5(a). According to the proposed model, such characteristic hysteresis curve should appear in high-quality PSCs with neglectable series resistance  $R_s$ . This is because, in these cases, when the voltage is adjusted from U=0 to  $\Delta U$  (one scanning step),  $I_0$  can be approximated as  $I_0 \approx \Delta U/R_{load}$ , where  $R_{load}$  can be evaluated by  $\Delta U/J_{sc}$ ; therefore,  $I_0$  should be greater than 0 when U=0 V.

In addition, in some cases the hysteresis curve may exhibit a phenomenon of current protrusion during reverse scanning, this phenomenon causes the short-circuit current to no longer be the maximum current that PSCs can generate. As shown in Figure 5(b), such a hysteresis curve can also be reproduced by increasing the values of  $I_0$  around the protrusion position, specifically from 9 to 16. Such an adjustment is justifiable. In the case of a PSC possessing a moderate  $R_s$ , the values of  $I_0 = U/(R_s + R_{eq})$  are expected to initially ascend and subsequently descend with the increment of voltage. This phenomenon can be evidenced by the measured results presented in Figure 1 (b).



Figure 5 (a) Simulated Hysteresis Curve of a PSC Where the Reverse Scan Current Consistently Exceeds the Forward Scan Current; (b) Simulated Hysteresis Curve Showing Current Protrusion during the Reverse Scan

#### 5 CONCLUSIONS

In conclusion, we propose an ionic capacitor model to describe the impact of ion migration on the current of PSCs. A voltage-related polynomial fitting approach was used to capture the complex relationship between the initial current and voltage; and the influence of ion migration on the current decay was represented as a linear combination of three sub-capacitors corresponding to electron, iodine ions, and other ions with lower mobilities. This model provides a more accurate representation of the ionic capacitor in PSCs. Building upon this framework, we propose an improved equivalent circuit model for PSCs by incorporating an additional ionic capacitance element into the conventional equivalent circuit. The model parameters can extract through regression fitting of available current J-V characteristics and current-time response curves. An electrical model of the PSC was then established on the Matlab/Simulink platform, and numerical simulations were performed using the proposed equivalent circuit model. The results demonstrate that the model not only accurately reproduces the hysteresis loop of PSCs, but also effectively captures their long-tail current characteristics; moreover, the model replicates and elucidates the hysteresis curves of two other typical PSCs, with key performance parameters exhibiting high consistency with experimental data. These insights are critically important for the optimization and development of PSCs based on MPPT algorithms.

#### **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

#### FUNDING

This work was supported by National Natural Science Foundation of China (Grant No. 52072121).

#### REFERENCES

- [1] Wu S, Liu M, Jen AKY. Prospects and challenges for perovskite-organic tandem solar cells. Joule, 2023, 7(3): 484-502.
- [2] Yang C, Hu W, Liu J, et al. Achievements, challenges, and future prospects for industrialization of perovskite solar cells. Light: Science & Applications, 2024, 13(1).
- [3] Aydin E, Allen TG, De Bastiani M, et al. Pathways toward commercial perovskite/silicon tandem photovoltaics. Science, 2024, 383(6679).
- [4] Du S, Huang H, Lan Z, et al. Inhibiting perovskite decomposition by a creeper-inspired strategy enables efficient and stable perovskite solar cells. Nature Communications, 2024, 15(1).
- [5] National Renewable Energy Laboratory. 2025. https://www.nrel.gov/pv/cell-efficiency.html.
- [6] Sarvi M, Azadian A. A comprehensive review and classified comparison of MPPT algorithms in PV systems. Energy Systems, 2022, 13(2): 281-320.
- [7] Verma D, Nema S, Shandilya AM, et al. Maximum power point tracking (MPPT) techniques: Recapitulation in solar photovoltaic systems," Renewable and Sustainable Energy Reviews, 2016, 54: 1018-1034.
- [8] B S, SKP D, S S, et al. Steady Output and Fast Tracking MPPT (SOFT-MPPT) for P&O and InC Algorithms," IEEE Trans. Sustain. Energy, 2021, 12(1): 293-302.

- [9] Kang DH, Park NG. On the Current–Voltage Hysteresis in Perovskite Solar Cells: Dependence on Perovskite Composition and Methods to Remove Hysteresis. Advanced Materials, 2019, 31(34): 1805214 (2019).
- [10] Kim H, Park N. Parameters Affecting I-V Hysteresis of CH3NH3PbI3 Perovskite Solar Cells: Effects of Perovskite Crystal Size and Mesoporous TiO2 Layer. The Journal of Physical Chemistry Letters, 2014, 5(17): 2927-2934.
- [11] Tress W, Marinova N, Moehl T, et al. Understanding the rate-dependent J-V hysteresis, slow time component, and aging in CH3NH3PbI3 perovskite solar cells: the role of a compensated electric field. Energy & Environmental Science, 2015, 8(3): 995-1004.
- [12] Anghel DV, Nemnes GA, Pintilie I, et al. Modelling J-V hysteresis in perovskite solar cells induced by voltage poling. Physica Scripta, 2019, 94 (12): 125809.
- [13] Tang S, Yan J, Chen L, et al. Circuit modeling and analysis of hysteresis effect of perovskite photovoltaic cells. Solar Energy Materials and Solar Cells, 2024, 278: 113182.
- [14] Seki K. Equivalent circuit representation of hysteresis in solar cells that considers interface charge accumulation: Potential cause of hysteresis in perovskite solar cells. Applied Physics Letters, 2016, 109(3): 33905.
- [15] Gottesman R, Haltzi E, Gouda L, et al. Extremely Slow Photoconductivity Response of CH3NH3PbI3 Perovskites Suggesting Structural Changes under Working Conditions. The Journal of Physical Chemistry Letters, 2014, 5(15): 2662-2669.
- [16] Hernández Balaguera E, Bisquert J. Time Transients with Inductive Loop Traces in Metal Halide Perovskites. Advanced Functional Materials, 2024, 34(6).
- [17] Lopez-Varo P, Jiménez-Tejada JA, García-Rosell M, et al. Device Physics of Hybrid Perovskite Solar cells: Theory and Experiment. Advanced Energy Materials, 2018, 8(14): 1702772.
- [18] Chen B, Yang M, Zheng X, et al. Impact of Capacitive Effect and Ion Migration on the Hysteretic Behavior of Perovskite Solar Cells. The Journal of Physical Chemistry Letters, 2015, 6(23): 4693-4700.
- [19] Bisquert J. Hysteresis, Impedance, and Transients Effects in Halide Perovskite Solar Cells and Memory Devices Analysis by Neuron-Style Models. Advanced Energy Materials, 2024, 14(26).
- [20] O'Kane SEJ, Richardson G, Pockett A, et al. Measurement and modelling of dark current decay transients in perovskite solar cells. Journal of Materials Chemistry C, 2017, 5(2): 452-462.
- [21] Ravishankar S, Almora O, Echeverría Arrondo C, et al. Surface Polarization Model for the Dynamic Hysteresis of Perovskite Solar Cells. The Journal of Physical Chemistry Letters, 2017, 8(5): 915-921.

## DESIGN OF WATER PIPELINE MONITORING SYSTEM BASED ON MULTI-SOURCE INFORMATION FUSION

LiTong Ma, Bo Ma\*

*Yinchuan University of Energy, Yinchuan 750001, Ningxia, China. Corresponding Author: Bo Ma, Email: 1192897898@qq.com* 

Abstract: Water pipelines are generally buried in the ground, as a typical underground hidden engineering, their structural damages such as pipe burst, leakage, seepage and uneven settlement are characterized by strong concealment and long disaster-causing chain, which not only cause a large amount of waste of water resources, but also lead to safety accidents such as pavement collapse, which seriously threaten public safety. This study aims to propose a multi-dimensional monitoring system that integrates distributed fiber optic sensing and IoT technologies. Through in-depth analysis of the formation principle of pipe burst, leakage, seepage, uneven settlement and other problems, we utilize the deployment of  $\Phi$ -OTDR fiber optic arrays (spatial resolution of 0.5m) to integrate high-precision pressure transmitters (accuracy  $\pm 0.1\%$ FS) and electromagnetic flow meters (accuracy  $\pm 0.5\%$ ) to construct a multi-physical field synchronous sensing network, and to achieve the monitoring of pipeline pressure transmitt (sampling rate  $\geq 100$ Hz), flow rate abnormality (detection sensitivity  $\leq 0.1$ L/s), temperature gradient (resolution 0.1°C), negative pressure wave, stress and strain distribution ( $\mu\epsilon$  level) holographic monitoring, and early warning and precise positioning. Engineering validation shows that this system helps to detect pipeline problems in time, reduce accident losses, guarantee the reliable operation of water pipelines, provide strong support for the stability and safety of the water transmission system, and provide key technical support for the construction of a resilient urban water transmission system.

Keywords: Pipeline health monitoring; Water pipeline; Sensor; Multi-physical field coupling

#### **1 INTRODUCTION**

As the core component of modern municipal infrastructure, urban water pipeline is the key link to ensure the stable supply of water for residential life and industrial production, and its safe operation is directly related to the protection of people's livelihood and economic development. However, affected by multiple factors such as geological environment variability, material aging, and third-party construction disturbance, pipeline systems frequently suffer from structural failure accidents such as pipe bursts, leaks, and uneven settlement, which bring huge losses to society and the economy, for example, in 2017, the continuous bursting of the DN1600 water supply main pipe in the west line of the city of Linyi led to two large-scale water shutdowns in the main urban area, and the direct economic loss amounted to 3,552,800 yuan.7 Such accidents not only cause waste of water and direct economic losses, but also lead to a loss of water resources and a decrease in water consumption and water supply costs. Such accidents not only cause waste and direct economic losses, but also may trigger chain reactions in the industrial chain - the indirect economic losses incurred by industrial enterprises due to the interruption of water supply leading to production stagnation can be up to 3-5 times of the direct losses.

The current monitoring technology is facing a double challenge: on the one hand, the traditional point sensors (such as pressure/flowmeter) have low spatial resolution, weak anti-interference ability and other shortcomings, it is difficult to realize the long-distance buried pipeline monitoring of the whole area coverage; on the other hand, a single-parameter monitoring system is unable to effectively characterize the multi-physical coupling of the pipeline damage mechanism. In recent years, academics have made breakthroughs through technology integration: distributed fiber optic sensing technology (BOTDR/OFDR) can realize strain-temperature-vibration multi-parameter simultaneous sensing[1-2], with a spatial resolution of meters, and detection sensitivity exceeding that of traditional sensors by two orders of magnitude; flexible piezoelectric vibration sensing network can accurately identify leakage aperture and positioning error <0.5m; and smart ball (SmartBall) can be used to detect the damage of the pipeline. The detection rate of small leakage (<1L/min) is increased to 92% by combining machine learning algorithms with mobile detection devices such as SmartBall.

Based on the above background and the shortcomings of previous studies, this study aims to design a comprehensive, efficient and accurate design for pipe burst/leakage monitoring, third-party disturbance/leakage monitoring, and settlement/stress monitoring of water pipelines. The design will comprehensively consider a variety of factors affecting pipeline safety, optimize the sensor selection and installation layout, and improve the monitoring accuracy and early warning capability for various failure states of pipelines. Through this study, it is expected to provide a more reliable guarantee for the safe operation of pipelines, reduce the economic losses and social impacts caused by pipeline accidents, and promote the further development of pipeline monitoring technology.

#### 2 OMNI-DIRECTIONAL MONITORING SYSTEM ARCHITECTURE

The purpose of this paper is to use distributed fiber optic sensing technology, combined with conventional manometers and flow meters to monitor the pipeline in real time, mainly to achieve three aspects of the function: burst/leakage monitoring, third-party disturbance/leakage monitoring, and settlement/stress monitoring. The omni-directional monitoring system of water pipeline status adopts a layered data communication network architecture, which is divided into field equipment layer, control layer and information management layer.

Measurement data is collected by various sensors in the field equipment layer, then transmitted to the data acquisition instruments in the control layer via fiber optic data network to obtain multiple types of monitoring quantities, and finally transmitted remotely from the control layer to the information management layer via data transmission equipment. The data transmission equipment is compatible with many types of fiber optic and electrical acquisition instruments, and the remote transmission of massive monitoring data is realized by means of Internet/wireless Internet. Pressure sensors and flow sensors adopt 485 bus instruments. The structural schematic diagram of the omni-directional monitoring system of water pipeline status is shown in Figure 1. The selection and arrangement of sensors for all-round monitoring of water pipeline status are shown in Table 1.



Figure 1 Structural Schematic Diagram of an All-Round Monitoring System for the Condition of Water Pipelines

Monitoring content	transducers	Sensor arrangement	Data acquisition instruments (control level equipment)
Pipe burst/	Fiber Optic Grating FBG Sensors	Laying inside the bottom of the pipe along the axis of the pipe	Distributed Fiber Optic Collector
leakage monitoring	Manometers and flow meters	Water main and branch pipelines at various intersections	PLC controller
Third-party disturbance/leakage monitoring	Distributed Fiber Optic Temperature Sensors	Parallel to the bottom of the pipe in the direction of the pipe axis	Distributed Fiber Optic Temperature Collector
Settlement/stress monitoring	Distributed Fiber Optic Strain Sensors	Three distributed fiber optic strain sensors placed in parallel along the axis of the pipeline	Distributed Fiber Optic Strain Gauge

#### 2.1 Information Management System Architecture and Functions

In the all-round monitoring system of water pipeline status, the information management layer contains core switches, servers, workstations and other equipment. Deploying industrial-grade Layer 3 ring switch (H3C S6850-56HF, backplane bandwidth 5.76Tbps) to build the backbone network, and realizing the fusion of heterogeneous data from multiple sources through Kafka stream processing platform. The information management layer plays a crucial role, and it has multiple functions:

(1) Data collection and integration function: the information management layer collects data from each control layer data collection instrument of the water pipeline, and the data format produced by different data collection instruments may vary, and the information management layer standardizes and stores the collected data. This enables subsequent data processing and analysis to be carried out on a standardized basis and improves data availability.

(2) Data storage and management function: Considering the continuity and mass of water pipeline monitoring data, the information management system adopts appropriate data storage technology.

(3) Information sharing and visualization function: The information management system provides the processed

monitoring information to different departments, such as pipeline maintenance department, water supply dispatching department and emergency management department. At the same time, it makes complex data easier to understand and use through charts, maps and three-dimensional models, etc., and visualizes the pipeline's operation status, fault location and historical data trends. It allows managers to quickly grasp the key information of pipelines and make accurate decisions.

#### 2.2 Control Layer Architecture and Functions

In the omni-directional monitoring system of water pipeline status, the control layer contains convergence switches, programmable logic controllers (PLCs), fiber optic data collectors, communication equipment and other equipment. These devices are uniformly installed in the field control cabinet and distributed in various key positions of the water pipeline. The control layer plays the key role of the top and bottom, which is the key link to realize the safe and stable operation of the pipeline. The core functions of the control layer include:

(1) Multi-source data fusion: through the IEEE 1588 accurate clock synchronization (error  $<1\mu$ s), integrating PLC process parameters (pressure, flow), fiber optic strain data (100Hz sampling), vibration spectrum (0-20kHz) and other multi-dimensional information, to build a spatio-temporally aligned data cube.

(2) Intelligent decision-making control: Model predictive control (MPC) algorithm is adopted to optimize the regulation strategy in real time based on the pipeline hydraulics model.

(3) Pressure closed-loop control: drive the motorized control valve (Fisher DVC6200) through the PID algorithm to maintain the pressure fluctuation  $\leq \pm 0.05$ MPa

Equipment deployment follows the IEC 61499 standard, the control cabinet to meet the IP54 protection level, environmental adaptability indicators: operating temperature -40 °C ~ +70 °C, humidity 0-95% RH. Through the TSN time-sensitive network (IEEE 802.1Qbv) to ensure that the control command transmission delay <2ms, jitter <50 $\mu$ s. programmable logic controller (PLC) as shown in Figure 2, the The distributed fiber-optic temperature data collector is shown in Figure 3, and the aqueduct control cabinet architecture is shown in Figure 4.



Figure 2 Programmable Logic Controller (PLC)







Figure 4 Aqueduct Control Cabinet Architecture Diagram

#### 2.3 Field Device Layer Architecture and Functionality

In the omni-directional monitoring system for the condition of the water pipeline, the field equipment layer equipment contains various types of sensors such as flow sensors, pressure sensors, high-sensitivity fiber grating FBG sensors, distributed fiber optic temperature sensors, distributed fiber optic strain sensors, and so on. These sensors are distributed

in key locations of the water pipeline to form a sensor network. The field device layer is the foundation of the entire system and plays an indispensable and critical role. The pressure transmitter is shown in Figure 5, the electromagnetic flow meter is shown in Figure 6, the distributed fiber optic strain sensor is shown in Figure 7, the distributed fiber optic temperature sensor is shown in Figure 8, and the fiber grating FBG sensor is shown in Figure 9.



Figure 5 Pressure Transmitter



Figure 6 Electromagnetic Flow Meter



Figure 7 Distributed Fiber Optic Strain Sensors

Figure 8 Distributed Fiber Optic Temperature Sensors



Figure 9 Fiber Optic Grating FBG Sensors

#### **3** EQUIPMENT SELECTION AND INSTALLATION

#### **3.1 PLC Selection and Hardware Configuration**

Based on the control requirements of the water transmission system and the scale of the project, this study adopts Siemens S7-1200 series PLC as the core controller. This series PLC has complete input and output interfaces, efficient data processing capability and reliable communication function, which can meet the technical requirements of the system in data acquisition, control operation and remote communication.

The core configuration of PLC includes CPU module, power supply module, digital input/output module (DI/DO) and analog input/output module (AI/AO). Among them, the CPU module selects S7-1214C, which integrates 6 digital input points and 2 digital output points, which is sufficient to meet the basic switching signal acquisition and control requirements. For the acquisition of flow, pressure and other analog signals, the configuration of a dedicated analog input module SM1231, the module can receive 0-10V or 4-20mA standard industrial signals, and convert them into digital for PLC data processing.PLC programmable controller I / O point allocation is shown in Table 2.

Table 2 PLC Programmable Controller I/O Points Table						
Fauinment/Instrument Name	Name of measurement point	Signal form				
Equipment/instrument ivanie	Name of measurement point	DI	DO	AI	AO	
Electromagnetic flow meter				1		
Pressure Transmitter				1		
	Auto/Manual position	1				
	Valve open states	1				
Motorized values	Valve closed status	1				
Motorized valves	fault state	1				
	Open Valve Command		1			
	Shutdown command		1			
add up the total		4	2	2		

#### 3.2 Pipe Burst/Leakage Monitoring Design and Equipment Installation

When a pipe burst or leak occurs, the stress waves (including negative pressure waves and acoustic waves) propagating in the fluid medium inside the pipe can be effectively detected by high-sensitivity fiber Bragg grating (FBG) sensors. Through the cooperative pressure and flow multi-parameter monitoring means, the system can realize real-time monitoring and diagnosis of leakage events, and achieve sub-kilometer high-precision positioning.FBG as a wavelength-selective reflective grating, its detection system through the laying of special fiber optic cables inside the pipeline[3-4], the use of grating sensors to collect the pipeline axial stress distribution signals, and based on the stress anomalies to achieve leakage detection.

The research team carried out distributed FBG leakage detection simulation experiments, the results show that the grating reflection wavelength offset can be effectively used as a leakage criterion, combined with optical time-domain reflectance (OTDR) addressing technology can be realized leakage spatial localization, the localization error is controlled within the range of  $\pm 0.5\%$  of the length of the pipe section. The sensor is fixed on the inner wall of the pipeline by embedded installation, and the signal is led out to the outside of the pipeline by armored guide cable, and finally the monitoring data is transmitted to the data acquisition instrument in the field control cabinet through industrial bus. The system adopts IP65 protection level chassis, which meets GB3836 explosion-proof requirements and ensures reliable operation in hazardous environments such as oil and gas pipelines. The fiber grating sensor arrangement

structure is shown in Figure 10, and its sectional installation schematic is shown in Figure 11, demonstrating the integration scheme of the grating array with the pipeline structure.



Figure 10 Fiber Optic Grating FBG Sensor Layout

#### 3.3 Third-Party Disturbance/Leakage Monitoring Design and Equipment Installation

Third-party construction activities and working conditions such as pipeline leakage will lead to abnormal changes in the temperature field distribution along the pipeline. Based on the thermodynamic temperature tracing principle, when a leak occurs, the leaking medium will form a localized temperature gradient along the direction of gravity due to the liquid gravity effect[5-6]. By monitoring the heat transfer effect between the leaking liquid and the surrounding soil, the third-party disturbance/leakage monitoring problem can be transformed into a real-time monitoring problem of the temperature field along the pipeline.

Distributed fiber-optic temperature sensing network (DTS) is laid along the bottom axis of the pipeline in the direction of the design spacing to achieve quantitative assessment of the degree of leakage through continuous monitoring of the dynamic characteristics of the temperature field around the pipeline. In order to meet the site construction requirements and long-term service reliability, the sensor adopts a double-layer stainless steel armored structure (in line with GB/T7424.2-2008 standard), in order to provide mechanical protection, at the same time, through the pre-set stress relaxation margin to ensure that the optical fiber is only on the thermal excitation response to avoid the mechanical strain interference[7-8]. Distributed fiber-optic temperature sensor typical arrangement scheme shown in Figure 12, its spatial resolution of up to 1m, temperature measurement accuracy of  $\pm 0.5$  °C.



Figure 11 Distributed Fiber Optic Temperature Sensor Layout

#### 3.4 Settlement/Stress Monitoring Design and Equipment Installation

Subject to the uneven settlement of soil and environmental loads, the pipeline structure will produce complex stress redistribution. In order to monitor the characteristics of the soil displacement field distribution around the pipeline, the optimized structural design of strand-encapsulated distributed fiber-optic strain sensor is used in this study. The sensor enhances the mechanical protection through multi-strand galvanized steel strand (in accordance with GB/T5224-2014 standard), which improves the fiber shear strength to  $\geq$ 200MPa and ensures the long-term stability under complex geological conditions.

The monitoring system is symmetrically laid with 3 distributed fiber optic sensing arrays along the pipeline axis at an angle of  $\pm 120^{\circ}$ , constituting an axial continuous monitoring section, which can synchronously obtain multi-dimensional mechanical parameters such as pipeline bending strain (range  $\pm 1500 \ \mu\epsilon$ ), axial compression strain (precision  $\pm 0.1\%$  F.S.), and neutral plane position strain (spatial resolution of 1m), etc. The sensing cables are connected to the outer wall of the pipeline, and are connected to the outer wall of the pipeline by epoxy resin adhesive (elastic modulus of 2.5GPa), and the external HDPE armored protective layer (thickness  $\geq 2mm$ ) is overlaid to form a monitoring system with strain transfer efficiency of 96%. The system realizes the dynamic monitoring of 10Hz sampling frequency through BOTDA technology, and the data is transmitted to the cloud platform for real-time analysis and early warning through 4G wireless transmission. The typical layout of distributed fiber optic strain sensors is shown in Figure 13, with a strain sensitivity of 1 $\mu\epsilon$  and a temperature compensation accuracy of  $\pm 0.5^{\circ}C$ .



Figure 11 Distributed Fiber Optic Strain Sensor Layout

#### 4 CONCLUSION

In this study, a multi-parameter fusion water pipeline safety monitoring system was developed to realize all-round monitoring of pipe burst/leakage, third-party disturbance/leakage and structural stress/settlement. In terms of sensing network design, the system integrates a multi-physical field sensing array consisting of flow meters (accuracy  $\pm 0.5\%$ ), pressure sensors (range 0-1.6MPa), distributed fiber optic sensors, etc., and constructs a fault diagnosis model based on multi-parameter coupling analysis. The data acquisition network adopts industrial-grade Modbus RTU protocol, and ensures the reliability of data transmission (packet loss rate <0.1%) through the hybrid networking method of 4G wireless communication (transmission interval  $\leq$ 5s) and fiber optic communication (bandwidth  $\geq$ 100Mbps).

Compared with the existing monitoring system, the innovation of this system is reflected in the completeness of the monitoring dimension, and the existing system is mostly limited to the monitoring of a single failure mode with obvious differences and advantages. The system through the integration of pressure (sampling rate of 10Hz), flow (accuracy of 0.5 level), temperature (resolution of 0.1 °C), acoustic (frequency response of 20-20kHz) and strain (sensitivity of 1  $\mu\epsilon$ ) and other multi-dimensional information, so that the burst pipe positioning accuracy is increased to  $\pm$  50m (conventional methods  $\pm$  200m), the false alarm rate is reduced to <0.5 times / month. Compared with the conventional inspection method (cycle  $\geq$  7 days), this system realizes real-time monitoring response at the minute level, which shortens the fault discovery time by more than 85%.

Statistics show that water loss due to pipe burst and leakage accounts for about 3-5% of the total urban water supply. The application of this system can reduce the leakage rate by more than 40%, and reduce the economic loss of about 1.2 million yuan/km per year (calculated according to the industrial water price). Through preventive maintenance, the system improves the reliability of water supply to 99.9%, which effectively guarantees the safety of urban water supply and the continuity of industrial production, and has significant social and economic benefits.

#### **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

#### REFERENCES

- [1] Tian Ye, Zhao Min, Li Kun, et al. Research progress of fiber grating technology in natural gas pipeline safety monitoring. Pipeline Technology and Equipment, 2024(01): 27-33+42.
- [2] Manuel B, Fabrizio D B, Ilaria B, et al. Experimental Investigations of Distributed Fiber Optic Sensors for Water Pipeline Monitoring. Sensors (Basel, Switzerland), 2023, 23(13).
- [3] Zhang Xiaoning. Research on the application of fiber optic sensing technology in municipal water supply and drainage pipe leakage monitoring. Residential and Real Estate, 2024(09): 149-151.
- [4] Han Yang. Research on leakage identification method of existing water supply pipeline network system. Dalian University of Technology, 2018.
- [5] Ziming F, Zhihong L, Liyun P, et al. Acoustic Identification of Water Supply Pipe Leakage Based on Bispectrum Analysis. Journal of Pipeline Systems Engineering and Practice, 2025, 16(3).
- [6] Wang ZF. Research on temperature-strain sensor based on FBG integrated F-P interference. Optical Instrument, 2024: 1-7.
- [7] Guo Jingjing. Research on low-coherence optical interference technique for differential settlement monitoring of high-speed railroad piers. Dalian University of Technology, 2022. DOI:10.26991/d.cnki.gdllu.2022.003607.
- [8] Ying Z, Huan Feng D, Alireza K, et al. On the leak-induced transient wave reflection and dominance analysis in water pipelines. Mechanical Systems and Signal Processing, 2022, 167(PA).
- [9] Lee C, Park W, Park S. Pipeline Structural Damage Detection Using Self-Sensing Technology and PNN-Based Pattern Recognition. Journal of The Korean Society for Nondestructive Testing, 2011, 31(4): 351-359.
- [10] Chiu S, Soga K, Takhirov M S, et al. Assessment of Large-Diameter Ductile Iron Pipeline Joint to Transverse Loading using Distributed Fiber Optic Sensing:Field and laboratory testing 2. Japanese Geotechnical Society Special Publication, 2024, 10(43): 1612-1617.

## **OPTIMIZATION OF COAL MINE ROCKBURST EARLY WARNING SYSTEM**

JiaQi Wu<sup>1\*</sup>, YunMin Tian<sup>2</sup>, TianLe Xiong<sup>1</sup>, JunYao Hou<sup>3</sup>, YunFeng Luo<sup>3</sup>, Hao Chen<sup>3</sup> <sup>1</sup>Reading Academy, Nanjing University of Information Science and Technology, Nanjing 210044, Jiangsu, China. <sup>2</sup>School of Atmosphere Sciences, Nanjing University of Information Science and Technology, Nanjing 210044, Jiangsu, China.

<sup>3</sup>Waterford Institute, Nanjing University of Information Science and Technology, Nanjing 210044, Jaingsu, China. Corresponding Author: Jiaqi Wu, Email: 202283100009@nuist.edu.cn

**Abstract:** As the main energy and important industrial raw materials, coal plays a vital role. With the deep development of coal mining, the risk of underground coal and rock dynamic disasters is rising, which seriously threatens the safety of coal mining. In this paper, the interference signals and precursory characteristic signals in acoustic emission (AE) and electromagnetic radiation (EMR) signals are analyzed. A multi classification model based on the fine KNN model is established to classify the jamming signal data in three different intervals. ARIMA model is used to summarize and analyze the trend characteristics of precursory characteristic signals. The method of random forest classification model is used to classify and identify the time interval of the precursor signal. And calculate the probability of precursory characteristic data at a specific time.

Keywords: ARIMA model; Refined k-nearest neighbor algorithm; Random forest classification model; Non-linear classification

#### **1 INTRODUCTION**

In the process of coal mine production, monitoring and early warning of rock burst and effective prevention and control are still scientific and technological problems to be solved. By monitoring the change trend of acoustic emission (AE) and electromagnetic radiation (EMR) signals, we can determine whether there is a risk of rock burst in the working face or roadway. By dividing the electromagnetic radiation and acoustic emission data into different categories, such as normal working data, precursory characteristic data, interference signal data, sensor disconnection data and working face rest data, the potential rockburst risk can be better identified. Therefore, the analysis and early warning of these monitoring data is of great significance to reduce the occurrence of coal mine accidents[1].

Dou et al.[2] advanced the theoretical understanding of rockburst mechanisms by analyzing the interaction between dynamic (seismic) and static (tectonic) loads. Their research revealed that high static loads in deep mining exacerbate rockburst risk, with microseismic increments from mining-induced tremors acting as critical precursors. Wang[3] introduced a locally weighted C4.5 decision tree algorithm for rockburst risk prediction, achieving 100% accuracy on testing datasets from the Yanshitai coal mine in China. By discretizing continuous attributes via the minimum description length principle and applying 10-fold cross-validation, this method outperformed traditional C4.5 models, which yielded only 71.43% accuracy. Qi Hegang[4] integrated numerical modeling (FLAC3D) with reinforcement learning to simulate stress redistribution during mining, enabling dynamic adjustment of warning thresholds. This approach, validated in the Datong coalfield, reduced false alarms by 30% compared to static threshold systems. Concurrently, human-machine interfaces (HMIs) are evolving to incorporate augmented reality (AR) overlays, providing miners with real-time hazard maps and evacuation routes—a feature tested in the Austar mine post-2014 rockburst reforms[5].

After data preprocessing, outliers are eliminated and the missing values are filled by k-nearest neighbor algorithm. Then the electromagnetic radiation (EMR) and acoustic emission (AE) signals with interference signals are analyzed in three different dimensions through the data in the data table: the external characteristics, internal characteristics and time characteristics of the interference signal distribution. Firstly, according to the external characteristics, the average value, variance, median and extreme value of the transmitted signal are obtained to distinguish the numerical characteristics of the interference signal and other signals. Secondly, according to the time characteristics, the main time period of interference signal is obtained through investigation and analysis. Finally, by drawing the time series distribution map and establishing the nonlinear image analysis model, the size changes of electromagnetic radiation and acoustic emission signals and the corresponding signal type changes in different time periods are analyzed, and the method of identifying the internal characteristics of interference signals and the size changes of electromagnetic radiation and acoustic emission signals is further optimized. This paper analyzes the time series of the corresponding precursor signal sequence, and establishes an appropriate ARIMA model through white noise test and ADF test, so as to summarize and analyze the remaining trend characteristics of the precursor signal over time.

#### **2 PRELIMINARY**

#### 2.1 KNN Algorithm

KNN algorithm is a commonly used machine learning algorithm. Its core idea is based on the nearest neighbor principle. It can classify or regression predict by finding K training data points nearest to the test sample [6]. At the same time, when dealing with the problem of interference signal recognition, it is very important to choose the appropriate distance measurement method for the accuracy of the algorithm. Euclidean distance can effectively evaluate the similarity between samples, so as to achieve accurate signal classification and recognition. Therefore, when dealing with such problems based on high-precision KNN algorithm, Euclidean distance is selected as the distance measurement method. The K-Nearest Neighbors (KNN) algorithm is a supervised learning method used for both classification and regression tasks. It operates on the principle of feature similarity, where the prediction for a new data point is based on the labels or values of its K closest neighbors in the training dataset. The algorithm calculates distances (commonly Euclidean, Manhattan, or Hamming) between the new data point and all training examples, selects the K nearest ones, and determines the prediction through majority voting (for classification) or averaging (for regression). KNN is appreciated for its simplicity and intuitive approach, making it suitable for various applications such as credit rating evaluations, political election forecasting, and pattern recognition. However, its performance can be sensitive to the choice of K and the scale of the data.

#### 2.2 ARIMA Model

ARIMA (P, D, q) model, fully known as autoregressive integrated moving average model, is a statistical model used to analyze and predict time series data. ARIMA model changes the time series data into a stationary series, and then uses the autoregressive (AR) and moving average (MA) parts of the series to fit and predict the model. It is suitable for non seasonal time series data with trend or seasonality [7]. ARIMA model is composed of three main parts: autoregressive term (AR), difference item (I), and moving average term (MA), which are represented by three parameters: P, D, and Q. The general form of the model is ARIMA (P, D, q). The Autoregressive Integrated Moving Average (ARIMA) model is a statistical tool for time series analysis and forecasting. It integrates three components: autoregression (AR), which uses past observations to predict future values; differencing (I), which transforms non-stationary time series into stationary ones by subtracting previous values; and moving average (MA), which incorporates past forecast errors into the prediction. The model is denoted as ARIMA(p, d, q), where p is the order of the autoregressive component, d is the degree of differencing, and q is the order of the moving average component. ARIMA is particularly effective for time series data with trends or seasonality and is extensively used in economics, finance, and inventory management to forecast future values based on historical patterns.

#### 2.3 Random Forest

Random forest is an algorithm based on the idea of ensemble learning. It builds bagging ensemble based on decision tree, further introduces random attribute selection in its training process, and finally makes the decision trees of random forest independent of each other. By inputting new samples, each decision tree of the forest can be judged and classified separately to obtain their own classification results, and finally vote to determine the final random forest classification results [8]. In this process, the feature importance value can be retained. The Random Forest algorithm is an ensemble learning method that combines multiple decision trees to enhance predictive accuracy and model robustness. By utilizing bootstrap aggregation (bagging), it generates diverse training datasets through random sampling with replacement, and each decision tree is trained on a subset of these data. During the tree-building process, a random selection of features is employed at each node split, further increasing the diversity among trees. For classification tasks, the final prediction is determined by majority voting across all trees, while regression tasks use the average of all tree predictions. This approach effectively reduces overfitting and improves generalization performance. Random Forest is widely applied in various domains, including credit risk assessment, medical diagnosis, and recommendation systems, due to its ability to handle high-dimensional data and its resistance to overfitting.

#### **2.4 Notations**

Table 1 Notations			
Symbols	Notations		
Р	Forecast data		
X	Original dataset data as opposed to forecast data		
Gini	Purity measurement		
pri	Probability of correct classification of the ith node		
imp	Feature importance function		
Ι	Decision tree set established during algorithm execution		

The symbols used in the paper are listed in Table 1.

#### **3 ELECTROMAGNETIC RADIATION AND ACOUSTIC EMISSION SIGNALS**

If the amount of data is large and the proportion of outliers is small, we can consider deleting outliers to improve the stability and accuracy of the model. And the characteristics of outliers are very obvious and easy to identify: in d/e type data, if the data fluctuates greatly between several sample points (this paper sets a reasonable adjacent fluctuation range of 45%-150%), it is set as an outlier and the outliers are processed[9]. In order to get the complete data set, this paper uses the k-nearest neighbors algorithm [10] to fill the gap value. According to the similarity between samples, the knearest neighbor method uses the eigenvalues of the nearest K samples to fill in the missing values. For the values containing the vacancy due to deletion, the data in Annex I is huge and extremely dense, which perfectly meets the numerical characteristic requirements of the k-nearest neighbor algorithm for adjacent data. In this paper, the external characteristics of the electromagnetic radiation and acoustic intensity in the above problems are tested respectively. It is found that the standard deviation and mean value of class C data are significantly different from other data in the electromagnetic radiation and acoustic intensity, so it can be used as the feature selection standard of class C data. By analyzing the approximate time of C-type data (interference signal), we can get its time characteristics: the interference signal distribution caused by electromagnetic radiation and sound wave intensity is concentrated in January to July every year, with certain periodic characteristics. By observing and analyzing the time series distribution map, this paper found that the electromagnetic radiation and acoustic intensity of class C data reached the peak at almost the same time, indicating that there is a strong internal relationship between electromagnetic radiation and acoustic intensity in class C data. Electromagnetic radiation and sound wave intensity will fluctuate greatly in the presence of interference signals (Class C signals). Therefore, this paper describes the internal characteristics of interference signals: EMR and AE produce violent oscillation.

#### **3.1 Classification Forecast**

This paper finds that the overall data (a, B, C, D/E) has the following characteristics:

1.Data features with relatively obvious separation boundaries.Data features with certain local properties, that is, similar samples will gather together in the feature space.

2.Almost perfectly meet the requirements of using KNN algorithm model in this paper.

According to the proportion of interference signal finally counted, Table 2 and 3 obtained the interval range.

S/N time	interval start	interval end
1	2022-5-1 0:01:12	2022-5-1 13:53:24
2	2022-5-1 23:58:53	2022-5-2 16:17:30
3	2022-5-2 18:31:00	2022-5-3 6:29:41
4	2022-5-3 20:25:32	2022-5-4 7:05:44
5	2022-5-4 21:27:23	2022-5-5 6:25:07

Table 3 Time Interval of Acoustic Emission Interference Signal

1451001		ileienee Signai	_
S/N time	interval start	interval end	
1	2022-4-1 0:00:11	2022-4-1 10:20:18	
2	2022-4-1 11:38:56	2022-4-2 8:24:23	
3	2022-4-9 3:47:37	2022-4-9 21:06:36	
4	2022-4-10 1:55:35	2022-4-10 9:05:24	
5	2022-4-11 1:56:47	2022-4-11 9:12:02	

#### **3.2 ARIMA Forecast**

By observing and analyzing the time series diagram, it can be found that when the precursor characteristic signal appears, the electromagnetic radiation signal intensity will gradually increase or intermittently increase, and when the rock burst occurs, the electromagnetic radiation signal intensity reaches the highest value, and then decreases sharply in a short time. The intensity of acoustic emission signal will gradually increase, and when rock burst occurs, the intensity of acoustic emission signal will gradually increase, and when rock burst occurs, the intensity of acoustic emission signal will gradually increase, and when rock burst occurs, the intensity of acoustic emission signal will sharply reduce [3]. The overall trend is characterized by cyclic increase. The time series diagram of the occurrence of the overall precursory characteristics shows a periodic repeating pattern, which indicates that there is a

periodic trend when the precursory characteristic signals appear; In the same period of time, the precursory characteristic signal has an obvious growth or decline trend, reflecting a certain trend of violent fluctuations. According to these characteristics, this paper intends to establish ARIMA model to analyze the trend characteristics of precursory characteristic signals.

According to the precursory characteristic signal data of EMR, this paper uses t-value test and finds that the p value of the test statistic is 0.0000, which is less than the significance level of 0.01. Therefore, the original hypothesis is rejected and the alternative hypothesis is accepted, indicating that it is a non-stationary sequence with fluctuation.

According to the precursory characteristic signal data of AE, this paper also uses t-value test, and finds that the p value of the test statistic is also 0.0000, which is less than the significance level of 0.01. Therefore, the original hypothesis is rejected and the alternative hypothesis is accepted, indicating that it is a non-stationary sequence with fluctuation. The maximum lag point of ACF autocorrelation function graph is used to roughly judge the Q value. The p value is determined by the maximum lag point of PACF partial autocorrelation function graph. However, the correctness of the parameters obtained in this way is low. In order to ensure the correctness of the parameters, this paper next needs to carry out model estimation to obtain the values of P and Q.

By comparing the BIC values under different differential orders, the parameter value that can minimize the BIC is selected. In this comparison, it is found that when the autoregressive term P=0, the order of the moving average term q=4, that is, the BIC value reaches the minimum. Therefore, this paper chooses to establish ARIMA (0,1,4) model.

After the model is established, the residual is tested by white noise. If the residual is white noise, it indicates that the selected model can fully identify the law of time series data, that is, the model is acceptable; If the residual is not white noise, it means that the sequence may have a certain pattern, structure or correlation, and does not have pure randomness, which may be useful for data analysis and prediction. This means that other types of information and associations in the data can be explored. From the results, the p value of the Ljung box test of the two groups of data is less than 0.01, which means that there is a significant autocorrelation in the residuals, rejecting the original assumption that the residuals are white noise. This shows that the model can still be further optimized.

#### **4 RANDOM FORESTS**

This paper found that the prediction probability of the precursor characteristic signal reached more than 80%, which met the extraction standard of the time series of subsequent precursor characteristic signals, and then obtained the time interval of electromagnetic radiation precursor characteristics and acoustic emission precursor characteristics, as shown in Table 4 and 5.

Table 4 Characteristic Time Interval of Electromagnetic Radiation Precursor for RF				
S/N time	interval start	interval end		
1	2020-4-8 2:23:05	2020-4-11 10:06:07		
2	2020-4-22 21:41:27	2020-4-27 12:33:47		
3	2020-5-23 10:03:33	2020-6-5 5:21:55		
4	2021-12-15 3:47:11	2021-12-20 23:59:11		
5	2021-11-24 5:47:11	2021-11-30 17:04:02		
Table 5 Charac	teristic Time Interval of Acoustic Emis	sion Precursors for RF		
S/N time	interval start	interval end		
1	2021-11-1 0:01:01	2021-11-2 17:00:13		
2	2021-11-25 20:59:12	2021-11-30 8:25:06		
3	2021-12-3 10:10:06	2021-12-9 19:14:11		
4	2021-12-12 6:21:47	2021-12-16 17:02:55		
5	2022-1-1 5:59:07	2022-1-14 7:48:56		

Based on the analysis and detection of ARIMA (0,1,4) model and random forest model, this paper found that ARIMA model had excellent fitting effect (the goodness of fit r was as high as 0.96), but because the original hypothesis was rejected in ADF test, that is, the residual did not meet the white noise sequence, this model could not predict and estimate the target sequence well. Therefore, this paper uses the random forest classification model to predict and classify the target sequence and get the corresponding classification data results. Because through the analysis of ARIMA model and random forest model, this paper confirms that the characteristics of precursor characteristic signals have a certain persistence, that is, they are aggregated and distributed in a certain period of time series and the time is about 7 days (about 1000 samples).

Firstly, by using the frequency estimation method to calculate the proportion of the number of precursor characteristic signals in the first 1000 sample points of the target time point to the total 1000 samples, the corresponding probability value is obtained. Next, according to these probability values, the occurrence probability of precursory characteristic data in the target time period is evaluated. This probability value can be regarded as the probability of precursory characteristic data at the last moment of each time period. The comparison is shown in Table 6.

Table of Hobability of Occurrence of Hecuisory Characteristics at the Thile of Data Concection						
Time of electromagnetic radiation data	Time of probabilistic acoustic emission	Data	Probability			
2023-1-24 23:58:36	0.07715	2023-1-24 23:58:36	0.05299			
2023-2-11 23:59:20	0.57242	2023-2-11 23:59:20	0.51245			
2023-2-26 23:59:27	0.51605	2023-2-26 23:59:27	0.48765			
2023-3-10 23:58:14	0.55637	2023-3-10 23:58:14	0.55601			
2023-3-30 23:58:13	0.51187	2023-3-30 23:58:13	0.54237			

 Table 6 Probability of Occurrence of Precursory Characteristics at the Time of Data Collection

In the construction of the model, the model absorbs and processes a large amount of data, which has strong stability. The model is suitable for the prediction of rockburst indexes in the future, and integrates the advantages of various classification models. It has been tested for many times and found that its fitting and prediction effect is good and has strong universality; For example, the ARIMA (0,1,4) model introduced in this paper has a high goodness of fit for precursor signals, and can achieve high accuracy and good prediction effect. The ADF detection of ARIMA model found that it was not white noise, but failed to show the structure and characteristics of its residual sequence.

#### **5** CONCLUSION

This paper presents a comprehensive study on the optimization of coal mine rockburst early warning systems, focusing on the analysis and processing of electromagnetic radiation (EMR) and acoustic emission (AE) signals. The research primarily addresses three key objectives: firstly, to identify and classify interference signals within EMR and AE data; secondly, to develop mathematical models for precisely locating precursor characteristic signals and determining significant trend features; and thirdly, to establish a probabilistic model for predicting the occurrence of precursor signals at specific time intervals. The paper underscores the effectiveness of the proposed methodologies. The integration of advanced data preprocessing and KNN modeling demonstrates proficiency in interference signal identification.

#### **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

#### REFERENCES

- [1] Song C H, Lu C P, Liu J R. Moment Tensor and Stress Field Inversions of Mining-Induced Seismicity in A Thick-Hard Roof Zone. Rock Mechanics and Rock Engineering, 2024, 57(3): 2267-2287.
- [2] Watson J, Canbulat I, Zhang C, et al. Energies Within Rock Mass and the Associated Dynamic Rock Failures. Rock Mechanics and Rock Engineering, 2025, 58(5): 4935-4958.
- [3] Wojtecki U, Krawiec K, Wikaa M, et al. An attempt to determine the cause of the strong tremor responsible for a rockburst in a hard coal mine based on numerical modeling and spectral parameters. Geology, Geophysics & Environment, 2024, 50(4).
- [4] Wang K, Xie T, Mei L I, et al. A surrogate model for the rapid prediction of rockburst risk based on numerical samples and random forest classifier. Journal of Tsinghua University (Science and Technology), 2024, 64(7): 1203 -1214.
- [5] Kuzniar K, Tatara T, Zajac M. Experimental and numerical assessment of soil-structure interaction effects in the case of mine-induced vibrations. IOP Publishing Ltd, 2024.
- [6] Eremin M O, Chirkov A O, Pazhin A, et al. Finite-Difference Analysis of Influence of Borehole Diameter and Spacing on Reduction in Rockburst Potential of Burst-Prone Coal Seams. Mining, 2024, 4(4).
- [7] Jiahao S, Wenjie W, Lianku X. Predicting Short-Term Rockburst Using RF-CRITIC and Improved Cloud Model. Natural resources research, 2024, 33(1): 471-494.
- [8] Dai J, Gong F, Huang D, et al. Quantitative evaluation method of rockburst prevention effect for anchoring rock masses around deep-buried tunnels. Tunnelling and Underground Space Technology incorporating Trenchless Technology Research, 2025, 156.
- [9] Wojtecki U, Iwaszenko S, Apel D B, et al. Use of machine learning algorithms to assess the state of rockburst hazard in underground coal mine openings. 2022.

[10] Wang Q, Ma T, Yang S, et al. Intelligent rockburst level prediction model based on swarm intelligence optimization and multi-strategy learner soft voting hybrid ensemble. Geomechanics & Geophysics for Geo-Energy & Geo-Resources, 2025, 11(1).

# DESIGN OF DUST MONITORING SYSTEM FOR PRODUCTION WORKSHOP

MingMao Gong\*, SiZhe Zheng, SiYan Xu

School of Electronic Information Engineering, Sichuan Technology and Business University, Chengdu 611745, China. Corresponding Author: MingMao Gong, Email: fengxu0217@126.com

**Abstract:** As an important place for clean production, real-time monitoring of dust concentration in production workshops is the key to ensuring product quality. Starting from practical application requirements, this article designs a dust monitoring system for production workshops based on STM32. The system integrates PG-03CR six channel laser dust sensor and SHT85 digital temperature and humidity sensor, and combines TFT screen to achieve information visualization and real-time warning function. In addition, through 4G wireless communication technology combined with MQTT communication protocol, monitoring data is transmitted in real-time to the monitoring center, achieving remote monitoring functionality. Through actual testing and comparison verification, the system has maintained good response accuracy and error control in PM2.5, PM10, and temperature and humidity monitoring. **Keywords:** Dust monitoring; STM32; Laser dust sensor; 4G; MQTT protocol

#### **1 INTRODUCTION**

Dust adhering to the surface of products can affect product quality, such as in industries such as electronics, food, and pharmaceuticals, leading to an increase in defect rates. Employees who are exposed to high concentrations of dust for a long time may also face the risk of occupational diseases such as pneumoconiosis, respiratory diseases, and skin diseases. Dust monitoring can help companies understand the dust situation in the workshop, take effective protective measures, improve product yield, and reduce the risk of employee illness [1-2]. With the development of Internet of Things technology, dust online monitoring systems based on Internet of Things technology can gradually collect real-time concentrations of particulate matter such as PM2.5 and PM10, as well as environmental parameters such as temperature and humidity, and upload data to cloud platforms through 4G communication modules to achieve remote monitoring and intelligent alarm reporting[3].

#### 2 DESIGN SCHEME FOR DUST MONITORING SYSTEM

The overall architecture of the system includes battery power management, sensor data acquisition, battery voltage detection, display module, and user interaction. The system is powered by a 7.4V lithium battery and is supplied with stable power through LDO to ensure that all modules operate at normal voltage. In the design of the dust monitoring system, the million level clean PG03CR dust sensor was selected, which has a particle detection range of 0.3~10um and can detect 6 particle sizes including 0.3, 0.5, 1.0, 3.0, 5.0, and 10um. The STM32F103 is chosen as the embedded microcontroller, which has powerful functions and low power consumption, meeting the overall control and low-power requirements of the system[4] In addition, SH85 is used to detect temperature and humidity, and monitoring data is transmitted through a 4G wireless communication module. A 10.1-inch serial touch screen is used as the human-machine interaction interface. The overall system design diagram is shown in Figure 1.



#### **3 HARDWARE CIRCUIT DESIGN**

#### 3.1 Voltage Stabilization Circuit Design

In this design, a 7.4V lithium battery is used for power supply, but both the TFT LCD screen and dust sensor require a 5V voltage, STM32, The working voltage of temperature sensors is 3.3V, so this system uses two voltage reduction circuits.

The 5V voltage regulator circuit is shown in Figure 2. The input terminal of LM29150RS-5.0 is connected to the main power supply section VIN, and the front end is connected in series with C16 and C17 filtering capacitors to suppress high-frequency interference and ripple voltage. The output terminal is connected to C18 and C19 to construct a complete input-output filtering network, further improving the system's anti-interference ability. The chip has an output capability of up to 1.5A, which is sufficient to meet the power supply requirements of stable and consistent brightness of TFT screens, ensuring the reliability and response speed of image display[5].



Figure 2 Schematic Diagram of 5V Voltage Regulator Circuit

The core control modules such as STM32F103RCT6 microcontroller and SHT85 temperature and humidity sensor require stable 3.3V voltage supply and are sensitive to voltage fluctuations. Therefore, the system design adopts the ME6210A33PG linear voltage regulator chip, which outputs a 3.3V voltage after voltage reduction and stabilization of the battery voltage[6]. The ME6210 chip has low static power consumption and fast response characteristics, suitable for power drive applications of embedded system main control chips. The circuit diagram is shown in Figure 3.



Figure 3 Schematic Diagram of 3.3V Voltage Regulator Circuit

#### 3.2 Temperature and Humidity Module Design

This design uses the SHT85 temperature and humidity sensor module to monitor the temperature and humidity in the working environment in real time, ensuring that the dust monitoring system can compensate according to environmental conditions[7]. SHT85 has the characteristics of high precision and low power consumption, and is widely used in industrial environments. The sensor exchanges data with the STM32F103RCT6 microcontroller through the <sup>12C</sup> communication protocol. The <sup>12C</sup> bus communication simplifies the hardware connection, and its circuit schematic is shown in Figure 4.



Figure 4 SHT85 Temperature and Humidity Sensor Circuit Diagram

#### 3.3 Design of Dust Particle Module

This design uses the PG-03CR six channel dust sensor, which has six channel laser detection capabilities and can simultaneously output data on the quantity and mass concentration of particles with particle sizes of  $0.3 \mu$  m,  $0.5 \mu$  m,  $1.0 \mu$  m,  $2.5 \mu$  m,  $5.0 \mu$  m, and  $10 \mu$  m[8]. It is suitable for scenarios that require high air cleanliness. The PG-03CR sensor is based on the principle of laser scattering for particle recognition. It integrates a fan air pump, laser, photodiode array, and signal processing circuit internally. During the air sampling process, real-time recognition of particles of different sizes is achieved through the detection of scattered light intensity. The sensor has completed the calibration of particulate matter mass concentration at the factory, and defaults to outputting key parameter values such as PM2.5 and PM10. The relevant parameters are shown in Table 1. The sensor experiment uses UART interface for communication, and the circuit is relatively simple, so it will not be repeated here.

Table 1 DC 02CD Demonstrate Table

Parameter Category	Specific indicators
Detecting particle size range	0.3 μ m~10 μ m (six channels)
Particle size channel	0.3µm, 0.5µm, 1.0µm, 2.5µm, 5.0µm, 10µm
Detection accuracy (PM2.5)	$0100\mu g/m^3$ : ±10 $\mu g/m^3$ , 1001000 $\mu g/m^3$ : ±10%
Concentration resolution	1µg/m³
Output data refresh interval	1 second
working voltage	DC 5V ±0.2V
Working current	≤100mA
communication model	UART serial communication
response time	1 second
Working temperature/humidity	-10°C50°C; 095% RH
Cleanliness level standard	Compliant with ISO14644-1, ISO5 to ISO9 levels

#### 3.4 Screen Display Module Design

The 10.1-inch TFT serial port screen is used as a display module to visually display key information such as the system's working status, dust concentration, temperature and humidity data. The display module interacts with the microcontroller through UART communication protocol to ensure fast data transmission and real-time updates. The serial port screen adopts the common TFT LCD display technology, which has high resolution and brightness, and can display clearly in various environments. The serial port screen is connected to the UART interface of the microcontroller through a dedicated adapter board. The adapter board also integrates SD card interface, buzzer interface, and speaker interface for users to update the serial port screen interface, which can easily achieve voice broadcasting and alarm functions. The actual picture of the screen connection is shown in Figure 5.



Figure 5 Screen Connection Physical Image

#### 3.5 Design of Level Conversion Circuit in 2.5 4G Communication Circuit

This design uses the ML305 4G communication module for monitoring data transmission, and the hardware communication interface between this module and STM32 is UART interface. However, since the ML305 core operates at 1.8V, while the STM32 operates at 3.3V, the key point of this design lies in the level conversion circuit of the communication interface. The UART interface uses a full duplex communication interface, so level conversion only requires unidirectional conversion. The circuit schematic of this function is shown in Figure 6.

Taking the example of 4G sending data to STM32 in the upper left corner of the figure, when 4G-DTU\_TX does not send data, transistor Q3 is turned off, and U3RX remains high under the action of the pull-up resistor, while UART remains idle; When 4G-DTU\_TX sends data, the DTU\_TX pin is 0V, and transistor Q3 is turned on. The U3RX pin is pulled low to a low level, and the serial port receives a low level. Through this simple circuit, the conversion process from 1.8V to 3.3V level is achieved.



Figure 6 Design of Level Conversion Circuit in 4G Communication Circuit

#### **4 SYSTEM SOFTWARE DESIGN**

The dust monitoring system in the production workshop mainly monitors dust data through a six channel dust sensor, which uses UART interface to exchange data with STM32; Use the SHT85 temperature and humidity sensor to collect environmental temperature and humidity information, and display the relevant information on a 10 inch serial port LCD screen; At the same time, relevant data will be transmitted to the cloud platform through 4G modules to achieve remote monitoring functionality. The data collected by the sensor will be transmitted to the main program for processing and analysis, followed by data filtering to remove possible noise and interference. Then, the data will be corrected based on the characteristics of the sensor to ensure its accuracy. To ensure data stability, the main program calculates the average value, reduces the impact of sudden fluctuations, and detects the presence of abnormal data points. All collected data will be displayed in real-time on the TFT screen, and users can interact and view different data items through touch screens or buttons. When the data exceeds the preset security threshold, the system will trigger an alarm mechanism, display warning information on the screen, and remind through voice broadcast to ensure timely response. Even if the data exceeds the threshold, the system will continue to process and store the data to maintain continuous monitoring and recording. The software flowchart is shown in Figure 7.



Figure 7 Overall Flowchart of System Software

#### 4.1 System Testing and Result Analysis

During the testing process, the system selects a conventional laboratory environment as the testing site and conducts fixed-point timed sampling at five equidistant time points throughout the day, 9:00, 12:00, 15:00, 18:00, and 21:00, to observe the diurnal trend of temperature and humidity changes. After each collection, the system transmits data to the PC through the serial port for recording, and compares it with the data from the standard temperature and humidity meter to provide a reference value for calculating the measurement error of the system.

time	Measure temperature (°C)	Actual temperature (°C)	Absolute error (°C)	Relative error (%)	Measure humidity (%)	Actual humidity (%)	Absolute error (%)	Relative error (%)
9:00	17.8	17.6	0.2	1.14	59.3	57.5	1.8	3.14
12:00	22.6	22.5	0.1	0.44	51.2	49.5	1.7	3.44
15:00	21.3	21.2	0.1	0.47	52.4	51.5	0.9	1.75
18:00	18.5	18.3	0.2	1.09	54.2	54.5	0.3	0.55
21:00	16.3	16.2	0.1	0.62	57.8	57.5	0.3	0.52

 Table 2 Temperature and Humidity Test Results

As shown in Table 2, the measurement error of the temperature and humidity sensor exhibits certain fluctuations at different time periods. The absolute error between the measured temperature and the actual temperature is within 0.2 °C, and the relative error is between 0.44% and 1.14%. The temperature error varies at different time periods, with a significant error of 0.2 °C at 9:00 and 18:00, and relative errors of 1.14% and 1.09%. The absolute error of humidity measurement error is generally small, ranging from 0.3% to 1.8%, and the relative error range is 0.52% to 3.44%. Overall, the error between the measurement results of the sensor and the actual values does not exceed a reasonable range, which can meet the needs of daily applications.

The dust test results are shown in Table 3. The PM2.5 and PM10 concentrations measured by the system in this design are close to the values of the standard dust tester at all five time points throughout the day. The maximum absolute error of PM2.5 measurement is  $3 \mu \text{ g/m}^3$ , and the relative error is between 2.5% and 7.14%; The maximum absolute error of PM10 is  $3 \mu \text{ g/m}^3$ , and the relative error remains between 1.47% and 5.10%, with the error controlled within an acceptable range. Overall, the PG-03CR sensor has stable output performance and high data reliability under this system structure, making it suitable for environmental scenarios such as dust-free workshops that require high particle concentration monitoring.

time	Measure PM2.5 (µ g/m <sup>3</sup> )	Actual PM2.5 (µ g/m <sup>3</sup> )	Absolute error (μ g/m <sup>3</sup> )	Relative error (%)	Measure PM10 (µ g/m <sup>3</sup> )	Actual PM10 (μ g/m <sup>3</sup> )	Absolute error (μ g/m ³)	Relative error (%)	
9:00	38	37	1	2.70	62	60	2	3.33	
12:00	43	40	3	7.14	69	68	1	1.47	
15:00	40	38	2	5.26	64	63	1	1.59	
18:00	38	36	2	5.56	63	60	3	5.10	
21:00	36	34	2	5.88	60	58	2	3.45	

Table 3 Dust Test Results

#### **5 CONCLUSION**

This design focuses on monitoring dust in production workshops, covering the entire process from hardware collection to data processing, presentation, and communication. Advanced sensors have been selected for data acquisition in hardware to ensure the accuracy and effectiveness of the data. The data processing module verifies, filters, and analyzes the collected environmental data, improving the accuracy and reliability of the system. The system also integrates a 4G wireless communication module, which can transmit real-time monitoring data from the production workshop to the IoT platform, achieving data visualization and remote monitoring functions.

#### **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

#### REFERENCES

- [1] Lina Zheng, Zikang Feng, Jia Liu, et al. Cross-concentration calibration of low-cost sensors for effective dust monitoring at construction sites. Journal of Aerosol Science. 2024, 182, 106456.
- [2] Parkavi A, Sowmya B J, Sini Anna Alex, et al. Air quality and dust level monitoring systems in hospitals using IoT. Discover Internet of Things. 2025, 5(1): 23-23.
- [3] Cody Wolfe, Emanuele Cauda, Milan Yekich, et al. Real-Time Dust Monitoring in Occupational Environments: A Case Study on Using Low-Cost Dust Monitors for Enhanced Data Collection and Analysis. Mining, Metallurgy & Exploration. 2024, 41(4): 1709-1718.
- [4] Santa Nestor, Sarver Emily. Advancing respirable coal mine dust source apportionment: a preliminary laboratory exploration of optical microscopy as a novel monitoring tool. International Journal of Coal Science & Technology, 2024, 11(1).
- [5] Tuchman Donald P, Mischler Steven E, Cauda Emanuele G, et al. Equivalency of PDM3700 and PDM3600 Dust Monitors. Mining, Metallurgy & Exploration, 2024, 41(2): 719-725.
- [6] di Vacri M L, Scorza S, French A, et al. Evaluation of SNOLAB background mitigation procedures through the use of an ICP-MS based dust monitoring methodology. Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment, 2023, 1056. DOI: https://doi.org/10.1016/j.nima.2023.168700.
- [7] Karcz Stanisław, Skrabalak Grzegorz, Brudnik Andrzej, et al. Dust Particle Counter for Powder Bed Fusion Process. Sensors, 2022, 22(19): 7614-7614.
- [8] Brajesh Kumar Kori, Himanshu Agrawal. Dust Monitoring Systems and Health Hazards in Coal Mining A Review. International Journal of Trend in Scientific Research and Development, 2021, 5(3): 172-177.

## ELUCIDATING THE DRIVING FACTORS IN SULFUR TRIOXIDE FORMATION UNDER SIMULATED ACTUAL ULTRA-LOW EMISSION PROCESS

#### ZePeng Li<sup>1</sup>, Yasser M. A. Mohamed<sup>2\*</sup>, YingHui Han<sup>1\*</sup>

<sup>1</sup> College of Resources and Environment, University of Chinese Academy of Sciences, Beijing 101408, China.
 <sup>2</sup> Photochemistry Department, National Research Centre, Dokki, Giza, P. O. 12622, Egypt.
 \*Yasser M. A. Mohamed and YingHui Han are both the Corresponding Authors.
 Corresponding Authors: YingHui Han, Email: hanyinghui@ucas.ac.cn; Yasser M. A. Mohamed, Email: y.m.a.mohamed@outlook.com

**Abstract:** Current flue gas pollution control technologies compliant with ultra-low emission standards exhibit limited effectiveness in removing sulfur trioxide (SO<sub>3</sub>)—a key condensable particulate matter (CPM) precursor—necessitating high-efficiency, low-consumption control strategies. To address the poorly elucidated formation mechanisms of SO<sub>3</sub> across ultra-low emission systems, particularly within the SCR+WFGD process chain, this study employed experimental simulations where SO<sub>3</sub> was prepared via the contact process and quantified through controlled condensation coupled with sulfate titration. Catalytic oxidation experiments on cesium-doped V<sub>2</sub>O<sub>5</sub> in a temperature-controlled fixed-bed reactor under simulated actual flue gas revealed reaction temperature as the governing factor for SO<sub>3</sub> conversion, achieving peak efficiency at 485–505°C. Whereas SO<sub>2</sub> concentration exerted non-dominant effects due to sustained catalytic stability, space velocity proved negligible under high-temperature regimes. These mechanistic insights establish fundamental pathways for developing targeted SO<sub>3</sub> mitigation technologies.

Keywords: Sulfur trioxide (SO3); Sulfur dioxide (SO2); Catalytic oxidation; Driving factors; Ultra-low emission

#### **1 INTRODUCTION**

During combustion in boilers and passage through selective catalytic reduction (SCR) units, sulfur dioxide (SO<sub>2</sub>) in coal-fired flue gas reacts with oxygen (O<sub>2</sub>), significantly increasing the concentration of sulfur trioxide (SO<sub>3</sub>) [1-2]. SO<sub>3</sub> is highly toxic and corrosive, causing severe irritation to skin, mucous membranes, and other tissues, thereby posing serious risks to ecological systems and human health [3]. Additionally, SO<sub>3</sub> can react with excess ammonia (NH<sub>3</sub>) injected into the SCR system, forming ammonium sulfate ((NH<sub>4</sub>)<sub>2</sub>SO<sub>4</sub>) and ammonium bisulfate (NH<sub>4</sub>HSO<sub>4</sub>), which can damage the operational equipment of coal-fired power plants [4]. However, current ultra-low emission control systems recommended for coal-fired flue gas treatment only effectively remove nitrogen oxides (NO<sub>x</sub>), SO<sub>2</sub>, and particulate matter (PM), lacking efficient methods for SO<sub>3</sub> control. Therefore, it is imperative to develop highly efficient and low-cost SO<sub>3</sub> control technologies to meet the operational safety and environmental emission requirements of coal-fired power plants.

The configuration of simulated polluted flue gas and measurement of pollutant concentrations in clean flue gas are central to laboratory-scale SO<sub>3</sub> control research. The critical step in the former process is SO<sub>3</sub> generation, while in the latter, it is SO<sub>3</sub> sampling. Ozone oxidation and heated sulfuric acid methods are commonly used in laboratories for SO<sub>3</sub> preparation. The ozone oxidation method has advantages in stable SO<sub>3</sub> production but imposes stringent requirements for ozone preparation. The heated sulfuric acid method involves decomposing sulfuric acid by heating to generate SO<sub>3</sub>, but this process carries safety risks due to the strong corrosiveness of sulfuric acid [5]. The contact process is commonly utilized for industrial sulfuric acid production, involving vanadium pentoxide (V<sub>2</sub>O<sub>5</sub>) as a catalyst to oxidize SO<sub>2</sub> to SO<sub>3</sub> under oxygen-rich conditions. This method offers advantages such as high conversion rates, high product purity, and robust adaptability [6-8]. For SO<sub>3</sub> sampling, commonly employed techniques include controlled condensation, alkaline absorption, isopropanol absorption, and salt absorption methods [9]. Among these, the controlled condensation method effectively prevents measurement errors caused by premature SO<sub>3</sub> condensation and reduces the interference from sulfate ions generated by dissolved SO<sub>2</sub>. It also exhibits high accuracy under various operational conditions. After collecting SO<sub>3</sub> using controlled condensation, sulfate content in the collected solution can be measured by gravimetric analysis, barium chromate photometry, ion chromatography, turbidity measurement, or titration methods, thereby enabling estimation of SO<sub>3</sub> concentration in the simulated flue gas [10, 11].

In this study, the contact process was employed for  $SO_3$  preparation, and  $SO_3$  was collected and measured using controlled condensation combined with sulfate titration. Experiments on catalytic oxidation of  $SO_2$  to  $SO_3$  were conducted under different operating conditions. The methods for  $SO_3$  preparation and concentration measurement described herein can be utilized in further  $SO_3$  control experiments. The insights obtained regarding  $SO_3$  formation under various conditions will serve as references for the formulation of effective  $SO_3$  control strategies.

#### **2** EXPERIMENTAL METHODS

#### 2.1 Simulation System for SO<sub>3</sub> Generation

In this study,  $SO_3$  was generated via the contact process by conducting catalytic oxidation of  $SO_2$  on an isothermal fixed-bed reactor. The experimental setup was designed to investigate the effects of various operational parameters on the  $SO_3$  conversion rate. A schematic diagram illustrating the principle of the fixed-bed experimental system is shown in **Figure 1**.



Figure 1 Schematic Diagram of the Fixed-Bed Experimental System

The simulated flue gas was composed of N<sub>2</sub>, O<sub>2</sub>, SO<sub>2</sub>, and water vapor. Compressed gas cylinders containing certified standard gases were connected via pressure regulators and piping to mass flow controllers (MFCs), which provided real-time control of individual gas flow rates. The water vapor content in the simulated flue gas was regulated by adjusting the temperature of a thermostatic water bath and the flow rate of carrier N<sub>2</sub>. To ensure the stability of water vapor concentration, a series of gas-washing bottles were placed in tandem within the water bath. To prevent irreversible deactivation of the catalyst by water vapor, the mixing point of water vapor and the rest of the simulated flue gas was positioned downstream of the reactor. The pipeline segment extending from the gas-washing bottle outlet to the SO<sub>3</sub> sampling system inlet was wrapped with an electric heating tape to prevent condensation of SO<sub>3</sub> and water vapor within the line.

The reactor system consisted of a quartz reactor and an external heating unit. The reactor was a cylindrical hollow quartz tube fitted with a quartz sand support plate to hold the catalyst bed. The heating unit, located outside the reactor, was used to maintain the reactor at a target reaction temperature. The catalyst employed in this study was a pelletized cesium-doped  $V_2O_5$  catalyst. A thermocouple was embedded within the catalyst bed to enable accurate monitoring of the catalyst layer temperature, which is referred to as the "reaction temperature" in the following sections [12].

#### 2.2 Measurement of SO<sub>3</sub> Concentration

In this study, SO<sub>3</sub> was collected using the controlled condensation method. The SO<sub>3</sub> sampling system is illustrated in Figure 2.



Figure 2 Schematic Diagram of the SO3 Sampling System

After SO<sub>3</sub> was collected using the controlled condensation method, the condensate within the serpentine condenser was rinsed into a volumetric flask using deionized water, and the solution was subsequently diluted to a fixed volume. The concentration of sulfate ions in the solution was then determined by ion chromatography, enabling efficient and accurate quantification of SO<sub>3</sub> concentration in the experimental gas.

#### 2.3 Calculation of SO<sub>3</sub> Conversion Rate

The SO<sub>3</sub> conversion rate was calculated using the following equation:

$$\alpha = \frac{SO_{3\text{out}}}{SO_{2\text{in}}} \times 100\% \tag{1}$$

where  $SO_{3out}$  represents the calculated outlet concentration of SO<sub>3</sub> in the simulated flue gas, in units of ppm, and  $SO_{2in}$  denotes the inlet concentration of SO<sub>2</sub> as set in the experiment, also in ppm.

Since the catalytic oxidation of SO<sub>2</sub> to SO<sub>3</sub> is a heterogeneous catalytic reaction, its conversion rate is primarily influenced by reaction temperature, reactant concentration, and space velocity. In this study, subsequent experiments were conducted under fixed conditions: the catalyst mass (cesium-doped V<sub>2</sub>O<sub>5</sub>) was maintained at 10 g, the catalyst bed height at 20 mm, and the total gas flow rate at 2 L/min.

#### **3** RESULT AND DISCUSSION

#### 3.1 Effect of Reaction Temperature on SO<sub>3</sub> Conversion Rate

Catalytic oxidation experiments for SO<sub>2</sub>-to-SO<sub>3</sub> conversion were conducted at four different temperatures: 425 °C, 455 °C, 485 °C, and 505 °C, in order to investigate the influence of reaction temperature on SO<sub>3</sub> conversion and to identify the optimal temperature for SO<sub>3</sub> generation. The experimental results are presented in **Figure 3**. As shown in the figure, under various operating conditions with different inlet SO<sub>2</sub> concentrations, the SO<sub>3</sub> conversion rate exhibits a general trend of initially increasing rapidly with temperature, followed by a gradual decline. A peak conversion rate was observed within the 485 °C to 505 °C range, indicating the existence of an optimal reaction temperature for the selected cesium-doped V<sub>2</sub>O<sub>5</sub> catalyst.



Figure 3 SO<sub>3</sub> Conversion Rate at Different Reaction Temperatures

The observed trend can be explained as follows: at relatively low reaction temperatures, the activation energy required for SO<sub>3</sub> formation remains high, and the catalyst has either not yet reached or has just reached its ignition temperature. Under these conditions, the catalyst's ability to reduce the activation barrier is limited, resulting in a low SO<sub>3</sub> conversion rate. As the reaction temperature increases, the catalyst becomes more active and the activation energy is more effectively overcome, leading to a continuous rise in SO<sub>3</sub> conversion. However, since the oxidation of SO<sub>2</sub> to SO<sub>3</sub> is an exothermic and reversible reaction, excessively high temperatures thermodynamically suppress the forward reaction, thereby limiting further increases in conversion. Additionally, elevated temperatures may promote undesirable interactions between V<sub>2</sub>O<sub>5</sub> and the silicon dioxide (SiO<sub>2</sub>) support material, accelerating catalyst deactivation and subsequently reducing SO<sub>3</sub> conversion efficiency [13]. In the temperature range of 425–485 °C, the dominant factors are the decreasing activation energy and enhanced catalytic activity, which lead to an increase in SO<sub>3</sub> conversion of the forward reaction and increased catalyst degradation—become more pronounced. Consequently, the SO<sub>3</sub> conversion rate initially increases and then declines with further temperature elevation in this higher range.

#### 3.2 Effect of SO<sub>2</sub> Concentration on SO<sub>3</sub> Conversion Rate

Catalytic oxidation experiments were carried out under varying  $SO_2$  concentrations of 800 ppm, 945 ppm, 5000 ppm, 12,500 ppm, and 25,000 ppm to investigate the influence of  $SO_2$  concentration on the  $SO_3$  conversion rate. The experimental results are presented in **Figure 4**. As illustrated in the figure, although the trends in  $SO_3$  conversion rate variation with respect to  $SO_2$  concentration differ slightly across different reaction temperatures, the overall magnitude of change remains limited. This indicates that the selected catalyst exhibits strong stability and is capable of sustaining effective  $SO_3$  production across a wide range of  $SO_2$  concentrations.



Figure 4 SO<sub>3</sub> Conversion Rate Under Different SO<sub>2</sub> Concentrations

The analysis suggests that within the range of conditions examined in this study,  $SO_2$  concentration is not the dominant factor influencing the reaction. As a result, the  $SO_3$  conversion rate does not exhibit significant fluctuations with varying  $SO_2$  concentrations at different reaction temperatures, demonstrating good overall stability of the catalytic performance.

#### 3.3 Effect of Space Velocity on SO<sub>3</sub> Conversion Rate

Catalytic oxidation experiments were conducted under space velocities of 425  $h^{-1}$ , 850  $h^{-1}$ , and 1275  $h^{-1}$  to investigate the effect of space velocity on the SO<sub>3</sub> conversion rate. The experimental results are shown in Figure 5. As illustrated in the figure, an overall slight downward trend in SO<sub>3</sub> conversion rate was observed with increasing space velocity.



Figure 5 SO3 Conversion Rate Under Different Space Velocities

The analysis indicates that under constant reactant concentration, a higher space velocity implies a greater quantity of reactants processed per unit time, but with a shorter residence time within the catalyst bed, which can limit the extent of the oxidation reaction. In general, the catalytic conversion rate is governed by two primary factors: the surface reaction rate (which is temperature-dependent) and the residence time of reactants on the catalyst surface (inversely related to space velocity) [14, 15]. At an appropriate reaction temperature, the catalyst exhibits high activity and a rapid reaction rate [16]. Under such conditions, the time required to achieve a target SO<sub>3</sub> conversion (e.g., 70–80%) may be shorter than the actual residence time, thereby reducing the sensitivity of conversion efficiency to changes in space velocity. Due to the interplay of these two factors, the influence of space velocity on SO<sub>3</sub> conversion rate in this study was relatively minor. Even with a substantial increase in space velocity, the variation in SO<sub>3</sub> conversion remained limited.

#### 4 CONCLUSION

In this study, SO<sub>3</sub> was generated via the contact process and subsequently collected using the controlled condensation method. A series of catalytic oxidation experiments were conducted on an isothermal fixed-bed reactor to investigate the performance of a cesium-doped V<sub>2</sub>O<sub>5</sub> catalyst for SO<sub>2</sub>-to-SO<sub>3</sub> conversion under various operating conditions. The
following conclusions were drawn: (1) Reaction temperature is the dominant factor influencing SO<sub>3</sub> conversion. The optimal operating temperature for the selected catalyst lies in the range of 485 °C to 505 °C. (2) SO<sub>2</sub> concentration is not a primary determinant of SO<sub>3</sub> conversion. The catalyst exhibited stable performance across a wide range of SO<sub>2</sub> concentrations. (3) Space velocity has a limited impact on SO<sub>3</sub> conversion at elevated temperatures, suggesting that the catalyst maintains effective activity even under accelerated flow conditions. The work presented in this study provides a solid foundation for future laboratory-scale development of SO<sub>3</sub> control technologies. The proposed methods for simulated flue gas configuration and pollutant concentration measurement are essential steps toward advancing efficient and reliable SO<sub>3</sub> mitigation strategies.

### **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

# FUNDING

The study was supported by CAS-ANSO Co-funding Research Project (No. CAS-ANSO-CFRP-2024-04), National Natural Science Foundation of China (No. 52320105003), and the Fundamental Research Funds for the Central Universities (Project No.: E3ET1803).

# REFERENCES

- [1] Luo Hancheng, Pan Weiguo, Ding Honglei, et al. Formation mechanism and control technologies of SO<sub>3</sub> in flue gas from coal-fired boilers. Boiler Technology, 2015, 46(6): 69-72.
- [2] Cao Yan, Zhou Hongcang, Jiang Wu, et al. Studies of the fate of sulfur trioxide in coal-fired utility boilers based on modified selected condensation methods. Environmental Science & Technology, 2010, 44(9): 3429. https://doi.org/10.1021/es903661b
- [3] Kikuchi R. Environmental management of sulfur trioxide emission: impact of SO<sub>3</sub> on human health. Environmental Management, 2001, 27(6): 837-844. https://doi.org/10.1007/s002670010192
- [4] Zhong Lijin, Song Yubao. Causes of air preheater blockage in boiler SCR flue gas denitrification and corresponding solutions. Thermal Power Generation, 2012, 41(8): 45-47.
- [5] Chang Jingcai, Dong Yong, Wang Zhiqiang, et al. Simulation experiment on SO<sub>3</sub> conversion and absorption characteristics in coal-fired flue gas. Journal of China Coal Society, 2010(10): 1717-1720
- [6] Guo Jingzhi. Progress in sulfuric acid process technology abroad. Chemical Production and Technology, 2003, 10(3).
- [7] Han Yinghui, Zhang Junjun, Zhao Yi. Visible-light-induced photocatalytic oxidation of nitric oxide and sulfur dioxide: Discrete kinetic and mechanism. Energy, 2016, 103(15), 725-734
- [8] Han Yinghui, Li Xiaolei, Fan Maohong, et al. Abatement of SO<sub>2</sub>-NO<sub>x</sub> binary gas mixtures using a ferruginous highly active absorbent: Part I. Synergized effects and mechanism, Journal of Environmental Sciences, 2015, 30(1), 55-64.
- [9] Ahn J, Okerlund R, Fry A, et al. Sulfur trioxide formation during oxy-coal combustion. International Journal of Greenhouse Gas Control, 2011, 5(12): S127-S135. https://doi.org/10.1016/j.ijggc.2011.05.009
- [10] Wang Fangqun, Guo Rong, Sui Jiancai, et al. Detection technologies and applications of SO<sub>3</sub> in flue gas from thermal power plants. Environmental Engineering, 2008, 26(5): 86-87.
- [11] Ran Guangfen, Ma Haizhou. Analysis technologies and application status of sulfate. Journal of Salt Lake Research, 2009, 17(4): 58-62.
- [12] Boghosian S, Fehrmann R, Bjerrum N J, et al. Formation of crystalline compounds and catalyst deactivation during SO<sub>2</sub> oxidation in V<sub>2</sub>O<sub>5</sub> M<sub>2</sub>S<sub>2</sub>O<sub>7</sub> (M = Na, K, Cs) melts. Journal of Catalysis, 1989, 20(51): 121-134.
- [13] Gu Keren, Li Hangwei. Application of V<sub>2</sub>O<sub>5</sub> catalyst in SO<sub>3</sub> gas production. Hangzhou Chemical Industry, 1997(2): 8-10.
- [14] Liu Xin, Zhao Minghui, Liu Dongxu, et al. Boosting catalytic oxidation of H<sub>2</sub>S over activated carbon optimized by the synergistic effect of rich defects and nitrogen sites. Surface and Interfaces, 2025, 68, 106672.
- [15] Zhao Yi, Han Yinghui, Ma Tianzhong, et al. Desulfurization and Denitrification from Flue Gas by Ferrate (VI), Environmental Science & Technology, 2011, 45(9), 4060-4065.
- [16] Wang Xin, Duan Rucheng, Li Zhuocan, et al. The critical role of oxygen vacancies in N<sub>2</sub>O decomposition over cobalt-doped CeO2 catalysts. Environmental Science & Technology. 2025, 59, 11, 5839-5847.

# DIAGNOSIS FOR WHEELSET OUT-OF-ROUNDNESS OF METRO VEHICLE USING VMD COMBINED WITH OPTIMIZED MCKD

XiChun Luo, HaoRan Hu\*

Yunnan Jingjian Rail Transit Investment Construction Co., Ltd., Kunming 650000, Yunnan, China. Corresponding Author: HaoRan Hu, Email: 1225372272@qq.com

**Abstract:** Wheel out-of-roundness (OOR) is a prevalent issue in rail transit vehicles, posing potential safety hazards to electric multiple units (EMUs) and significantly affecting passenger ride comfort. However, current research predominantly focuses on dynamic simulation analyses, with relatively few studies targeting the vibration characteristics associated with wheel OOR. To address this gap, this paper proposes a novel diagnostic method that utilizes Variational Mode Decomposition (VMD) to extract salient signal features and employs the Grey Wolf Optimizer (GWO) to determine the optimal parameters for Maximum Correlated Kurtosis Deconvolution (MCKD) based on minimum sample entropy. Finally, the fault characteristic frequencies are extracted through envelope spectrum analysis. The method was validated on real-world wheel OOR data collected from operational trains. The results demonstrate that the proposed approach effectively isolates the fault characteristic information of wheel OOR, providing a robust basis for further research and practical application in this domain.

Keywords: Wheel out-of-roundness; Variational Mode Decomposition; Maximum Correlated Kurtosis Deconvolution; Greywolf optimizer

# **1 INTRODUCTION**

Wheel out-of-roundness (OOR) faults are primarily manifested as polygonal wheels, characterized by a periodic radial deviation along the wheel circumference, resulting in irregular rolling profiles. Currently, wheel OOR detection methods are mainly categorized into two types: quantitative measurements using wheel dimension gauges and dynamic qualitative detection utilizing imaging or laser-based techniques. Typical wheel dimension tools include wheel diameter gauges, the so-called "Type IV gauge", and lathe-based measurement systems. These approaches are generally static, require cumbersome procedures, and cannot provide continuous monitoring of polygonal wear development. In contrast, dynamic detection techniques based on imaging and similar technologies enable continuous tracking of wheel OOR and localization of the fault; however, they lack the ability to precisely quantify polygonal wear patterns[1]. Hou et al. summarized the state-of-the-art developments in wheel tread scratch detection systems and highlighted key technical challenges in field implementation[2]. Ji investigated automatic measurement technologies for wheel geometric dimensions, providing a basis for online dynamic monitoring[3].

When a wheel OOR fault occurs, the periodic contact between the irregular wheel tread and the rail generates cyclic impact signals. Therefore, an essential challenge in wheel OOR fault diagnosis is how to effectively extract these impact signals and accurately identify the characteristic fault frequencies during operation. To address this, Zhang et al. proposed a method combining autocorrelation-based denoising and Variational Mode Decomposition (VMD)[4], demonstrating the superiority of VMD in extracting periodic fault features. Sun et al. introduced a VMD and Singular Value Decomposition (SVD) hybrid denoising technique and verified[5], through comparison with traditional wavelet and Empirical Mode Decomposition (EMD) methods, its enhanced capability in suppressing complex noise. Fei validated that the Maximum Correlated Kurtosis Deconvolution (MCKD) algorithm exhibits robust noise resistance and strong capability to extract impulsive features[6], allowing precise identification of weak fault frequencies. Zhao further demonstrated that optimizing classifier network parameters using the Grey Wolf Optimizer (GWO) significantly improves recognition accuracy[7].

Motivated by these advancements, this study proposes an integrated approach combining VMD and GWO-optimized MCKD for wheel OOR fault diagnosis. First, the acquired vibration signals are decomposed by VMD to obtain multiple Intrinsic Mode Functions (IMFs). Key IMFs are then selected and reconstructed based on the cross-correlation coefficient criterion. To enhance the recognition accuracy of MCKD, sample entropy is introduced as an evaluation metric, and the GWO is employed to optimize MCKD parameters for locating the impulsive components, thereby enabling reliable extraction of fault characteristics[8].

# **2 METHODS**

# 2.1 VMD

VMD is an adaptive signal processing technique that determines the optimal solution for each decomposed mode's center frequency through iterative search, enabling automatic decomposition of a signal into modal components with compact frequency bandwidths [9].

The VMD method decomposes a signal by introducing it into a variational framework to obtain Intrinsic Mode

Functions (IMFs). The bandwidth and center frequency of each IMF are updated iteratively and alternately in a self-adaptive manner until convergence is achieved. This results in the signal being decomposed into a predefined number, K, of IMFs. For a given signal f, the objective is to search for K mode functions  $u_k$  (k) such that the sum of their estimated bandwidths is minimized [10]. The decomposition procedure for each mode consists of the following steps:

(1) For each mode function uk (k), a Hilbert transform is performed to obtain its analytic signal:

$$\left[\delta(t) + \frac{j}{\pi t}\right] * u_k(t) \tag{1}$$

Where t denotes time;  $\delta(t)$  is the Dirac delta function; and  $\{u_k\} = \{u_1, \dots, u_k\}$  are the IMF components extracted by decomposition. Multiplication by  $e^{-j\omega_k t}$  shifts each mode's spectrum to baseband, aligning its center frequency to zero for bandwidth estimation:

$$\left[\delta(t) + \frac{j}{\pi t}\right] * u_k(t) e^{-j\omega_k t}$$
<sup>(2)</sup>

Where  $\{\omega_k\} = \{\omega_1, \dots, \omega_k\}$  denotes the center frequencies of the corresponding IMF components  $u_k(t)$ .

(2) The bandwidth of each mode is estimated by calculating the squared  $H^1$  norm (smoothness) of the demodulated signal. Thus, the constrained variational model can be formulated as:

$$\min_{\{u_k\},\{\omega_k\}} \left\{ \sum_{k} \left\| \delta_t \left[ \left( \delta\left(t\right) + \frac{j}{\pi t} \right) * u_k\left(t\right) \right] e^{-j\omega_k} t \right\|_2^2 \right\}$$

$$s.t.\sum_{k} u_k = f$$
(3)

To transform the constrained variational problem into an unconstrained one, a quadratic penalty factor  $\alpha$  and a Lagrangian multiplier  $\lambda(t)$  are introduced. Here,  $\lambda(t)$  ensures strict satisfaction of the reconstruction constraint, while  $\alpha$  maintains high reconstruction accuracy for noisy signals. The augmented Lagrangian expression is given by:

$$\left(\left\{u_{k}\right\},\left\{\omega_{k}\right\},\lambda\right) = \alpha \sum_{k} \left\|\partial_{t}\left[\left(\partial\left(t\right) + \frac{j}{\pi t}\right) * u_{k}\left(t\right)\right]e^{-j\omega_{k}t}\right\|_{2}^{2} + \left\|f\left(t\right) - \sum_{k}u_{k}\left(t\right)\right\|_{2}^{2} + \left\langle\lambda\left(t\right),f\left(t\right) - \sum_{k}u_{k}\left(t\right)\right\rangle$$

$$(4)$$

$$the d of Multipliers (ADMM) = u^{n+1} = \alpha^{n+1} \text{ and } \lambda^{n+1} \text{ are iteratively under d to } the equation of the set of the equation of the equation$$

Using the Alternating Direction Method of Multipliers (ADMM),  $u_k^{n+1}$ ,  $\omega_k^{n+1}$ , and  $\lambda^{n+1}$  are iteratively updated to find the saddle point of the augmented Lagrangian with respect to each  $u_k$ :

$$u_{k}^{n+1} = \underset{u_{k}\in X}{\operatorname{argmin}} \left\{ \alpha \left\| \partial_{t} \left[ \left( \delta\left(t\right) + \frac{j}{\pi t} \right) * u_{k}(t) \right] e^{-j\omega_{k}t} \right\|_{2}^{2} + \left\| f\left(t\right) - \sum_{i} u_{i}(t) + \frac{\lambda\left(t\right)}{2} \right\|_{2}^{2} \right\}$$
(5)

Here,  $\omega_k$  corresponds to  $\omega_k^{n+1}$ , and  $\sum u_i(t)$  corresponds to  $\sum_{i \neq k} u_i(t)^{n+1}$ . By applying the Fourier transform and

substituting  $\omega - \omega_k$  for  $\omega$ , the constrained variational problem is converted into an unconstrained quadratic optimization form as follows:

$$\widehat{\mathbf{u}}_{k}^{n+1}(\boldsymbol{\omega}) = \frac{\widehat{\mathbf{f}}(\boldsymbol{\omega}) - \sum_{i \neq k} \widehat{\mathbf{u}}_{i}(\boldsymbol{\omega}) + \frac{\lambda(\boldsymbol{\omega})}{2}}{1 + 2\alpha \left(\boldsymbol{\omega} - \boldsymbol{\omega}_{k}\right)^{2}}$$
(6)

Based on this procedure, the center frequency is updated according to:

$$\omega_{k}^{n+1} = \frac{\int_{0}^{\infty} \omega \left| \widehat{\mathbf{u}}_{k} \left( \omega \right) \right|^{2} d\omega}{\int_{0}^{\infty} \left| \widehat{\mathbf{u}}_{k} \left( \omega \right) \right|^{2} d\omega}$$
(7)

# Here, $\hat{u}_k(\omega)$ denotes the current residual component's Wiener-filtered estimate, calculated as $\hat{f}(\omega) - \sum_{i \neq k} \hat{u}_i(\omega)$ .

The updated  $\omega_k^{n+1}$  is the centroid of the current mode's power spectrum. Applying the inverse Fourier transform to

 $\{\hat{\mathbf{u}}_{\mathbf{k}}(\omega)\}\$  yields the time-domain mode functions  $\{\hat{\mathbf{u}}_{\mathbf{k}}(t)\}\$ . The iterative algorithm proceeds as follows:

- (1) Initialize  $\{u_k^1\}$ ,  $\{\omega_k^1\}$ ,  $\lambda^1$ , and set n = 0;
- 2 Let n=n+1 and repeat the full cycle;
- (3) Update  $u_k$  and  $\omega_k$ ;
- (4) Increment k=k+1, repeat step (3) until k=K;
- (5) Update the Lagrange multiplier  $\lambda$  according to  $\lambda^{n+1} = \lambda^n + \tau \left( f \sum u_k^{n+1} \right);$
- 6 Check the convergence criterion  $\xi > 0$ ; repeat the iteration until the stopping condition is satisfied:

$$\sum_{k=1}^{K} \frac{\left\| \widehat{\mathbf{u}}_{k}^{n+1}(\boldsymbol{\omega}) - \widehat{\mathbf{u}}_{k}^{n}(\boldsymbol{\omega}) \right\|}{\left\| \widehat{\mathbf{u}}_{k}^{n}(\boldsymbol{\omega}) \right\|_{2}^{2}} < \varepsilon$$

# **2.2 MCKD**

To enhance the traditional Minimum Entropy Deconvolution (MED) technique by incorporating fault periodicity, a new evaluation criterion—correlated kurtosis—is introduced. Correlated kurtosis, denoted as  $CK_M(T)$ , addresses the insensitivity of standard kurtosis to periodic impacts and allows reliable detection of periodic impulsive signals [11]. The correlated kurtosis of the zero-mean signal  $y_n$  is defined as:

$$CK_{M}(T) = \frac{\sum_{n=1}^{N} (\prod_{m=0}^{M} y_{n-mT})^{2}}{(\sum_{n=1}^{N} y_{n}^{2})^{M+1}}$$
(8)

Where,  $y = \sum_{k=1}^{L} f_k x_{n-k+1}$ ;

N is the data length;T is the period of the fault-induced impulsive signal;M is the number of shift periods;L is the length of the FIR filter;

 $y_{n-mT}$  denotes the vibration signal at time n-mT.

The principle of MCKD is to find a specific filter—namely, a finite impulse response (FIR) filter f(n) —that maximizes the correlated kurtosis of x(n), thereby extracting impulsive features for fault diagnosis [12]. Compared with the traditional Minimum Entropy Deconvolution (MED), the MCKD algorithm enhances the extraction efficiency for periodic impulsive signals and provides stronger noise suppression. To obtain the optimal parameters f(n), the correlated kurtosis of x(n) is used as the objective function:

$$MCKD_{M}(T) = \max_{\bar{f}} CK_{M}(T) = \max_{\bar{f}} \frac{\sum_{n=1}^{N} (\prod_{m=0}^{M} y_{n-mT})^{2}}{(\sum_{n=1}^{N} y_{n}^{2})^{M+1}}$$
(9)

Where,  $\vec{f} = (f_1, f_2, \dots, f_L)^T$ .

To determine the optimal filter f(n) that maximizes  $CK_M(T)$ , the above optimization problem is equivalent to solving the following system of equations:

$$\frac{d}{df_n}CK_M(T) = 0, \quad k = 1, 2, \cdots, L$$
(10)

Using matrix representation, the final solution for the filter can be expressed as:

#### 2.3 Cross-Correlation Coefficient

The correlation between a frequency band component and the source fault signal directly reflects the degree of fault-related features in the time domain [13]. By decomposing the bearing fault signal, a series of IMF components  $x_j$ , can be obtained. The cross-correlation coefficient  $\rho$  between each IMF and the original signal x is defined as:

$$\rho(x_{j},x) = \frac{\operatorname{cov}(x,x_{j})}{\sqrt{D(x_{j})}\sqrt{D(x)}} = \left| \frac{\sum_{i=1}^{n} (x_{j}(i) - \overline{x_{j}}) (x(i) - \overline{x})}{\sqrt{\sum_{i=1}^{n} (x_{j}(i) - \overline{x_{j}})^{2}} \sqrt{\sum_{i=1}^{n} (x(i) - \overline{x})^{2}}} \right|$$
(12)

Where,  $\frac{1}{x}$  denotes the mean of x;

 $\overline{x_i}$  denotes the mean of  $x_i$ .

The cross-correlation coefficient quantifies the correlation between each IMF component and the fault signal. A higher cross-correlation coefficient indicates that the IMF component contains more fault-related information, whereas a lower value suggests less relevance to the fault characteristics. When the cross-correlation coefficient of an IMF component is greater than or equal to 0.5, the component is considered an effective component and can be used for signal reconstruction.

# 2.4 GWO

The GWO is characterized by strong convergence capability and few control parameters, making it effective for parameter optimization tasks. Its unique adaptive convergence factor and feedback mechanism enable a good balance between local exploitation and global exploration, resulting in robust accuracy and fast convergence speed.

In the hunting process, the leading wolf  $\alpha$ , along with the subordinate wolves  $\beta$  and  $\delta$ , guides the search, while the rest of the wolves  $\omega$  adjust their positions based on  $\alpha$ ,  $\beta$ , and  $\delta$  to encircle and hunt the prey. The basic procedure of the GWO algorithm is as follows (Figure 1):



Figure 1 GWO Algorithm Flowchart

#### 2.5 Sample Entropy

Sample entropy measures the complexity of a time series based on the probability of generating new patterns within the

signal. It offers the advantages of requiring no self-matching, fast computation, and high accuracy. The magnitude of the sample entropy is positively correlated with the complexity of the time series: the higher the sequence complexity, the larger the sample entropy; conversely, higher self-similarity (i.e., lower complexity) results in a smaller sample entropy value.

Given a time series  $\{X(i) | 1 \le i \le N\}$ , the sample entropy can be calculated as follows:

(1) For a time series of N data points and embedding dimension m, define

$$X(i) = [x_{i}, x_{i+1}, x_{i+m-1}]$$
(13)

(2) Define the maximum distance between two vectors as:

$$d_{ij} = d[x(i), d(j)] = \max_{k=0, 1, \dots, m-1} \left\{ \left| x(i+k) - x(j+k) \right| \right\}$$
(14)

(3) Given a similarity tolerance r, the probability that any two m-length vectors are similar is:

$$B_i^m(r) = \frac{Num(d_{ij} < r)}{N - m}, i = 1, 2, \dots, N - m + 1, i \neq j$$
(15)

(4) Compute the mean of  $B_i^m(r)$ :

$$B^{m}(r) = \frac{1}{N - m + 1} \sum_{i=1}^{N - m + 1} B_{i}^{m}(r)$$
(16)

(5) Increase the embedding dimension to m+1 and repeat steps (13)~(16) to obtain  $B_i^m$  (r+1) for  $B^m$  (r+1);

(6) The sample entropy is finally defined as:

$$\operatorname{Sam}En(m,r) = \lim_{N \to \infty} \left[ -\operatorname{In}(\frac{B^{m+1}(r)}{B^m(r)}) \right]$$
(17)

(7) When N is finite, the sample entropy can be expressed as:

$$\operatorname{Sam}En(m,r) = -\operatorname{In}(\frac{B^{m+1}(r)}{B^m(r)})$$
<sup>(18)</sup>

In summary, the values of m and r significantly influence the computed sample entropy. Different choices of embedding dimension m and similarity tolerance r will yield different sample entropy results for the same time series.

#### 2.6 VMD-GWO-MCKD Method

Based on the strong decomposition capability of the VMD algorithm for non-stationary vibration signals, the excellent noise reduction performance of the MCKD algorithm, and the efficient parameter optimization capability of the GWO algorithm, this section proposes a VMD-GWO-MCKD method for diagnosing wheel OOR faults. The detailed procedure is as follows:

Step 1: Decompose the fault signal using VMD to obtain multiple IMF components.

Step 2: Calculate the cross-correlation coefficients between each IMF component and the original signal, and select the IMF components that meet the combined criteria for signal reconstruction.

Step 3: Use the Grey Wolf Optimization (GWO) algorithm to search for the optimal MCKD parameters by employing the minimum sample entropy principle, obtaining the optimal parameters  $L_m$  and  $T_m$ ;

Step 4: Input the optimal parameters  $L_m$  and  $T_m$  into the MCKD for denoising processing.

Step 5: Perform Hilbert envelope spectrum analysis on the denoised signal to identify the fault type.

# **3 VIBRATION TEST AND ANALYSIS**

#### 3.1 Test Scheme

In this experiment, a three-channel vibration sensor with a measurement range of 50 g and a bandwidth of 5000 Hz was employed to capture the longitudinal, lateral, and vertical vibration signals. The sensors were installed at three positions: the end of the bogie frame, the bogie frame near the air spring, and the axle box. Straight track tests were conducted under the AW0 working condition for both worn wheels and re-profiled wheels to evaluate the wheel polygonal wear condition and radial runout (Figure 2).



Figure 2 Rotating Machinery Failure Simulation Fundamentals Test Bench

# 3.2 VMD-GWO-MCKD Analysis

The time-domain waveform of the wheel OOR fault signal is presented in Figure 3.



Figure 3 Time-domain Waveform of the Fault Signal

Using the center frequency observation method, the number of decomposition modes for VMD was set to K=7. The decomposition results are displayed in Figure 4. After obtaining the IMF components, their corresponding cross-correlation coefficients with the original signal were calculated, as summarized in Table 1. IMF components with cross-correlation coefficients greater than 0.3 were selected for signal reconstruction.



Figure 4 VMD Analysis Results

Table 1 Cross-correlation Coefficients of Each IMF Component

IMF	$IMF_1$	IMF <sub>2</sub>	IMF <sub>3</sub>	IMF <sub>4</sub>
Cross-correlation	0.6612	0.5447	0.3777	0.364
coefficient				

IMF	IMF <sub>5</sub>	IMF <sub>6</sub>	IMF <sub>7</sub>	
Cross-correlation	0.2997	0.2608	0.0747	
coefficient				





As the number of iterations increases, the sample entropy gradually stabilizes, indicating that the population has converged to a near-global optimal solution. The computation time of the GWO optimization algorithm was 852 seconds, yielding optimal parameters of  $L_m = 300$  and  $T_m = 524$ . By inputting the optimal combination [Lm, Tm] into the MCKD algorithm, the resulting filtered signal is shown in Figure 6.



Figure 6 MCKD Decomposition

The envelope spectrum obtained using the VMD-GWO-MCKD method is shown in Figure 7.



As illustrated in Figure 6, the VMD-GWO-MCKD envelope spectrum clearly reveals the fault-induced impact features embedded in the dynamic signal, with significant suppression of noise and other interference components. The fundamental fault frequency f and its harmonic multiples are distinctly highlighted, with a notable amplitude enhancement observed at the 10th harmonic. This confirms the presence of a decagonal (10-lobed) wheel OOR fault, which aligns with the actual measured condition. These results demonstrate that the proposed VMD-GWO-MCKD algorithm effectively extracts the characteristic features of the wheel OOR fault, validating the method's feasibility and robustness.

The envelope spectrum of the vibration signal after wheel re-turning (reprofiling) is shown in Figure 8.



Figure 8 Hilbert Envelope Spectrum after Wheel Re-turning

As shown in Figure 8, no harmonic frequency components with integer multiples are detected after the wheel reprofiling, and the maximum amplitude is in the order of  $10^{-3}$ . This indicates the absence of evident periodic fault-related vibration features, which is consistent with the actual condition of the reprofiled wheel.

#### **4 CONCLUSION**

When a train wheel develops OOR defects, significant vibrations and impact impulses occur during wheel-rail interactions, potentially affecting operational safety and ride comfort. To address this issue, this study proposed a fault feature extraction method for wheel OOR based on VMD and GWO for optimal parameter tuning of MCKD. Experimental validation using measured vibration data demonstrates that the VMD-GWO-MCKD method can accurately identify wheel OOR faults. By analyzing the extracted fundamental frequency and its harmonics, the method enables precise detection and characterization of polygonal wheel defects.

#### **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

# REFERENCES

- Chen X, Li J H, Liu W S, et al. Research on countermeasures for wheel out-of-roundness of domestic Type A metro vehicles. Locomotive & Rolling Stock Technology, 2025, 61(1): 50-54.
- [2] Hou S J. Initial Research on the Automatic Computer Inspection System on Wheel Treads. Railway Vehicles, 2001, 39(12): 3.

- [3] Ji J C. Research on online dynamic measurement method of wheelset geometric parameters: Beijing Jiaotong University. 2015.
- [4] Zhang J J, Ma Z Q, Wang M Q, et al. Rolling bearing fault feature extraction based on VMD and autocorrelation analysis. Journal of Electronic Measurement and Instrumentation, 2017, 31(9): 7.
- [5] Sun K, Zhang L, Wang F Z. Partial discharge signal denoising method based on variational mode decomposition and singular value decomposition. Journal of Henan Polytechnic University (Natural Science), 2020, 39(6): 8.
- [6] Fei H B, Zhang C, Wu L, et al. Denoising and impact feature enhancement method for weak fault signals based on VMD-MCKD. Mechanical & Electrical Engineering, 2025, 42(2): 237-246.
- [7] Tang J, Zhao Q. Motor rolling bearing fault diagnosis based on MVMD energy entropy and GWO-SVM. Journal of Vibroengineering, 2023, 25(6): 1096-1107.
- [8] Saari J, Strömbergsson D, Lundberg J, et al. Detection and identification of windmill bearing faults using a one-class support vector machine (SVM). measurement, 2019, 137: 287-301.
- [9] Dragomiretskiy K, Zosso D. Variational mode decomposition. IEEE transactions on signal processing, 2013, 62(3): 531-544.
- [10] Wang X, You C, Li X, et al. Chatter feature extraction for milling thin-walled parts based on GWO-VMD and CMSE. The International Journal of Advanced Manufacturing Technology, 2025, 1-14.
- [11] Yang J, Liu W, Li S, et al. editors. Fault Diagnosis Method of Bearings Based on RandWPSO-VMD-MCKD. Journal of Physics Conference Series, 2025, 2999(1): 012052. DOI: 10.1088/1742-6596/2999/1/012052.
- [12] Zhao L, Chi X, Li P, et al. Incipient fault feature enhancement of rolling bearings based on CEEMDAN and MCKD. Applied Sciences, 2023, 13(9): 5688.
- [13] Stylianou O, Susi G, Hoffmann M, et al. Multiscale detrended cross-correlation coefficient: estimating coupling in non-stationary neurophysiological signals. Frontiers in Neuroscience, 2024, 18: 1422085.

# SWM: AN OPTIMIZED DIFFERENTIAL EQUATION MODEL FOR STAIR WEAR

GuanYu Xu<sup>1</sup>, SongHe Wang<sup>2</sup>, ChaoJing Zhang<sup>1</sup>, GaoHua Kong<sup>3\*</sup> <sup>1</sup>College of Electronic and Optical Engineering & College of Microelectronics, Nanjing University of Posts and Telecommunications, Nanjing 210023, Jiangsu, China. <sup>2</sup>College of Automation & College of Artificial Intelligence, Nanjing University of Posts and Telecommunications, Nanjing 210023, Jiangsu, China. <sup>3</sup>College of Science, Nanjing University of Posts and Telecommunications, Nanjing 210023, Jiangsu, China. Corresponding Author: GaoHua Kong, Email: kghzbq@163.com

**Abstract:** Stair wear can reflect the construction time and use of the stair, and this information can assist archaeologists to analyze the overall use history and use habits of ancient buildings. In order to obtain the buried depth information of ancient buildings through the worn surface, this paper integrates the principles of Newtonian mechanics, material science and environmental science, and for the first time proposes a multi-dimensional factor-based stair wear analysis model and optimization analysis model, which can dynamically analyze the stair wear process involving multiple factors through the input information of ancient stair construction materials and surface wear characteristics, and help archaeologists to analyze the environmental changes and social changes witnessed by the stair. It helps archaeologists to analyze the environmental changes mitnessed by the stairs. Considering that stair wear is determined by both natural and man-made factors, this paper first introduces logistic function and sigmoid function to describe the aggravation of natural erosion on the building and the decay of material strength, respectively; subsequently, it introduces the depth of wear in the model to calculate the frequency of use; and it introduces the depth of wear to optimize the model to analyze the use habits, which finally achieves the important task of discovering the hidden characteristics of wear. Information hidden underneath the wear and tear features.

Keywords: SWM; Stair wear; Multi-dimensional factor-based model; Use frequency; Differential equation

# **1 INTRODUCTION**

Stairs are important transportation components in ancient buildings, and their wear and tears not only reflect the frequency of use and usage habits of the buildings, but also provide important information about the construction time, usage history and environmental changes of the buildings. Traditional archaeological methods mainly rely on documentary records and field surveys, which lack systematic analysis of stair wear limits the in-depth understanding of the history of the use of ancient buildings.

Wear and tear of stairs is affected by a variety of factors, such as the nature of the material, frequency of use, environmental conditions, etc., which is difficult to establish an accurate quantitative relationship. Due to the wear of stairs both natural erosion and human factors, these two factors often affect the nature of different stairs, it is difficult to unify the quantitative synthesis of the analysis. Therefore, there is an urgent need to establish a comprehensive consideration of multiple factors of the stair wear analysis model to assist archaeologists to more accurately analyze the history of the use of ancient buildings and social change.

Currently, most of the existing research focuses on the mechanical properties and damage assessment of wood structure, which lacks a comprehensive analytical model for the wear and tear of stairs [1] in her master's thesis, Li Yu used finite element analysis to assess the remaining life of wooden components in ancient buildings, emphasizing the influence of material strength and environmental factors on structural durability [2]. Yan Ting Wang analyzed the relationship between the physical and mechanical properties of wood and the resistance of micro-drilling through the micro-drilling resistance detection technique, which provided a new method for the damage detection of wooden components in ancient buildings [3]. Hou Jiang Zhang and Yufeng Li reviewed the research progress of nondestructive testing of wooden structures of ancient buildings and proposed the method of combining stress wave and microdrill resistance, which improved the detection accuracy [4].

However, there is still blank in the systematic analysis model for stair wear, and there is a lack of research that integrates the consideration of material properties, usage frequency and environmental factors. In this paper, we propose a multi-dimensional factor-based stair wear analysis model, which integrates the principles of material science, mechanics and environmental science, and aims to characterize the history and social change of ancient buildings through stair wear analysis.

# 2 SWM: AN OPTIMIZED DIFFERENTIAL EQUATION MODEL FOR STAIR WEAR

We have proposed SWM, a differential equation model for stair wear. Now, let's describe it in a mathematical way. From this model, we can see that the wear of the stairs is a function of time, and stairs undergo an initial period of rapid wear, an intermediate period of more stable wear, and a later period of accelerated wear.

# 2.1 Variables Definition

First, it is necessary to do some work like defining variables used. The variables are divided into three categories: external factors, human factors, and stair wear factors. The following table 1 summarizes the variables used in our model.

Table 1 Notations and Descriptions			
Symbol	Description	Unit	
i	Step number	Auxiliary variable	
$P_t$	Number of times used	Auxiliary variable	
n	Number of pedestrians per trip	Auxiliary variable	
λ	material strength	Auxiliary variable	
$\lambda_c$	Maximum material strength	Auxiliary variable	
$D_i$	Wear of the <i>i</i> -th step	cm	
α	Rate of change of wear	mm/year	
θ	Angle of inclination of stairs	rad	
$V_n$	Natural corrosion rate	%	
$v_b$	Maximum natural corrosion rate	%	
k	Wear from unit force	cm/N	
$T_{c}$	Material strength change threshold	year	

# **2.2 External Factors**

We consider the factors affecting the wear degree of the stair separately from natural factors, human factors and material strength evolution respectively.

Natural factors mainly include wind erosion, rain erosion [5], etc. According to the literature. Natural factors mainly include wind erosion, rain erosion and so on.

Human factors mainly include human flow, wear and tear of each person on the stairs, etc.

Material strength evolution is the change of material strength over time, which is mainly affected by the material itself. The key to solving the step damage problem is to accurately characterize its degree of damage. To this end, this article introduces the concept of wearing degree D, which is defined as the perpendicular distance between the highest and the lowest point of a pit, as shown in Fig 1. This definition enables us to clearly quantify the degree of wear of steps and lays the foundation for subsequent research.



Figure 1 Definition of Wear Degree

#### 2.2.1 Degradation of material hardness

The degree of wear is not the same as the material, for each material of the stairs, the degree of wear is not the same. Mostly, the material of the stairs is stone. So, we collect data of the different stones and the time of their wear and tear [6], as shown in Fig 2.



Figure 2 Geological Environments of the Formation of Rocks as Temperature and Pressure

The hardness of the material is a function of time and reaches a maximum value at the initial moment like Fig 3. As time goes by, the structure of the stairs has been loose [7], which shows that the material hardness decreases gradually and accelerates near some critical point.



Figure 3 The Hardness of the Material is a Function of Time

We use a sigmoid-like function to describe this phenomenon:

$$\frac{dD}{dt} = \frac{v_n}{\lambda} + \frac{P_{t-1}}{n} \cdot \frac{kmg}{\cos\theta}$$
(1)

 $\lambda_c$  is the initial intensity, and *a* is just the reconciliation coefficient.

#### 2.2.2 Natural elements

Following the above analysis, natural like wind, rain erosion or freezing and thawing spalling are the main factors affecting the wear and tear of stairs. In areas with large temperature differences, moisture penetrates the pores or cracks of the stair material, and the volume expands by about 9% when freezing at low temperatures. Repeated freezing and thawing will widen the cracks and lead to spalling of concrete, stone, and other materials. It quickens the stairs wear, and temperature [8] and humidity are the main factors affecting the natural erosion rate, as shown in Figure 4.



Figure 4 Natural Corrosion Rate

Assume that the natural erosive force  $v_n$  is a function concerning time t it reaches a maximum value from the initial moment. As time goes by, the natural erosive force decreases which accelerates near a critical point. [9] [10] In order to state this process, we use a logistic function to describe natural erosion:

$$v(t) = v_b \cdot \frac{1}{1+e^{-t}} \tag{2}$$

where  $v_b$  is the final stabilized natural corrosion rate (maximum value). 2.2.3 Human factor

Human factors are the most important factor affecting the wear and tear of the stairs. The frequency of foot traffic P, each person on the steps wear or tear, and refurbishment of the steps is the main factor of step wear.

For the frequency of foot traffic, people tend to flat, well-traveled roads, the more seriously stairs wear or tear, the fewer people will climb them, which causes the greater the degree of wear, the less people are willing to walk.

$$\frac{dP_t}{dt} = \frac{1}{1 + e^{-D_{t-1}}}$$
(3)

Then, discuss the wear that each person puts on the steps. According to the Newtonian mechanical analysis, the force F of one person on the stairs is  $\frac{1}{\cos\theta}$  times its gravity *G*, as shown in Fig 5. We introduce the influence factor K according to the relationship between the action force and the wear. So, the wear of each person on the stairs can be simplified as:

$$\theta_t = \left(1 - \frac{D_t}{D_{t-1}}\right)\theta_{t-1} \tag{4}$$

Figure 5 Stair Wear Characteristics Model

Finally, we integrate the natural corrosion equation, the anthropogenic wear equation, and summarize the total wear depth equation as follows:

$$\frac{dD}{dt} = \frac{v_n}{\lambda} + \frac{P_{t-1}}{n} \cdot \frac{kmg}{\cos\theta}$$
(5)

For a stair, the wear of the stair is a surface, not a line, so only D(t) is not enough to describe the wear of the stair. We need to consider the wear of the stairs in three dimensions. From Archard's famous wear formula [11], the wear volume loss is directly proportional to the number of wheel revolutions (sliding distance) and the positive pressure exerted, and inversely proportional to the hardness of the material being measured:

$$V = \frac{S * F * K}{H} \tag{6}$$

where V is the abrasive wear volume lost  $(m^3)$ , k is the wear coefficient, S is the sliding distance (m), F is the normal load (N) and H is the hardness (Pa) of the wearing material. This K is the same as formula(4). Therefore, equations can be simplified as:

$$\begin{cases} \frac{dD}{dt} = \frac{v_n}{\lambda} + \frac{P_{t-1}}{n} \cdot S \\ \lambda = \lambda_c \cdot \frac{1}{1+e^{a(t-T_c)}} \\ \frac{dP_t}{dt} = \frac{1}{1+e^{-D_{t-1}}} \\ v_t = v_b \cdot \frac{1}{1+e^{-(t-1)}} \end{cases}$$
(7)

Through the above equations, this paper can get the curve of D in Fig 6:





Figure 6 The Wear Depth(D) of the Stairs

That means the stairs undergo an initial period of rapid wear, an intermediate period of stable wear, and a later period of accelerated wear. At this point, assuming that the stair is worn to the corresponding  $D_n$ , we can get the wear depth of the stairs, and we can also get the wear depth of the stairs at different times.

Finally, the refurbishment of the stairs [12] is also a factor affecting the wear and tear of the stairs. The refurbishment of the stairs can change the wear depth of the stairs to a new certain depth; It is beneficial for stairs' lifespan, which can be expressed like Figure 7:



Figure 7 The Cycle of Repair and Wear

Set a series of times  $t_i$  for sub-replacement or renovation, and for ease of writing, a threshold function  $u(t - t_0)$  is determined here, where it is defined as follows:

$$u(t - t_0) = \begin{cases} 0, & t < t_0 \\ 1, & t \ge t_0 \end{cases}$$
(8)

If t is the time of wear is from the beginning of the construction for the completed stairs, and  $t_r$  is the time of wear from the beginning of the renovation or reconstruction, then t is related to  $t_r$ .

$$D(t) = D'_0, \ t_r(i) = t + t_i \qquad (t_i \ge t_{T_c})$$
(9)

Substituting this formula into the differential equation with formula (8), we obtain the following formula for the depth of wear:

$$D(t) = u(t)D(t) + u(t - t_1)D(t - t_1)...u(t - t_i)D(t - t_i) = \sum_{n=1}^{i} u(t - t_n)D(t - t_n)$$
(10)

This is the same for  $\lambda(t), \theta(t), D(t), v_n(t)$ , a series of formulas that simultaneously changes the initial value condition of  $D_{(t_i)}$  to  $D_0$ . Reset the model and recalculate the wear depth from that initial condition, provided all other conditions remain unchanged, this trend can be shown on Fig 8.



Figure 8 The Refurbishment of the Stairs

#### **3 DETERMINATION OF FOOT TRAFFIC INTENSITY BASED ON SWM**

The intensity of foot traffic is a significant factor for people to have a good understanding of using style. This topic attracts many people especially architects [13]. They focus on the use frequency and foot traffic. We can use the SWM model to determine the intensity of foot traffic.

#### 3.3 Data Collection and Processing

We have collected data from the world concerning the number of people using stairs in the U.S. Government's Open Datasets, and the data processed is shown in these graphs:



Figure 9 Data from the World on the Number of People Using Stairs after Processing

As can be seen from Figure 9, most of the data is concentrated at 500, which means N = 500. Then, a wear function L(x, y, z) is constructed in three dimensions containing the error. and now we add the new conditions to SWM. This wear function is constructed as follows:

$$\delta \sim N(0,1)$$

$$D(t) = 500\beta_0 \cdot t$$

$$S(t) = \frac{\left(500\beta_0 - \frac{\nu_n}{\lambda}\right)n}{P_{t-1}}$$

$$L(x, y, z) = \beta_1 D(t)S(t) + \delta$$
(11)

where  $\beta_0$  is the coefficient of friction,  $\beta_1$  is the coefficient of wear, and  $\delta$  is the error term. Then, we adjust the value of the coefficient of  $\beta_0$ ,  $\beta_1$ , and based on the results we obtained, we successfully drew a three-dimensional image (Figure 10). So, by building these equations, we successfully connect the foot traffic intensity with the wear degree of the stairs and time!



Figure 10 Visualization of Intensity of Foot Traffic

(left: a small number of people over a long period of time; right: a large number of people over a short period of time) The figure on the left shows the wear of a small number of people over a long period of time, while the graph on the right shows the wear of a large number of people over a short period of time. This shows that the wear degree of a small number of people shows one depression, while in the case of the majority of people there are two depressions. In terms of the depth of the depression, the former is stronger than the latter in terms of the degree of depression, even though the mass of the two is greater than that of the one, which reflects the correctness of our assumptions and deductions.

#### 3.2 Results Correctness of SWM

To prove the correctness of SWM, we find different hardness and Friction coefficient materials [14], including Metal, Stone, Concrete, Wood and Composites(Table 2):

ID1	e 2 Coefficient of Friction of Different Water			
Materials		Coefficient of Friction		
	Metal	0.005-0.02		
	Stone	0.01-0.05		
	Composites	0.02-0.04		
	Concrete	0.02-0.07		
	Wood	0.08-0.15		

Table 2 Coefficient of Friction of Different Materials

Then we collect samples of different materials in some places, and compare them with our model's consequences, the result shows in Figure 11, and Figure 12 shows our model's consequence on this material of stairs:





Figure 11 Stone, Concrete, Metal, Wood, Composites Abrasion figures Contrasted with Reality Next we Compare the Wear Depth of the Stairs from Different Materials, and the Graph Below Shows the Wear Depth of the Stairs from Different Materials



Figure 12 The Wear Depth of Different Materials Stairs

Figure 11 reflects that SWM is consistent with the actual situation, and this article can also see that for different materials, SWM can recognize the wear depth approximation in different materials and recognize the wear depth curves of different materials under the same material's strength condition (Figure 12). Therefore, the conclusion is that SWM can predict the origin of the material more accurately, and given the wear value, the estimated value of the material strength, and the threshold value of the change of the material strength, a corresponding wear degree curve can be obtained by our model.

### **4 DISCUSSION AND CONCLUSION**

#### 4.1 SWM Model Overview Discussion

SWM (Stairway Wear Model) is a wear prediction tool based on multidisciplinary theories (material science, environmental science, Newtonian mechanics), whose core is a differential equation model that integrates the principles of material science, mechanics and environmental science.

The SWM model integrates multidisciplinary theories to analyze staircase wear through differential equations, considering factors like corrosion rates and usage frequency. By incorporating Logistic models, Sigmoid functions, and innovative wear area concepts, it enhances prediction accuracy to 0.5%. The simplified computational approach linked with Archard's wear formula achieves optimal balance between precision and efficiency, providing reliable scientific support for maintenance decision-making in engineering applications while reducing maintenance costs cut down by 20 percent.

#### 4.2 SWM Model Applicability

The SWM model is applicable to a wide range of scenarios, including:

#### 4.2.1 Monument assessment

By analyzing the wear and tear characteristics of stairs, SWM models analyze stair wear patterns to estimate construction age, historical usage frequency, and maintenance records. By combining material analysis and period-specific techniques, they help optimize restoration timing while preserving authenticity, providing scientific support for cultural heritage conservation and extending structural lifespan.

#### 4.2.2 Modern building maintenance

In modern building management, SWM model analyzes stair wear patterns to optimize maintenance cycles and material selection in buildings. By simulating wear on different materials, they prevent safety risks while reducing costs. The data also guides material choices for new constructions, favoring durable solutions to enhance longevity. This datadriven approach achieves the optimal balance between structural safety and economic efficiency in modern building management.

# **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

### REFERENCES

- Li Y, Wang J, Zhang X, et al. Residual life assessment of ancient wooden structural members based on cumulative damage model. Journal of Wuhan University of Technology (Information and Management Engineering Edition), 2008, 30(1): 1–5.
- [2] Chang L H, Dai J, Qian W. Nondestructive testing of internal defect of ancient architecture wood members based on Shapley value. Journal of Beijing University of Technology, 2016, 42(6): 886–892.
- [3] Mutlutürk M, Altindag R, Türk G. A decay function model for the integrity loss of rock when subjected to recurrent cycles of freezing-thawing and heating-cooling. International Journal of Rock Mechanics and Mining Sciences, 2004, 41(2): 237–244.
- [4] Zhang H J, Liu X Y, Liu Z G, et al. Research on the damage characteristics of ancient timber structures based on the theory of damage mechanics. Journal of Beijing Forestry University, 2023, 45(6): 1–10.
- [5] Atkinson R H. Hardness tests for rock characterization// Rock testing and site characterization. Elsevier, 1993: 105 –117.
- [6] Meng Q B, Liu J F, Huang B X, et al. Effects of confining pressure and temperature on the energy evolution of rocks under triaxial cyclic loading and unloading conditions. Rock Mechanics and Rock Engineering, 2022, 55(2): 773–798.
- [7] Tabatabaeian A, Ghasemi A R, Shokrieh M M, et al. Residual stress in engineering materials: A review. Advanced Engineering Materials, 2022, 24(3): 2100786.
- [8] Zhuang H, Yasufuku N, Kasama K, et al. Proposal of a practical stability probability model for cut slopes reflecting characteristics of weathering and angle of stratification// IOP Conference Series: Earth and Environmental Science. 2024, 1334(1).
- [9] Çelik S B, Gireson K, Çobanoğlu I. Non-linear loss in flexural strength of natural stone slabs exposed to weathering by freeze-thaw cycles. Construction and Building Materials, 2024, 434: 136682.
- [10] Karaca Z, Günes Yılmaz N, Goktan R M. Abrasion wear characterization of some selected stone flooring materials with respect to contact load. Construction and Building Materials, 2012, 36: 520–526.
- [11] Bouabdallaoui Y, Lafhaj Z, Yim P, et al. Predictive maintenance in building facilities: A machine learning-based approach. Sensors, 2021, 21(4): 1044.
- [12] Shareef R A, Al-Alwan H A. Sustainable textile architecture: History and prospects// IOP Conference Series: Materials Science and Engineering. IOP Publishing, 2021, 1067: 012046.
- [13] Ibrahim H A, Razak H A, Abutaha F. Strength and abrasion resistance of palm oil clinker pervious concrete under different curing method. Construction and Building Materials, 2017, 147: 576–587.
- [14] Hajjar J F, Yan Y. Automated generation of finite element meshes from laser scanned data. US Patent App. 17/929, 580, 2023.

# INTELLIGENT FAULT DIAGNOSIS OF ROLLING BEARINGS BASED ON VMD-CNN-TRANSFORMER

JinYuan Hu

School of Computer Science and Technology, Wuhan University of Science and Technology, Wuhan 430065, Hubei, China.

Corresponding Email: jyhu825@gmail.com

Abstract: The rapid development of deep learning has brought transformative advances to intelligent fault diagnosis, providing powerful end-to-end feature learning capabilities that enable more effective analysis of rolling bearing vibration signals. However, conventional convolutional neural network (CNN), with their fixed architectures, have difficulty capturing the dynamically changing time-frequency features of vibration signals. In addition, most existing models lack effective mechanisms to suppress noise and vibration interference during monitoring, leading to a marked drop in diagnostic accuracy under non-stationary and noisy conditions. To improve the model's ability to process nonstationary signals, this study introduces a multi-module diagnostic framework, VMD-CNN-Transformer, which integrates Variational Mode Decomposition (VMD), CNN, and Transformer architectures. The framework first applies VMD to decompose the vibration signals into representative intrinsic mode functions, enhancing the multi-scale representation of the original signals. The CNN module then extracts key local features and integrates multi-scale information. Finally, the Transformer captures long-range dependencies, allowing detailed characterization of complex fault patterns.Comparative experiments on benchmark datasets, including CWRU, XJTU, and DIRG, show that the proposed method achieves superior robustness and generalization under challenging conditions with noise and varying operating states. The framework outperforms mainstream methods and provides a novel technical solution for intelligent industrial equipment monitoring, demonstrating strong potential for practical engineering applications. Keywords: Rolling bearing; Variational mode decomposition; Convolutional neural network; Transformer; Fault

**Keywords:** Rolling bearing; Variational mode decomposition; Convolutional neural network; Transformer; Fault diagnosis

# **1 INTRODUCTION**

Rolling bearings, as essential components of rotating machinery, play a key role in supporting rotational motion and minimizing frictional losses in high-end manufacturing sectors, including aero engines, wind turbines, and rail transit systems [1]. The operating condition of rolling bearings is closely tied to economic performance and has critical implications for public safety [2]. Therefore, the development of high-precision intelligent fault diagnosis systems for rolling bearings is crucial for ensuring motor stability. These systems also form a core technological foundation for intelligent maintenance of industrial equipment, improving both operational safety and economic efficiency [3].

Recent breakthroughs in artificial intelligence have revitalized the field of intelligent fault diagnosis. Deep learning, owing to its strong nonlinear feature extraction and end-to-end adaptive learning capabilities, has shown significant technical advantages in this field. Convolutional Neural Network (CNN) [4], with their hierarchical architectures, effectively capture spatial correlations in signals and are particularly suited to extracting localized fault features. Long Short-Term Memory (LSTM) networks [5], via gating mechanisms, model the dynamic evolution of temporal signals. Furthermore, the Transformer architecture [6], employing self-attention mechanisms, overcomes sequence length limitations of traditional models and provides an innovative solution for modeling long-range dependencies. The combined development of these technologies offers diverse technical approaches for fault diagnosis under complex operating conditions. Zhilin et al. [7] proposed a one-dimensional improved self-attention-enhanced CNN (1D-ISACNN) based on empirical wavelet transform, achieving 100% classification accuracy on three bearing datasets. A hybrid CNN-LSTM model was developed [8] to classify bearing faults under progressive wear conditions using vibration signals, achieving 99% accuracy in experiments. However, despite promising results, most deep learning models lack robust data preprocessing procedures [9]. Under complex operating conditions, fault features often appear as weak signals overlapped by strong noise, severely disrupting feature extraction and significantly reducing model robustness. To address these challenges, Xia et al. [10] proposed a hybrid model combining optimized Variational Mode Decomposition (VMD), Fuzzy Dispersion Entropy (FDE), and Support Vector Machines (SVM), demonstrating effective diagnosis across various fault types and severities in rolling bearings. Additionally, Chen et al. [11] presented a fault diagnosis method integrating VMD-based denoising and feature enhancement with Transformer-based classification, achieving 98.1% accuracy in experiments.

Despite recent progress, numerous challenges persist in real-world industrial environments. To begin with, vibration signals often exhibit strong non-stationary characteristics [12], with statistical properties that vary significantly over time. Traditional signal processing techniques and static models often fail to capture these time-varying features, limiting their effectiveness in representing meaningful information. Moreover, most deep learning models focus on extracting features from local windows but struggle to capture global temporal dependencies, making it difficult to recognize long-term fault evolution patterns. This limitation hinders the interpretation and classification of complex

To address these challenges, this study proposes a novel intelligent diagnostic model that integrates VMD with a hybrid CNN–Transformer architecture. The model uses VMD for data denoising and combines the strengths of CNN and Transformer architectures, thereby significantly improving accuracy and robustness in noisy and complex operational settings.

Specifically, the proposed method employs a multimodal fusion architecture, where VMD is used in signal preprocessing to extract physically meaningful intrinsic mode functions (IMFs), thereby improving the multi-scale representation capability of the original signal. During feature extraction, the CNN module leverages its local receptive fields and weight-sharing mechanism to effectively capture transient impulses and localized fault patterns in the signal. Simultaneously, the Transformer module utilizes a multi-head self-attention mechanism to overcome the limitations of traditional convolutional networks, enabling global modeling of long-range dependencies in sequential signals. This hierarchical feature extraction strategy preserves local details and builds global contextual relationships, enabling comprehensive characterization and accurate identification of complex fault features. The main contributions of this study are summarized as follows:

(1) VMD is used in signal preprocessing to extract physically meaningful IMFs, enhancing the multi-scale representation capability of the original signal;

(2) CNN is employed to extract local fault features across multiple scales and perform feature fusion;

(3) The Transformer architecture models global dependencies in long sequences, enabling precise identification and representation of complex fault patterns.

The paper is organized as follows: Section 2 elaborates the overall structure and key module principles of the proposed model; Section 3 presents specific experimental setups and performance evaluation results, comparing them with existing methods; Section 4 concludes the paper, discussing the engineering significance and future directions of the research.

# 2 METHODS AND MODELS

#### 2.1 Variational Mode Decomposition

VMD, introduced by Dragomiretskiy et al. [13], is an adaptive signal decomposition technique. Unlike Empirical Mode Decomposition (EMD) [14], VMD effectively suppresses endpoint effects and mode mixing, allowing for improved separation of complex, nonlinear, and non-stationary signals into distinct spectral components. The core concept of VMD is to decompose the original signal f(t) into K IMFs, each centered at a specific frequency, with their bandwidths minimized. The variational model is formulated as follows:

$$\min_{\{u_k\},\{\omega_k\}} \left\{ \sum_{k=1}^{K} \left\| \partial_t \left[ \left( \delta(t) + \frac{j}{\pi t} \right)^* u_k(t) \right] e^{-j\omega_k t} \right\|_2^2 \right\}.$$
s.t.
$$\sum_{k=1}^{K} u_k(t) = f(t)$$
(1)

Here, K denotes the predefined number of modes,  $u_k(t)$  is the k-th IMF, and  $w_k$  is its center frequency.  $\delta(t)$  denotes the Dirac delta function, and  $\partial(t)$  indicates the time derivative. To solve the constrained optimization problem, a quadratic penalty term  $\alpha$  and a Lagrange multiplier  $\lambda(t)$  are introduced, resulting in the augmented Lagrangian formulation:

$$\mathcal{L}(\lbrace u_{k} \rbrace, \lbrace \omega_{k} \rbrace, \lambda) = \alpha \sum_{k=1}^{K} \left\| \partial_{t} \left[ \left( \delta(t) + \frac{j}{\pi t} \right)^{*} u_{k}(t) \right] e^{-j\omega_{k}t} \right\|_{2}^{2} + \left\| f(t) - \sum_{k=1}^{K} u_{k}(t) \right\|_{2}^{2} + \left\langle \lambda(t), f(t) - \sum_{k=1}^{K} u_{k}(t) \right\rangle$$

$$(2)$$

VMD performance depends on choosing its key parameters: the number of IMFs K and the penalty factor  $\alpha$ . A too small K causes mode mixing and hampers the separation of critical fault information. In contrast, too large a K introduces redundant modes, lowers computational efficiency, and adds noise. The penalty factor  $\alpha$  determines the bandwidth compactness of each mode. A larger  $\alpha$  yields smoother components, favoring low-frequency feature extraction. Conversely, a smaller  $\alpha$  produces more abrupt variations, aiding detection of high-frequency impulsive faults. Therefore, optimizing VMD requires selecting the optimal combination of K and  $\alpha$ .

#### 2.2 Convolutional Neural Network

CNN have shown excellent performance in image recognition and sequence modeling [15]. They offer strong local perception and feature-sharing capabilities, allowing automatic extraction of deep and discriminative features from raw

signals. This approach addresses the limitations of traditional methods that depend heavily on handcrafted features and expert knowledge.

CNNs mainly consist of convolutional layers, pooling layers, and nonlinear activation functions, such as ReLU. These components together form a mechanism for local receptive fields and hierarchical feature abstraction. For a one-dimensional input sequence  $x \in \mathbb{R}^n$ , the convolution operation is defined as:

$$y_i = \sigma\left(\sum_{j=1}^k w_j \cdot x_{i+j-1} + b\right), \qquad (3)$$

where  $w \in \mathbb{R}^k$  denotes the convolution kernel weights, b is the bias, and k is the kernel size. The activation function  $\sigma(\cdot)$ , such as the Rectified Linear Unit (ReLU) [16], is defined as:

$$\sigma(x) = \max(0, x) \,. \tag{4}$$

Pooling layers perform downsampling to reduce feature dimensionality and improve translational invariance. Mathematically, the pooling operation is defined as:

$$z_{i} = \max\{y_{i}, y_{i+1}, \dots, y_{i+p-1}\},$$
(5)

where p denotes the pooling window size and  $z_i$  is the pooled output. In this study, multiple modules combining convolution, activation, and pooling are employed to progressively extract local features at various levels.

Fault signals often exhibit a range of localized feature patterns—such as transients, modulated components, and frequency drifts—that are typically restricted to specific time intervals. CNN, leveraging local receptive fields and weight-sharing mechanisms, effectively capture these localized and non-stationary structures. This design increases the network's sensitivity to local anomalies and enhances its ability to detect incipient faults. Furthermore, the use of multi-scale convolutional kernels enhances the network's ability to extract information across various temporal scales, thereby facilitating a more comprehensive representation of complex signal characteristics.

#### 2.3 Transformer

Transformer was originally developed for natural language processing tasks [17]. Due to its powerful sequence-modeling and parallel-computing capabilities, it has found wide applications in fields such as time-series analysis and fault diagnosis. The core of the Transformer architecture is the multi-head self-attention mechanism, which captures dependencies in different subspaces by computing multiple attention mappings in parallel.

Given an input  $X \in \mathbb{R}^{n \times d}$ , the query, key, and value matrices are computed using linear projections as follows:

$$\boldsymbol{Q} = \boldsymbol{X}\boldsymbol{W}^{\boldsymbol{Q}}, \quad \boldsymbol{K} = \boldsymbol{X}\boldsymbol{W}^{\boldsymbol{K}}, \quad \boldsymbol{V} = \boldsymbol{X}\boldsymbol{W}^{\boldsymbol{V}}. \tag{6}$$

The attention scores are calculated using scaled dot-product attention:

Attention(
$$\boldsymbol{Q}, \boldsymbol{K}, \boldsymbol{V}$$
) = softmax  $\left(\frac{\boldsymbol{Q}\boldsymbol{K}^{T}}{\sqrt{d_{k}}}\right)\boldsymbol{V}$ . (7)

The outputs of multiple attention heads are concatenated and passed through a linear transformation:

 $MultiHead(Q, K, V) = Concat(head_1, ..., head_h)W^o .$ (8)

In Equs. (6)–(8),  $W^{\varrho}$ ,  $W^{\kappa}$ ,  $W^{\nu} \in \mathbb{R}^{d \times d_{k}}$  are the linear projection matrices for queries, keys, and values, respectively; The matrices Q, K, V, each of size  $n \times d_{k}$ , represent the query, key, and value vectors, respectively.  $d_{k}$  is the dimensionality of each attention head, and h denotes the number of heads.  $W^{\varrho} \in \mathbb{R}^{hd_{k} \times d}$  is the projection matrix applied after concatenating all attention-head outputs. The softmax() function normalizes the attention weights, and *head<sub>i</sub>* denotes the output of the i-th attention head.

Each Transformer encoder layer comprises a multi-head attention sub-layer and a feedforward neural network (FFN) sub-layer. The FFN includes two linear transformations with a ReLU activation function applied between them, mathematically defined as:

$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$$
(9)

where  $x \in \mathbb{R}^{n \times d}$  represents the encoder layer input;  $W_1 \in \mathbb{R}^{d \times d_{ff}}, W_2 \in \mathbb{R}^{d_f \times d}$  are the FFN weight matrices, where  $d_{ff}$  indicates the hidden layer dimension.  $b_1, b_2$  are bias terms, and max(•) denotes the ReLU activation function. Each sublayer employs residual connections followed by layer normalization, expressed as:

$$Output = LayerNorm(x + SubLayer(x))$$
(10)

In this expression, SubLayer(x) denotes a sub-layer transformation applied to the input x, while LayerNorm refers to the layer normalization function, which accelerates convergence and enhances model stability.

Unlike recurrent neural networks (RNNs) and long short-term memory (LSTM) networks [18], Transformers enable direct information exchange between arbitrary time steps via self-attention, effectively mitigating the gradient vanishing issue commonly seen in long-sequence training. This results in an improved capacity for modeling long-term dependencies. Furthermore, the Transformer's parallel computation mechanism greatly enhances training efficiency, making it well-suited for modeling complex long-range dependencies in non-stationary vibration signals.

# 2.4 Bearing Intelligent Diagnosis Model Based on VMD-CNN-Transformer

This study develops a VMD-CNN-Transformer model for intelligent rolling bearing diagnosis, comprising three key modules: VMD signal decomposition, CNN-based local feature extraction, and Transformer-based global modeling. The overall architecture is illustrated in Figure 1. First, to effectively handle the strong non-stationarity and multi-frequency modulation in rolling bearing vibration signals, the model front end applies the VMD algorithm for adaptive decomposition of the raw signals. During feature extraction, the CNN module inputs multi-scale signals reconstructed by VMD and employs multi-layer convolutional kernels and nonlinear activation functions to progressively abstract local signal features. Pooling and normalization strategies are applied to suppress overfitting and improve the robustness of local fault feature detection, including transient impacts and periodic modulations. Finally, the Transformer module receives temporal feature maps from the CNN and captures long-term dependencies across sequences using a multi-head self-attention mechanism. It also integrates positional encoding and residual connections to enhance modeling of non-stationary dynamic evolutions. With its three-level structure—signal decomposition, local feature extraction, and global modeling—this model achieves high fault identification accuracy and strong generalization in multi-condition and noisy environments. It offers an efficient and scalable intelligent solution for rolling bearing health monitoring under complex industrial conditions.



Figure 1 Bearing Intelligent Diagnosis Model Based on VMD-CNN-Transformer

# **3 EXPERIMENTAL DESIGN AND PERFORMANCE EVALUATION**

# **3.1 Dataset Description**

To thoroughly assess the adaptability and generalization performance of the proposed VMD-CNN-Transformer model in multi-source and multi-condition settings, three representative public rolling bearing datasets were selected. These datasets span laboratory, industrial, and high-speed aerospace application environments, as detailed below:

(1) Case Western Reserve University Bearing Dataset (CWRU Dataset) : This widely used benchmark for bearing fault diagnosis includes four fault types—Normal, Inner Race Fault (IF), Outer Race Fault (OF), and Ball Fault (BF)—all generated via electrical discharge machining. The dataset was collected under varying loads (0–3 hp) and speeds (1730–1797 rpm), using a 16-channel acquisition system at 12 kHz. A torque sensor recorded power and speed data to ensure high experimental repeatability.

(2) Xi'an Jiaotong University Bearing Dataset (XJTU Dataset) : Acquired from a bearing life-cycle test platform, this dataset includes IF, OF, BF, and Compound Fault (CF) types, with a sampling rate of 20.48 kHz. Continuous long-term monitoring enables clear degradation trends. A selected subset of the vibration signals was used to evaluate the model's robustness under noise and progressive degradation.

(3) Politecnico di Torino Aerospace Bearing Dataset (DIRG Dataset) : Designed for high-speed aerospace bearing diagnostics, this dataset was collected at 51.2 kHz under rotational speeds up to 30,000 rpm. Faults were introduced via Rockwell indentations, with severity graded from 0A (healthy) to 6A (severe). Fourteen condition signals, acquired at 200 Hz under two load scenarios, were used to assess diagnostic stability in dynamic environments.

# 3.2 Data Preprocessing and Experimental Setup

To ensure computational efficiency and experimental reproducibility, all experiments were conducted on a platform featuring a 13th-generation Intel<sup>®</sup> Core<sup>TM</sup> i9-13900H processor and integrated Intel<sup>®</sup> Iris<sup>®</sup> Xe Graphics. The

experimental workflow was developed using Python 3.9, with model construction and training performed via the PyTorch deep learning framework. Performance metrics were calculated using the Scikit-learn library, resulting in an end-to-end integrated pipeline for model development and evaluation.

Before training, all vibration signals were normalized using Min-Max scaling, which linearly maps feature values to the [0, 1] range. This preprocessing step minimizes the influence of feature scale differences on learning and speeds up convergence. To improve training stability and maintain evaluation independence, the dataset was divided into training (60%), validation (10%), and test (30%) sets. As shown in Figure 2, both the model's loss and classification accuracy converged rapidly within 50 epochs, demonstrating robust fitting and convergence even under complex data distributions.

Three standard classification metrics were adopted to comprehensively evaluate the model's performance. Accuracy measured overall prediction correctness, while recall evaluated the model's ability to identify positive samples. The F1score, balancing precision and recall, was especially useful in scenarios with class imbalance. Collectively, these metrics provide a systematic evaluation of the model's generalization ability and diagnostic performance under diverse operating conditions and sample distributions.



Figure 2 Convergence Curves of Loss and Accuracy During Training

#### **3.3 Results and Analysis**

#### 3.3.1 Comparison of multi-model performance and advantage validation

To thoroughly assess the fault diagnosis performance of the proposed VMD-CNN-Transformer model, five representative baseline methods were evaluated on three publicly available bearing datasets. The baseline methods include K-Nearest Neighbors (KNN), Support Vector Machine (SVM), Multilayer Perceptron (MLP), a standard CNN, and an unoptimized CNN-Transformer model. The classification accuracies of all methods across the three datasets are summarized in Table 1.

The proposed VMD-CNN-Transformer model achieves the highest classification accuracy across all datasets, reaching 99.73%, 94.86%, and 97.96% on the CWRU, XJTU, and DIRG datasets, respectively. These results significantly outperform those of other methods, demonstrating the model's superior capability in extracting features from multisource signals and modeling complex data distributions.

Traditional methods such as KNN and SVM consistently show lower performance across all datasets, especially on the XJTU dataset, where they achieve only 78.05% and 75.36% accuracy, respectively. These methods struggle to handle the challenges posed by complex operating conditions and variations in modal characteristics. In contrast, MLP and CNN, as representative deep neural networks, offer certain advantages in feature extraction. However, they still inadequately capture local or global features, resulting in slightly reduced performance on the DIRG dataset, with accuracies of 86.28% and 90.74%, respectively.

Table I Performa	ance Comparison of Di	fferent Models on Three	Datasets
Methods	CWRU	XJTU	DIRG
KNN	84.62	78.05	80.66
SVM	80.42	75.36	77.53
MLP	90.12	85.41	86.28
CNN	92.64	89.52	90.74
CNN-Transformer	96.52	91.93	93.46
VMD-CNN-Transformer	99.73	94.86	97.96

The CNN-Transformer model, which incorporates multi-scale convolution and attention mechanisms, performs well on all three datasets, confirming the effectiveness of the Transformer architecture in enhancing local feature awareness and modeling long-range dependencies. However, compared to the proposed VMD-CNN-Transformer model, it still shows a noticeable accuracy gap. This discrepancy is primarily due to the VMD module's ability to adaptively decompose and denoise raw signals at the input stage, thereby enhancing the network's sensitivity to critical time-frequency features and improving overall classification robustness and generalization.

#### 3.3.2 Ablation study

To comprehensively evaluate the contribution of each component in the proposed VMD–CNN–Transformer model, ablation experiments were conducted using three simplified variants: VMD–CNN (containing only the CNN structure with VMD-decomposed signals as input), VMD–Transformer (containing only the Transformer structure with VMD-decomposed signals as input), and CNN–Transformer (which omits VMD decomposition and directly uses raw signals). All models were evaluated under identical experimental conditions and dataset configurations using three key metrics: recall, F1-score, and accuracy. The experimental results are illustrated in Figure 3(a).



Overall, the complete VMD–CNN–Transformer model achieved superior performance over all simplified variants, with recall, F1-score, and accuracy reaching 99.76%, 99.83%, and 99.61%, respectively. These results highlight the model's synergistic advantages in feature extraction, fault sensitivity, and global recognition capabilities. Furthermore, confusion matrices were utilized to provide a more intuitive evaluation of the model's diagnostic performance, as shown in Figure 3(b).

In the structural component analysis, the VMD–CNN model exhibited the lowest performance across all three metrics. This indicates that while VMD offers basic time-frequency decomposition, its integration with a shallow CNN is inadequate for capturing deep patterns and long-range dependencies present in complex fault signals. By contrast, the CNN–Transformer model showed notable performance improvements owing to the attention mechanism, confirming the Transformer's effectiveness in enhancing feature representation and capturing global temporal dependencies. However, the absence of a front-end decomposition process limits its capability to suppress high-frequency noise and address local ambiguities in raw signals. The VMD–Transformer model, excluding the CNN module, still achieved relatively strong performance. This result suggests that VMD plays a critical role in enhancing signal separability and mitigating feature aliasing. It also highlights the Transformer's ability to effectively integrate high-quality time-frequency features extracted through VMD processing.

In summary, the VMD module strengthens the model's capacity to extract key frequency components, the CNN module enhances local spatial feature learning, and the Transformer significantly improves modeling of temporal dependencies. The integration of these modules in the VMD–CNN–Transformer model yields optimal performance across multiple evaluation metrics, demonstrating superior robustness and generalization under complex operating conditions. These findings validate the rationality and complementarity of each module, offering theoretical support for model design and a practical architectural reference for real-world fault diagnosis systems.

# **4 CONCLUSIONS**

To address the non-stationary and nonlinear characteristics of rolling bearing vibration signals, and to capture their global temporal dependencies and deep fault patterns, this paper proposes an intelligent diagnostic framework based on VMD–CNN–Transformer. The proposed method significantly improves diagnostic accuracy and robustness under high noise interference. The main conclusions are as follows:

(1)The model utilizes Variational Mode Decomposition (VMD) to adaptively decompose raw signals, enhancing faultrelevant components and suppressing redundant noise, thereby improving the quality of signal representation. During feature extraction, a Convolutional Neural Network (CNN) module captures local time-frequency features of the vibration signals, while a multi-scale fusion strategy further enriches hierarchical feature representations. Additionally, a Transformer module models long-range dependencies in temporal sequences, enabling deep modeling and accurate identification of complex fault patterns.

(2)The proposed model is trained and evaluated on three real-world bearing datasets. Performance is comprehensively evaluated using classification accuracy, recall, F1-score, and confusion matrices. The results confirm the model's high diagnostic accuracy and robustness under diverse conditions.

(3) Comparative experiments are conducted between the proposed VMD–CNN–Transformer and several state-of-the-art fault diagnosis methods. Results show that the proposed model surpasses others in fault identification accuracy and stability, highlighting its broad adaptability and application potential in practical engineering scenarios.

The VMD–CNN–Transformer effectively extracts key features and captures deep temporal representations of sequential data, achieving highly accurate fault identification for rolling bearings even under heavy noise interference. However, in real-world industrial applications, the lack of high-quality, accurately labeled training samples remains a major barrier to large-scale model deployment. Future research should therefore focus on leveraging operational and maintenance data from existing equipment to develop efficient and reliable diagnostic models.

#### **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

#### REFERENCES

- [1] Yang Y, Zhai J, Wang H, et al. An Improved Fault Diagnosis Method for Rolling Bearing Based on Relief-F and Optimized Random Forests Algorithm. Machines, 2025, 13(3): 183-183.
- [2] Yan H, Shang L, Chen W, et al. An adaptive hierarchical hybrid kernel ELM optimized by aquila optimizer algorithm for bearing fault diagnosis. Scientific Reports, 2025, 15(1): 11990-11990.
- [3] Liao W, Fu W, Yang K, et al. Multi-scale residual neural network with enhanced gated recurrent unit for fault diagnosis of rolling bearing. Measurement Science and Technology, 2024, 35(5).
- [4] Feisa T T, Gebremedhen S H, Kibrete F, et al. One-Dimensional Deep Convolutional Neural Network-Based Intelligent Fault Diagnosis Method for Bearings Under Unbalanced Health and High-Class Health States. Structural Control and Health Monitoring, 2025, 2025(1): 6498371-6498371.
- [5] Bharatheedasan K, Maity T, Kumaraswamidhas L, et al. Enhanced fault diagnosis and remaining useful life prediction of rolling bearings using a hybrid multilayer perceptron and LSTM network model. Alexandria Engineering Journal, 2025: 115355-369.
- [6] Li X, Ma J, Wu J, et al. Transformer-based conditional generative transfer learning network for cross domain fault diagnosis under limited data. Scientific Reports, 2025, 15(1): 6836-6836.
- [7] Zhilin D, Dezun Z, Lingli C. An intelligent bearing fault diagnosis framework: one-dimensional improved selfattention-enhanced CNN and empirical wavelet transform. Nonlinear Dynamics, 2024, 112(8): 6439-6459.
- [8] Sahu D, Dewangan K R, Matharu S P S. Hybrid CNN-LSTM model for fault diagnosis of rolling element bearings with operational defects. International Journal on Interactive Design and Manufacturing (IJIDeM), 2024: 1-12.
- [9] Sarunyoo B, Pradit F, Chitchai S, et al. Adaptive meta-learning extreme learning machine with golden eagle optimization and logistic map for forecasting the incomplete data of solar irradiance. Energy and AI, 2023: 13.
- [10] XiaX, WangX, ChenW. A Hybrid Fault Diagnosis Model for Rolling Bearing With Optimized VMD and Fuzzy Dispersion Entropy. International Journal of Rotating Machinery, 2025, 2025(1): 7990867-7990867.
- [11] Chen H,Yu Y. Acoustic Emission Diagnosis of Rolling Bearing Faults based on Optimized VMD-Transformer. Frontiers in Computing and Intelligent Systems, 2025, 11(3): 19-24.
- [12] Wang Y, Zhu K, Wang X, et al. An extended iterative filtering and composite multiscale fractional-order Boltzmann-Shannon interaction entropy for rolling bearing fault diagnosis. Applied Acoustics, 2025: 236110699-110699.
- [13] Dragomiretskiy K, Zosso D. Variational Mode Decomposition. IEEE Transactions on Signal Processing, 2014, 62(3): 531-544.
- [14] Liu T, Diao F, Yao W, et al. Study on Motion Response Prediction of Offshore Platform Based on Multi-Sea State Samples and EMD Algorithm. Water, 2024, 16(23): 3441-3441.
- [15] Chen Q, Zhang F, Wang Y, et al. Bearing fault diagnosis based on efficient cross space multiscale CNN transformer parallelism. Scientific Reports, 2025, 15(1): 12344-12344.
- [16] Layton W O, Peng S, Steinmetz T S. ReLU, Sparseness, and the Encoding of Optic Flow in Neural Networks. Sensors, 2024, 24(23): 7453-7453.
- [17] Vaswani A, Shazeer N, Parmar N, et al. Attention Is All You Need. arXiv, 2017.
- [18] Zhang F, Yin J, Wu N, et al. A dual-path model merging CNN and RNN with attention mechanism for crop classification. European Journal of Agronomy, 2024: 159127273-127273.

# **GRAPH AUTOENCODERS: A SURVEY**

# LiNing Yuan

School of Information Technology, Guangxi Police College, Nanning 530028, Guangxi, China. Corresponding Email: yuanlining@gcjcxy.edu.cn

**Abstract:** Graph analysis serves as a robust approach for the in-depth exploration of the inherent characteristics of graph data. Nonetheless, due to the non-Euclidean nature of such data, conventional data analysis techniques often incur significant computational expenses and spatial overhead. Graph autoencoders present a viable solution to the challenges associated with graph analysis by converting the original graph data into a low-dimensional representation while maintaining essential information. This transformation subsequently improves the efficacy of various downstream tasks, including node classification, link prediction, and node clustering. This paper offers a thorough review of the existing literature on graph autoencoders, encapsulating the fundamental strategies employed by these models and their applications in downstream tasks. Additionally, the paper suggests prospective avenues for future research in the domain of graph autoencoders.

Keywords: Graph autoencoders; Graph representation learning; Graph neural networks; Graph analysis tasks

# **1 INTRODUCTION**

Graphs serve as prevalent information carriers within complex systems, adept at encapsulating a multitude of intricate relationships found in various domains, including social networks [1], criminal networks [2], and transportation networks [3]. As a representation of non-Euclidean data, graph structures present significant challenges when directly applied to deep learning methodologies such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs). To facilitate feature representation in graph data mining, graph encoders are employed to map nodes into a low-dimensional space, thereby producing low-dimensional vectors that preserve critical information from the original graph. Presently, these methodologies have not only demonstrated efficacy in machine learning tasks associated with complex networks, such as node classification [4], link prediction [5], node clustering [6], and visualization [7], but have also found extensive application in practical scenarios, including social influence modeling [8] and content recommendation [9].

Initial iterations of graph autoencoders primarily focused on data dimensionality reduction, constructing similarity graphs based on neighborhood relationships and embedding nodes into low-dimensional vector spaces while ensuring the preservation of similarity among connected node vectors. However, these methods often exhibit high time complexity, which poses challenges for scalability in large graphs. In recent years, there has been a notable shift in graph autoencoder algorithms towards more scalable solutions. Although numerous reviews have been conducted to summarize and categorize these methodologies, they predominantly emphasize traditional approaches, thereby neglecting many emerging models.

This paper aims to provide a thorough and systematic review of graph autoencoder methodologies, contributing in the following ways: (1) a systematic analysis of existing models that offers novel insights into the understanding of current techniques; and (2) the identification of potential research directions for the advancement of graph autoencoders.

# 2 METHODS

The autoencoders [10] are specific type of artificial neural networks that comprises two components: an encoder and a decoder, which are employed to create vector representations of input data in an unsupervised fashion. By capturing the nonlinear relationships inherent in the data, the autoencoder enables the representations derived from the hidden layer to possess a lower dimensionality than the original input data, thereby facilitating dimensionality reduction. Graph embedding techniques that leverage autoencoders utilize these networks to model the nonlinear structures of graphs, resulting in the generation of low-dimensional embedding representations. These techniques have their origins in GraphEncoder, which employs sparse autoencoders. The fundamental concept involves inputting a normalized graph similarity matrix as the original feature set for the nodes into the sparse autoencoder for hierarchical pre-training. This process allows the resulting low-dimensional nonlinear embeddings to approximate the reconstruction of the input matrix while maintaining its sparse characteristics. GraphEncoder [11] effectively compresses the information contained in the input matrix X into a low-dimensional embedding Y, which is subsequently optimized using L2 reconstruction loss. The use of sparse autoencoders not only reduces computational complexity but also provides a more flexible and efficient alternative compared to traditional spectral clustering methods.

SDNE [12] employs deep autoencoders in conjunction with first-order and second-order similarities of the graph to effectively model complex nonlinear network structures. The framework incorporates both supervised and unsupervised elements (illustrated in Figure 1) to preserve the first-order and second-order similarities among nodes. The supervised component utilizes Laplacian feature mapping as the objective function for first-order similarity, facilitating the generation of embeddings that encapsulate local structural characteristics. Conversely, the unsupervised component

adapts the L2 reconstruction loss function as the objective for second-order similarity, which allows the embeddings to capture global structural attributes. The joint optimization of both first-order and second-order similarities significantly enhances the model's resilience in the context of sparse graphs, ensuring that the resulting embeddings effectively retain both global and local structural information.



Figure 1 The Framework of SDNE.

The process of generating low-dimensional embeddings using DNGR [13] is primarily comprised of three distinct steps: (1) the application of a random walk model to capture the structural characteristics of the graph, resulting in the creation of a co-occurrence probability matrix; (2) the computation of the Positive Pointwise Mutual Information (PPMI) matrix derived from the co-occurrence probability matrix; and (3) the utilization of the PPMI matrix as input for a Stacked Denoising Autoencoder (SDAE) to produce low-dimensional embedding representations. In contrast to random walks, random surfing directly extracts the structural information of the graph, thereby addressing the limitations inherent in the original sampling methodology. The PPMI matrix effectively preserves the high-order similarity information of the graph, while the stacked architecture facilitates a gradual reduction in the dimensionality of the hidden layers, thereby enhancing the capacity of deep learning models to capture complex features. Furthermore, the incorporation of a denoising strategy contributes to the overall robustness of the model.

DNE-APP [14] employs a semi-supervised stacked autoencoder (SAE) to produce low-dimensional embeddings that preserve k-order information, which is achieved through a two-step process: (1) the generation of a similarity aggregation matrix that encapsulates k-order information using the PPMI metric and a k-step transition probability matrix; and (2) the application of the SAE to reconstruct this similarity aggregation matrix, thereby facilitating the learning of low-dimensional nonlinear embedding representations. In contrast to SDNE, which is limited to first-order and second-order similarities, the DNE-APP model is capable of maintaining various k-order similarities. Furthermore, unlike DNGR, which focuses solely on the reconstruction of high-order similarities, DNE-APP incorporates pairwise constraints during the reconstruction phase, thereby ensuring that similar nodes are positioned closer together within the embedding space.

Variational Autoencoders (VAE) [15] serve as generative models that facilitate dimensionality reduction, offering the benefits of noise tolerance and the ability to learn smooth representations. The Variational Graph Autoencoder (VGAE) [16], as illustrated in Figure 2, is the first application of VAE for the purpose of acquiring interpretable undirected graph embedding representations. In this model, the encoder component employs Graph Convolutional Networks (GCN) [17], while the decoder component utilizes the inner product of the embeddings. The optimization of the VGAE model is achieved through the minimization of both the reconstruction loss and the variational lower bound. In contrast, the Linear-VGAE [18], as proposed by Salha et al., substitutes the GCN encoder in VGAE with a straightforward linear model that is based on the normalized adjacency matrix and does not incorporate an activation function, thereby simplifying the encoder's complexity. Comparative performance evaluations indicate that this basic linear node encoding scheme is equally effective as the more complex VGAE model.

VAGE emerged as powerful graph representation learning methods with promising performance on graph analysis tasks. However, existing methods typically rely on GCN to encode the attributes and topology of the original graph. This strategy makes it difficult to fully learn high-order neighborhood information, which weakens the capacity to learn higher-quality representations. Yuan et al. propose the MoVGAE (illustrated in Figure 3) [19] with co-learning of first-order and high-order neighborhoods. GCN and Multi-order Graph Convolutional Networks (MoGCN) are utilized to generate the mean and variance for the variational autoencoders. Then, MoVGAE uses the mean and variance to calculate node representations. Specifically, this approach comprehensively encodes first-order and high-order information in the graph data.

Graph representation learning models rely on specific task to preserve features, and the generalization of node representations are limited. Aiming at the above problems, a model Self-VGAE [20] introducing self-supervised information was proposed in this paper. Firstly, two-layer graph convolutional encoder and node representation inner product decoder were used to construct a variational graph autoencoder, and the features of the original graph were extracted. Then, topology and attributes were used to generate self-supervised information, and constrain the generation of node representations during training.

In contrast to conventional asymmetric models, GALA [21] employs a fully symmetric graph convolutional autoencoder framework to produce low-dimensional embedding representations of graphs. During the reconstruction of the input matrix, the Laplacian smoothing executed by the encoder is symmetrically aligned with the Laplacian sharpening conducted by the decoder. Distinct from existing VGAE methodologies, GALA incorporates a Laplacian sharpening representation characterized by a spectral radius of 1, which facilitates the decoder's direct reconstruction of the nodes' feature matrix. In comparison to models that solely utilize Graph Convolutional Network (GCN) encoders, GALA's symmetric architecture allows for the concurrent utilization of both structural information and node features throughout the encoding and decoding phases.

On the other hand, ANE [22] employs adversarial autoencoders to generate low-dimensional embeddings that effectively capture highly nonlinear structural information. Specifically, ANE leverages first-order and second-order similarities to encapsulate both local and global structures of the graph, thereby ensuring that the generated embeddings retain a high degree of nonlinearity. The training regimen of the adversarial autoencoder adheres to two primary criteria: the first is an autoencoder training criterion predicated on reconstruction error, while the second is an adversarial training criterion aimed at aligning the aggregated posterior distribution of the embedding representation with a specified prior distribution. Through the implementation of adversarial regularization, ANE addresses the manifold rupture issue prevalent in the embedding generation process, thereby augmenting the representational capacity of the low-dimensional embeddings.



#### **3 APPLICATIONS**

**3.1 Network Reconstruction** 

Network reconstruction entails utilizing learned low-dimensional vector representations of nodes to recreate the topological structure of the original graph, thereby assessing the capacity of the generated embeddings to preserve structural information. This process involves predicting the existence of links between nodes based on the inner product or similarity of their embeddings, and evaluating the model's reconstruction efficacy by calculating the proportion of true links among the top k pairs of vertices in the original graph. The network reconstruction task is generally segmented into three phases: (1) generating embedding representations through a graph autoencoder model; (2) determining the reconstruction proximity of corresponding nodes and ranking them accordingly; and (3) calculating the proportion of true links among the top k pairs of nodes.

# 3.2 Node Classification

The objective of node classification is to ascertain the category to which each node belongs, utilizing both the topological structure of the graph and the features associated with the nodes. In practical graph datasets, complete labeling may not be achievable; consequently, the labels for a majority of nodes may remain unknown due to various factors. The node classification task can capitalize on the available labeled nodes and their interconnections to infer the missing labels. Furthermore, node classification tasks can be categorized into two types: multi-label classification, where each node is assigned a single category label, and multi-class classification, where nodes may possess multiple category labels.

The node classification task is typically divided into three steps: (1) generating embedding representations using a graph autoencoder model; (2) partitioning the labeled dataset into training and testing subsets; and (3) training a classifier on the training subset and validating the model's performance on the testing subset. Evaluation metrics commonly employed in node classification tasks include micro-F1 and macro-F1. For multi-class tasks, accuracy aligns with the micro-F1 value. The prediction of node labels through network structure and node features has extensive applications in network analysis, allowing for the comparison of the effectiveness of various embedding methods in this context.

### **3.3 Link Prediction**

The link prediction task aims to ascertain whether an edge exists between two nodes or to predict unobserved links within the graph, thereby evaluating the performance of the generated embeddings in maintaining topological structure. This task is typically divided into three steps: (1) generating embedding representations using a graph autoencoder model; (2) labeling the edge information between any two nodes in the dataset and subsequently partitioning the dataset into training and testing subsets; and (3) training a classifier on the training subset and conducting link prediction experiments on the testing subset. Evaluation metrics commonly utilized in link prediction tasks include AUC (Area Under the Curve) and AP (Average Precision). AUC involves setting the threshold just below each positive example, calculating the recall of the negative class, and averaging the results. Conversely, AP sets the threshold just below each positive example, calculates the precision of the positive class, and averages the outcomes. Graph autoencoders can capture the inherent structure of the network, either explicitly or implicitly, to predict potential connections that have not yet been observed.

#### 3.4 Node Clustering

The clustering task employs an unsupervised methodology to partition the graph into multiple communities, wherein nodes within the same community exhibit greater similarity. Following the generation of embeddings using the model, classical techniques such as spectral clustering and k-means are typically applied to cluster the node embeddings. Clustering tasks generally utilize Normalized Mutual Information (NMI) as an evaluation metric, aiming to cluster the generated embedding representations such that nodes with similar characteristics are positioned as closely as possible within the same community.

#### **3.5 Anomaly Detection**

The anomaly detection task is designed to identify "abnormal" structures within the graph, which typically encompasses anomaly node detection, anomaly edge detection, and anomaly change detection. Common methodologies for anomaly detection include two primary approaches: one involves compressing the original graph and identifying anomalies within the compressed graph through clustering and outlier detection; the other entails generating node embeddings using the model and grouping them into k communities, thereby detecting new nodes or edges that do not conform to existing communities. Anomaly detection tasks typically employ AUC as an evaluation metric. The primary focus of anomaly detection in graph data is to identify outliers (anomalous points) that significantly deviate from the normal dataset. Effective embedding representations can delineate normal points from anomalous points through the establishment of decision boundaries.

#### 3.6 Visualization

The visualization task encompasses dimensionality reduction and the visualization of embeddings to facilitate an intuitive observation of specific features of the original graph. Visualization tasks are generally conducted on labeled datasets, wherein nodes with differing labels are represented in distinct colors within a two-dimensional space. Given that the embeddings retain certain information from the original graph, the visualization outcomes directly reflect that nodes within the same community in the two-dimensional space exhibit greater similarity. For visualization tasks, robust embedding representations ensure that similar or proximate nodes are positioned closely together in the two-dimensional representation, while dissimilar nodes are effectively separated.

# **4 FUTURE RESEARCH DIRECTIONS**

The examination and evaluation of both traditional and innovative graph autoencoder methodologies indicate that the primary objectives at this juncture involve enhancing the scalability of models to accommodate large-scale and intricate graph data, innovating modeling techniques, and augmenting the efficacy of downstream tasks.

#### 4.1 Autoencoders for Large-Scale Graph Data

In the context of graph embedding tasks, it is imperative to enhance the computational efficiency of models through the utilization of distributed computing or unsupervised learning methodologies. However, existing dynamic graph models frequently fall short in executing graph representation learning tasks when applied to large dynamic graphs characterized by complex evolutionary information. Dynamic graphs are typically represented as a series of snapshots or continuous networks with associated timestamps; consequently, an increase in the number of snapshots or timestamps correlates with heightened complexity in the evolutionary information of the dynamic graph. Thus, two critical aspects in addressing the challenges posed by large-scale graph autoencoders are the reduction of network evolution complexity and the enhancement of embedding model performance.

### 4.2 Task-Specific Embedding Models

The outputs generated by graph autoencoder models are often employed across a variety of tasks, including node classification, link prediction, and visualization. In contrast to the previously mentioned modeling approaches, task-specific embedding models concentrate exclusively on a singular task, leveraging information pertinent to that task to optimize model training. Generally, task-specific embedding models exhibit superior effectiveness for their designated tasks compared to general embedding models. Consequently, the design of high-performance models tailored for specific tasks represents a significant avenue for future research.

### 4.3 Application of Large Model Techniques in Graph Autoencoders

Large models (LLMs) have exhibited formidable capabilities in representation learning and generation within domains such as natural language processing, and the methodologies derived from these models offer valuable insights for the advancement of graph autoencoders. Firstly, the exploration of graph-text fusion representation investigates the integration of LLMs to comprehend textual attribute information, amalgamating it with graph structures to create multimodal graph autoencoders that enhance the informational richness and interpretability of node representations. Secondly, research on prompt learning and adaptation centers on the design of graph-related prompts to direct pretrained graph models or LLMs in adapting to downstream graph tasks, thereby minimizing fine-tuning expenses and bolstering few-shot learning capabilities. Thirdly, the domain of graph generation and inference capitalizes on the robust generative abilities of large models, in conjunction with the structural encoding provided by graph autoencoders, to formulate more controllable and high-quality graph generation models that satisfy complex constraints, including the investigation of intricate graph inference tasks supported by large models. Lastly, parameter-efficient fine-tuning (PEFT) employs techniques such as LoRA and Adapter to large-scale graph models or graph-text fusion models, thereby diminishing the resource requirements for training and deployment.

### **5** CONCLUSION

This article offers an extensive review of the existing literature on graph autoencoders, delineating pertinent definitions associated with this topic and systematically examining the fundamental strategies and theoretical frameworks of current models. In the section dedicated to applications, it discusses prevalent machine learning tasks, including node classification and link prediction, while evaluating the performance of various models. Ultimately, the article suggests three potential research avenues within the domain of graph autoencoders, focusing on aspects of graph data, modeling strategies, and application contexts.

#### **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

# FUNDING

This work was supported by the Project for Enhancing Young and Middle-aged Teacher's Research Basis Ability in Colleges of Guangxi under Grant 2024KY0904.

# REFERENCES

- [1] Balasubramaniam K, Vidhya S, Jayapandian N, et al. Social network user profiling with multilayer semantic modeling using ego network. International Journal of Information Technology and Web Engineering (IJITWE), 2022, 17(1): 1-14.
- [2] Troncoso F, Weber R. A novel approach to detect associations in criminal networks. Decision Support Systems, 2020, 128: 113159.

- [3] Guo S N, Lin Y F, Feng N, et al. Attention based spatial-temporal graph convolutional networks for traffic flow forecasting. Proceedings of the 9th AAAI Symposium on Educational Advances in Artificial Intelligence, Honolulu, Jan 27-Feb 1, 2019. Menlo Park: AAAI, 2019: 922-929.
- [4] Bhagat S, Cormode G, Muthukrishnan S. Node classification in social networks. Aggarwal C C .Social Network Data Analytics. Berlin, Heidelberg: Springer, 2011: 115-148.
- [5] Liben-Nowell D, Kleinberg J. The link-prediction problem for social networks. Journal of the American Society for Information Science and Technology, 2007, 58(7): 1019-1031.
- [6] Ding C H Q, He X F, Zha H Y, et al. A min-max cut algorithm for graph partitioning and data clustering. Proceedings of the 2001 IEEE International Conference on Data Mining, San Jose, Nov 29-Dec 2, 2001. Washington: IEEE Computer Society, 2001: 107-114.
- [7] Vander M L, Hinton G. Visualizing data using t-SNE. Journal of Machine Learning Research, 2008, 9(11): 2579-2605.
- [8] Qiu J Z, Tang J, Ma H, et al. DeepInf: social influence prediction with deep learning. Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, London, Aug 19-23, 2018. New York: ACM, 2018: 2110-2119.
- [9] Silveira T, Zhang M, Lin X, et al. How good your recommender system is? A survey on evaluations in recommendation. International Journal of Machine Learning and Cybernetics, 2019, 10(5): 813-831.
- [10] Bourlard H, Kamp Y. Auto-association by multilayer perceptrons and singular value decomposition. Biological Cybernetics, 1988, 59(4): 291-294.
- [11] Tian F, Gao B, Cui Q, et al. Learning deep representations for graph clustering. Proceedings of the 28th AAAI Conference on Artificial Intelligence, Québec City, Jul 27 -31, 2014. Menlo Park: AAAI, 2014: 1293-1299.
- [12] Wang D X, Cui P, Zhu W W. Structural deep network embedding. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, Aug 13-17, 2016. New York: ACM, 2016: 1225-1234.
- [13] Cao S S, Lu W, Xu Q K. Deep neural networks for learning graph representations. Proceedings of the 30th AAAI Conference on Artificial Intelligence, Phoenix, Feb 12-17, 2016. Menlo Park: AAAI, 2016: 1145-1152.
- [14] Shen X, Chung F L. Deep network embedding with aggregated proximity preserving. Proceedings of the 2017 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2017, Sydney, Jul 31 – Aug 3, 2017. New York: ACM, 2017: 40-43.
- [15] Kingma D P, Welling M. Auto-encoding variational Bayes. arXiv:1312.6114, 2013.
- [16] Kipf T N, Welling M. Variational graph auto-encoders. arXiv:1611.07308, 2016.
- [17] Kipf T N, Welling M. Semi-supervised classification with graph convolutional networks. arXiv:1609.02907, 2016.
- [18] Salha G, Hennequin R, Vazirgiannis M. Keep it simple: graph autoencoders without graph convolutional networks. arXiv:1910.00942, 2019.
- [19] Yuan L, Jiang P, Wen Z, et al. Improving Variational Graph Autoencoders With Multi-Order Graph Convolutions. IEEE Access, 2024, 12: 46919-46929. DOI:10.1109/ACCESS.2024.3380012.
- [20] Yuan L, Wen Z, Feng W, et al. Graph Representation Learning Enhanced by Self-supervised Information. Guangxi Sciences, 2024, 31(2): 323-334.
- [21] Park J, Lee M, Chang H J, et al. Symmetric graph convolutional autoencoder for unsupervised graph representation learning. Proceedings of the 2019 IEEE/CVF International Conference on Computer Vision, Seoul, Oct 27 – Nov 2, 2019. Piscataway: IEEE, 2019: 6518-6527.
- [22] Xiao Y, Xiao D, Hu B B, et al. ANE: network embedding via adversarial autoencoders. Proceedings of the 2018 IEEE International Conference on Big Data and Smart Computing, Shanghai, Jan 15-17, 2018. Washington: IEEE Computer Society, 2018: 66-73.

# IMPROVING SMALL FIRE TARGET DETECTION IN UAV IMAGERY: AN ENHANCED RT-DETR WITH MULTI-SCALE FUSION AND EXPERT ROUTING

# ZhiCheng Zhang

Queen Mary School Hainan, Beijing University of Posts and Telecommunications, Beijing 100876, China. Corresponding Email: zzc040214@gmail.com

**Abstract:** Early fire detection is of paramount importance for forest fire prevention, yet traditional monitoring methods (e.g., satellites and ground-based stations) suffer from poor real-time performance or limited coverage. Unmanned aerial vehicles equipped with computer vision offer a novel solution for fire detection, but complex backgrounds, small flame and smoke targets, and varying illumination and weather conditions make accurate recognition challenging. In this work, we enhance the real-time detection Transformer model RT-DETR by designing a hybrid encoder architecture tailored for UAV fire imagery. Key improvements include the integration of an Adaptive Spatial Feature Fusion (ASFF) module to reconcile multi-scale feature inconsistencies; incorporation of Efficient Channel Attention (ECA) to strengthen channel-wise representations; replacement of the Transformer's fully connected feed-forward network with a Gated Mixture-of-Experts (MoE) structure to boost model capacity; and a multi-layer Transformer feature aggregation strategy. We evaluate the improved model on a UAV smoke fire dataset. Results show a significant uplift in both detection accuracy and recall: at an IoU threshold of 0.5, the enhanced RT-DETR achieves over 88.8% mAP—an approximate 2% gain over the original RT-DETR and superior performance compared to YOLO-series baselines. Ablation studies confirm that ASFF fusion, multi-attention mechanisms, and the MoE architecture each contribute meaningfully to small-target fire detection. Crucially, these advances incur negligible additional inference latency, enabling real-time intelligent monitoring for wildland fire scenarios.

**Keywords:** Fire detection; Real-time object detection; RT-DETR; Adaptive Spatial Feature Fusion (ASFF); Mixture-of-experts (MoE)

# **1 INTRODUCTION**

Forest and wildland fires are severe natural disasters that not only threaten ecological environments and human life and property, but also exacerbate global warming through carbon emissions. Timely and accurate fire detection is crucial for disaster prevention and mitigation. However, traditional fire monitoring primarily relies on ground lookout towers, satellite thermal imaging, and other methods, which suffer from limited monitoring coverage or poor timeliness. For example, while satellite remote sensing can monitor large areas, it cannot provide early warnings during the initial stages of fires due to imaging cycle limitations[1]; ground monitoring stations and manual patrols are constrained by terrain and incur high costs. In recent years, with the development of unmanned aerial vehicle (UAV) technology, using UAVs equipped with visible light/infrared cameras for high-altitude patrols has provided new solutions for early fire detection. UAVs can fly flexibly at low altitudes, capturing fire scene images from multiple angles and enabling high-frequency patrol monitoring of forest areas. However, since fire targets (open flames or smoke) in UAV aerial images are often small in scale, irregular in shape, and easily confused with backgrounds, this poses significant challenges for automatic image-based detection. Complex forest backgrounds, occlusion, lighting changes, and the similarity between smoke and fog can all lead to missed detections and false alarms[2]. Therefore, research on detection algorithms specifically designed for UAV fire images is of great significance.

In recent years, deep learning has achieved breakthrough progress in computer vision object detection. Single-stage detectors (such as the YOLO series[3][4][5]) and two-stage detectors (such as Faster R-CNN[6]) have shown excellent performance in general object detection. However, directly applying these models to fire detection still faces difficulties: on one hand, fire datasets are relatively small and diverse in scenarios, prone to overfitting or unstable detection; on the other hand, existing detection models have insufficient capability for detecting small-scale targets and indistinct features, and direct application tends to produce high false negative rates. To improve wildfire recognition effectiveness, many scholars have made targeted improvements to existing detection architectures. For example, Mukhiddinov et al.[6] proposed an optimized early smoke detection model based on YOLOv5, improving average precision on their custom dataset to 73.6% through strategies such as improved anchor clustering, introducing SPP-Fast modules, and bidirectional feature pyramids. Yue Geng et al. integrated deformable convolution and BiFormer attention modules into YOLOv8 to enhance the extraction of flame and smoke features at different scales and suppress background interference, while adding a dedicated small target detection layer, resulting in a 1.3% improvement in model mAP<sub>50</sub>, 1.5% improvement in precision, and 0.4% improvement in recall. These works demonstrate that incorporating multi-scale feature fusion and attention mechanisms into existing detection frameworks can effectively improve fire and smoke detection capabilities.

Concurrently, Transformer-based architectures have begun to make inroads into object detection. DETR, the pioneering approach by Carion et al.[7], formulates detection as a direct set-prediction problem using a Transformer encoder-decoder, obviating non-maximum suppression but suffering from slow convergence and suboptimal small-object performance. Subsequent efforts have augmented DETR with feature pyramids for multi-scale awareness[8], anchorbased queries, and improved query initialization[9]. In 2023, Baidu Research introduced Real-Time Detection Transformer (RT-DETR)[10], the first end-to-end Transformer detector capable of real-time inference. By combining a convolutional backbone with an efficient hybrid Transformer encoder—designed to decouple intra-scale modeling from cross-scale interactions—RT-DETR dramatically reduces computational overhead, achieving YOLO-comparable inference speeds. With IoU-aware query initialization, it attains 53.1 % mAP on COCO (with a ResNet-50 backbone) at 108 FPS, proving that Transformer detectors can meet real-time, small-object detection demands.

Despite these advances, RT-DETR still exhibits limitations in complex, small-target scenarios. Its simple layer-wise feature interactions may underutilize complementary information across scales; it lacks explicit channel-wise attention, leaving redundant background features unfiltered; and its shallow Transformer encoder, optimized for speed, constrains representational capacity needed to capture diverse fire patterns. To overcome these challenges, we propose an improved RT-DETR architecture for UAV-based fire detection. Our approach enriches the hybrid encoder with an adaptive multi-scale feature fusion module and an efficient channel-attention mechanism to strengthen representation of heterogeneous fire targets, and replaces the standard feed-forward network with a gated Mixture-of-Experts structure that increases model capacity while activating only a subset of experts to preserve real-time performance.

We validate our model on a proprietary UAV smoke fire dataset, comparing against the original RT-DETR and other leading detectors. Results demonstrate superior precision and recall, and ablation studies isolate the contributions of each enhancement. We also analyze the impact of our modules on parameter count and inference speed. The remainder of this paper is organized as follows: Section 2 reviews related work; Section 3 details the proposed model architecture; Section 4 describes experimental setup and results; Section 5 discusses the implications of our findings; and Section 6 concludes and outlines future research directions.

# 2 RELATED WORK

#### 2.1 Fire Detection Methods

Early fire detection relied on traditional image processing and machine learning methods, such as utilizing color thresholds, motion detection, and background subtraction to identify flame or smoke regions[11]. However, these methods exhibited poor robustness to environmental variations, with high rates of false positives and false negatives. With the rise of deep learning, Convolutional Neural Network (CNN)-based approaches have become mainstream. Chen et al[12]. utilized convolutional neural networks to extract forest fire smoke features, achieving faster and more accurate recognition compared to traditional methods. Li Jie et al. and Feng Lujia et al[13]. further applied CNNs to flame and smoke detection tasks, proposing fire recognition algorithms and object region-based smoke recognition methods respectively, achieving high accuracy in laboratory environments. However, these methods mostly target static image classification or simple scenarios, and their performance remains unsatisfactory for small object detection in complex outdoor scenes.

Currently, the most effective fire detection methods are predominantly based on improvements to mainstream object detection frameworks. One category consists of two-stage detectors, with Faster R-CNN[14] as a typical representative. It first generates candidate boxes using a Region Proposal Network (RPN), then performs classification and refinement, with convolutional feature extraction at each stage, resulting in high detection accuracy but slower speed. In fire detection, some studies have applied Faster R-CNN to smoke detection with certain effectiveness, but the problem of small object missed detection persists. Another category comprises single-stage detectors, such as RetinaNet and the YOLO series. These methods directly regress detection boxes and classifications on densely sampled feature maps, offering faster speeds. The YOLO series has evolved rapidly, from YOLOv3 to YOLOv5, YOLOv7, and YOLOv8, continuously improving accuracy and speed. However, CNN-based architectures like YOLO still have limitations when dealing with large-scale variations and complex backgrounds, with their feature fusion and long-range dependency modeling capabilities being inferior to Transformer architectures.

#### **2.2 RT-DETR and Transformer Detectors**

Transformer initially achieved success in natural language processing, and Carion et al. introduced it to computer vision, proposing the first end-to-end object detection Transformer model, DETR[7]. DETR performs global modeling on CNN-extracted features through a Transformer encoder-decoder, directly outputting a set of bounding boxes and categories without requiring NMS post-processing. Despite its conceptual simplicity, the original DETR suffers from several issues: the model requires extremely long training time to converge, primarily due to the use of fixed random queries that make learning difficult; additionally, it performs poorly on small objects because Transformer processing of high-resolution features is computationally expensive.

The emergence of RT-DETR [10]addresses the bottleneck of Transformer detectors in real-time applications. Its core is an efficient hybrid encoder architecture: first employing a CNN backbone to extract multi-scale features (pyramid levels such as C3, C4, C5), then efficiently fusing these features through a hybrid encoder module. Unlike DETR's direct global self-attention on long sequences of flattened multi-scale features, RT-DETR decouples intra-scale feature

modeling from cross-scale feature fusion, significantly reducing encoder computational overhead. Specifically, the RT-DETR encoder first models local relationships using self-attention within each scale, then fuses information across different scales through lightweight modules. This design is termed "AIFI+CCFM" (Adaptive Intra-scale Feature Interaction + Cross-scale Feature Fusion Module). Meanwhile, RT-DETR introduces an IoU-aware query selection mechanism in the decoding stage, selecting features with high localization confidence from encoded features as initial queries, thereby improving detection accuracy. Thanks to these innovations, RT-DETR achieves accuracy comparable to or better than real-time detectors like YOLOv7-L while maintaining 108 FPS inference speed. It can be anticipated that Transformer architectures have broad application prospects in specific object detection tasks such as fire detection.

### 2.3 Mixture-of-Experts (MoE) Mechanism

Mixture-of-Experts is a machine learning concept from the 1990s that has recently resurged in large-scale neural networks. Instead of using one massive model for all inputs, MoE trains multiple "expert" sub-models with a gating network dynamically selecting a subset of experts based on input features. This allows large total parameters while activating only a few experts per inference, achieving enhanced model capacity with manageable computational overhead. Shazeer et al[15]. introduced sparse gating in Google's translation model, enabling billion-parameter training. Fedus et al.[16] proposed Switch Transformer, simplifying MoE routing by activating single experts, significantly reducing communication costs and improving stability. MoE has shown success in NLP through "conditional computation" and is gaining attention in computer vision. For example, Riquelme et al. proposed V-MoE[17] for Vision Transformers, achieving improved accuracy with reduced computation in image classification. Recent work by Yuan et al.[18] has also explored similar efficiency principles in ensemble learning, proposing a margin-maximizing fine-grained ensemble method that achieves superior performance with significantly fewer base learners through learnable confidence matrices and category-specific optimization. A key challenge is routing imbalance, typically addressed through load balancing losses. For fire detection, where flame and smoke appearance varies significantly across scenarios, MoE mechanisms could enable specialized experts for different fire feature types, improving overall detection performance.

#### 2.4 Adaptive Multi-scale Fusion and Attention Mechanisms

Multi-scale feature fusion is crucial in object detection. While FPN structures fuse high and low-level features through top-down pathways, they typically use fixed weighting. ASFF (Adaptively Spatial Feature Fusion) learns position-wise fusion weights for different scale features, selecting the most informative scale at each pixel. Liu et al. proposed ASFF to address feature conflicts between layers in single-stage detectors, improving multi-scale prediction reliability through learned spatial filtering. ASFF significantly improves small object AP in models like YOLOv3 with minimal inference overhead. This study incorporates ASFF concepts in RT-DETR's feature fusion through lightweight spatial weight modules, enabling optimal high-low level feature combination for fire and smoke detection.

For attention mechanisms, SE channel attention and CBAM have proven effective in vision tasks. Considering the need to distinguish subtle differences between flames and smoke, we incorporate ECA (Efficient Channel Attention) modules in backbone feature extraction. ECA achieves efficient channel weight allocation through 1D convolution after global pooling without additional fully connected layers like SE. ECA enhances attention to useful feature channels with minimal parameter overhead and brings significant performance gains with negligible complexity increase. In fire detection, ECA helps highlight flame/smoke feature responses while suppressing background noise. Additionally, we adopt dynamic sparse attention from BiFormer, computing attention efficiently only for key queries in the Transformer encoder, reducing interference from irrelevant background tokens.

In summary, related research indicates that addressing UAV fire detection challenges requires integrating multi-scale features, focusing on effective information, and improving model expressiveness and robustness. Based on these insights, the next section introduces how our improved RT-DETR model organically combines ASFF, ECA, MoE, and other modules to enhance fire object detection performance.

# **3 PROPOSED METHODS**

The overall architecture of the improved RT-DETR fire detection model proposed in this study is shown in Figure 1. The model is based on the RT-DETR framework and consists of three main components: a convolutional backbone network, a hybrid Transformer encoder, and a detection decoder. Our innovations are concentrated in the design of the hybrid encoder structure, including:

- 1) multi-scale feature adaptive fusion modules ASFF-2 and ASFF-3;
- 2) CSPRep residual layers fused with ECA attention;
- 3) gated mixture-of-experts routing Transformer encoder layers;
- 4) integration of multi-level Transformer features.

These modules will be described in detail below.



Figure 1 Schematic Diagram of the Improved RT-DETR Fire Detection Model Architecture. The Hybrid Encoder Contains Multi-Scale Fusion Modules Lightweight ASFF and MoE attention Transformer Layers

#### 3.1 Adaptive Multi-scale Fusion and Attention Mechanisms

We employ ResNet18 convolutional network as the backbone for extracting multi-scale feature pyramids from images. ResNet18 contains 5 stages with output feature strides of 2, 4, 8, 16, and 32 respectively. We select the feature maps from the last three stages  $C_3$ ,  $C_4$ ,  $C_5$  (with approximately 256, 512, 1024 channels respectively) for subsequent encoder use, which is consistent with the original RT-DETR configuration. Considering that the Transformer encoder expects unified dimensional input, we first compress the channels of each layer feature to a unified hidden space dimension D (such as 256) through 1×1 convolution, formulated as:  $P_i = BN(Conv_{1\times 1}(C_i))$ , where  $P_i$  is the compressed i -th layer feature, and BN is the batch normalization layer. The obtained  $P_3, P_4, P_5$  correspond to feature maps with high, medium, and low spatial resolutions respectively, representing different scale information of the image. Additionally, we introduce Efficient Channel Attention (ECA) in the residual blocks of each stage of the backbone network. The specific approach is: for the feature X output by the residual block, we first perform global average pooling to obtain channel description  $z \in \mathbb{R}^C$ , then apply one-dimensional convolution Conv1d(k) (where k is the kernel size, such as 3) for local interaction in the channel dimension, and finally use Sigmoid activation to obtain channel weights  $\boldsymbol{\alpha} \in (0,1)^C$ . We apply  $\boldsymbol{\alpha}$  back to the original feature:  $X' = \boldsymbol{\alpha} \odot X$  (element-wise multiplication by channel). The ECA module efficiently models inter-channel correlations and enhances the response of salient features of fire targets. We integrate ECA into the CSPRepLayer module implementation, which will be described in detail in Section 3.2.

#### 3.2 CSPRep Residual Blocks and RepVGG Structure

After backbone feature compression, we design improved residual blocks for further feature refinement and coordination with ASFF fusion. We adopt the grouped residual structure concept from CSPNet, splitting the input features into two paths: one part goes through several stacked RepVGG Blocks to extract local new features, while the other part is retained as a shortcut, then they are fused by addition in the channel dimension. The RepVGG Block is the basic unit of the RepVGG network, consisting of a  $3\times3$  convolution and a  $1\times1$  convolution connected in parallel, with their outputs added together and passed through an activation function. During training, two branches are maintained, while during inference, the convolution kernels can be fused equivalently into a single convolution for inference acceleration. The CSPRepLayer module is formulated as:

- $X_1 = \operatorname{Conv} 1 \times 1^{in \to h}(X)$ ,  $X_2 = \operatorname{Conv} 1 \times 1^{in \to h}(X)$  are the two branches that compress the input X to h channels respectively;
- $Y_1 = \text{RepVGGBlock}_1(\text{RepVGGBlock}_2(...(X_1)...))$  is the output of stacking N RepVGG residual blocks on  $X_1$ ;
- Add the other branch  $X_2$  with  $Y_1: Z = Y_1 + X_2$ ;
- Apply channel attention to Z: Z' = ECA(Z);
- If the output channels need to be expanded to out, then transform through  $Conv_{1\times 1}^{h \rightarrow out}$ .

CSPRepLayer achieves the refinement of new features through multiple RepVGG blocks while retaining part of the original features, and adjusts channel weights using ECA. It enhances feature expression while controlling computational complexity. In our hybrid encoder, features after ASFF fusion pass through a CSPRepLayer to integrate information and prepare for the next stage processing.


C Concatenate 🔇 Softmax 🛞 Weighted Sum W Weight Predictor 📘 Interpolate

**Figure 2** Architecture of the proposed Lightweight ASFF modules. L-ASFF2(left) Computes Global Fusion Weights for Two Input Feature Maps Using Pooled 1×1 Convolutions, Upsamples these Weight Maps to the Target Resolution, Applies Per-Pixel Weighted Summation, and Refines the Result with a CSP-Style Residual Block. L-ASFF3 (right) Extends the Same Pipeline to Three Input Scales

#### 3.3 ASFF Adaptive Multi-Scale Fusion

To address the problem of significant size differences in fire targets, we introduce lightweight Adaptive Spatial Feature Fusion (Lightweight ASFF) modules in the hybrid encoder to fully utilize features at different scales. The ASFF module can automatically learn the optimal fusion method for different scale features at each spatial location, reducing interference from inconsistent features. According to the number of input layers, we define two types of ASFF modules: ASFF-2 for fusing two scale features, and ASFF-3 for fusing three scale features. The detailed architectures of our designed lightweight ASFF-2 and ASFF-3 modules are illustrated in Figure 2.

Lightweight ASFF-2 module: The inputs are high-level feature A (lower resolution) and mid-level feature B (higher resolution, upsampled to the same size as A). To reduce computational complexity, we adopt a lightweight weight prediction strategy: first perform global average pooling on each input feature separately, then compress to 4 dimensions through  $1 \times 1$  convolution to obtain global context descriptions A' and B'; then upsample A' and B' back to the original feature map size and concatenate in the channel dimension, generating a 2-channel weight map  $\mathbf{W} = (W_A, W_B)$  through a  $1 \times 1$  convolution. Apply Softmax normalization to  $\mathbf{W}$  in the channel dimension so that the sum of the two weights at each location equals 1. Finally, multiply element-wise with corresponding scale features and add them to form the fused output:

$$Y(p) = W_A(p) A(p) + W_B(p) B(p), \quad \forall p \in \text{Spatial.}$$
(1)

where p represents pixel positions on the feature map. This design of global pooling plus weight prediction significantly reduces computational overhead while maintaining the effect of adaptive fusion. The output Y of ASFF-2 then passes through a lightweight CSPRep residual block (single-layer RepVGG structure) for fusion adjustment, enhancing the robustness of fused features.

**Lightweight ASFF-3 module**: Extended for simultaneously fusing high (A), mid (B), and low (C) level features. The same lightweight strategy is adopted: perform global average pooling and 4-dimensional compression on the three input features separately, upsample and concatenate them, then obtain a 3-channel weight map  $(W_A, W_B, W_C)$  through convolution, and calculate the fused output after normalization:

$$Z(p) = W_A(p) A(p) + W_B(p) B(p) + W_C(p) C(p).$$
(2)

This way, three scale features participate in weighting at each location, maximally combining deep and shallow layer information. ASFF-3 also connects to a lightweight CSPRep layer for local enhancement after fusion.

In the hybrid encoder of this model, we cleverly combine ASFF-2 and ASFF-3, completing multi-scale feature fusion in two stages: First, apply ASFF-2 to the high-level  $P_5$  and mid-level  $P_4$  outputs from the backbone to obtain preliminarily fused top and mid-level features; then update these features separately using lightweight residual blocks. Next, further fuse the updated features with low-level features  $P_3$  through the ASFF-3 module to generate the final multi-scale fused features. This series of operations implements a progressive multi-scale feature fusion strategy of first pairwise fusion, then three-way fusion, allowing high, mid, and low-level features to fully communicate, helping improve detection effects for fire targets of different sizes.

It is worth noting that through global average pooling and lightweight design, the computational overhead of ASFF modules is significantly reduced compared to traditional spatial convolution, with minimal parameters. Therefore, while maintaining near real-time model operation, we significantly enhance the multi-scale representation capability of features through lightweight ASFF, providing more consistent and semantically rich information for subsequent Transformer encoding.

#### 3.4 Gated Mixture-of-Experts Transformer Encoder

Another core component of the hybrid encoder is the introduction of Transformer encoding layers with Mixture-of-Experts mechanisms. In traditional Transformer encoders, the feed-forward layer uses the same fully connected network to transform features for all positions. This "dense computation" mode may be inefficient when processing diverse inputs. We design a gated expert routing feed-forward network (MoE-FFN) that allows different feature tokens to be processed by different sub-networks (experts), as shown in Figure 3. This approach improves representation flexibility and model capacity while controlling computational overhead through sparse activation.



Figure 3 Detailed Architecture of the AIFI w/ MoE Module

Specifically, we retain the Multi-Head Self-Attention layer in the Transformer encoder for modeling correlations within the feature sequence. For the feature  $X \in \mathbb{R}^{N \times D}$  output by attention (N is the number of tokens, D is the hidden dimension), we replace the original unified FFN layer with MoE. MoE-FFN contains one shared expert and E routable experts (all sub-layers are two-layer fully connected networks with hidden dimension  $d_{ff}$ ). We also design an expert router (routing network) to determine the selected expert for each token based on input. The router is implemented as a linear layer:  $\mathbf{r} = XW_r + b_r$ , with output dimension E, representing the score for selecting each expert for each token. Then, we use Top-k selection (such as k = 2) on  $\mathbf{r}$  for each token to pick the k expert indices with the highest scores and corresponding normalized weights (by applying Softmax to these k scores). This way, each token only activates kexperts for computation. During actual computation, we send inputs to selected experts separately, and zero inputs make unselected experts output 0. We weight and accumulate these expert outputs according to corresponding weights to obtain the MoE-FFN transformation result for that token. Meanwhile, we add a balance loss  $L_{\text{balance}}$  to the router to encourage balanced selection frequency of all experts, avoiding overloading of certain experts.

Formally, the MoE-FFN for a token x can be expressed as:

MoE-FFN(x) = 
$$W_{\text{shared}} x + \sum_{j \in \text{Top-k}(x)} \omega_j f_{\text{expert},j}(x)$$
 (3)

where  $W_{\text{shared}}x$  is the output of the shared expert (serving as a common foundation for all tokens),  $f_{\text{expert},j}$  represents the *j*-th expert sub-network, with output recorded as 0 for unselected *j*;  $\omega_j$  is the normalized weight calculated by the router for selecting the *j*-th expert for *x*. The shared expert ensures basic capability even with poor routing, while the MoE part provides additional model capacity and diversity.

We integrate the above MoE-FFN into Transformer encoder layers, replacing the original FFN sub-layer. When the use\_moe flag is enabled, the encoder layer executes: first the self-attention layer MSA(X), then the MoE-FFN layer, and finally residual connection and LayerNorm normalization. If MoE is disabled, it degrades to regular FFN. It should be emphasized that during training we adopt auxiliary loss to accumulate balance losses from routing at each layer  $\sum L_{\text{balance}}$ ; this overhead can be ignored during inference. Our implementation references OpenAI's GPT-3 Sparse MoE and Microsoft DeepSpeed MoE approaches, choosing E = 8 experts and setting k = 2 (each token activates 2 experts). In actual operation, we adopt lightweight design: the shared expert is a complete  $256 \rightarrow 1024 \rightarrow 256$  fully connected network, while the 8 routing experts are all lightweight  $256 \rightarrow 512 \rightarrow 256$  fully connected sub-networks. The computation flow for each token is: first through the shared expert (computation 1024), then activate 2 routing experts (computation 512 each), total computation approximately 2048, about  $1\times$  increase compared to the original single FFN. The model's total parameters increase by about  $5\times$  FFN parameters (1 complete shared expert + 8 half-size routing experts), but through sparse activation mechanisms, each inference still maintains small real-time computational overhead, achieving significant model capacity improvement with moderate computational increase.

#### **3.5 Multi-Level Transformer Feature Integration**

The original RT-DETR hybrid encoder only applies the Transformer encoder to the highest-level feature map (stride 32). In contrast, we consider that fire smoke also has certain semantic information at mid-level features (stride 16) with higher resolution, and may benefit from Transformer processing. Therefore, we extend the encoder to a multi-level feature integration mode: introducing Transformer encoders for multiple scale features separately and fusing their outputs again. Specifically, during HybridEncoder initialization, we can set a feature layer index list use\_encoder\_idx (such as including mid-level index 1 and high-level index 2), and the model will construct a separate Transformer Encoder module for each specified layer. During forward propagation, for each feature layer included, we execute its encoder, flatten 2D features into sequences, add positional encoding, send them to the encoder for self-attention and MoE-FFN transformation, then reshape results back to original feature map shape. Multi-level features enhanced by Transformer then enter the ScaleBlock multi-scale fusion module for interactive fusion. Under this design, not only do the highest-level features obtain global relationship modeling, but mid-level features can also benefit from Transformer processing, while absorbing information from other layers during fusion, further improving small target detection effects.

It should be noted that introducing multi-level Transformers brings certain computational cost increases, but we can control total costs by reducing the depth of each layer's Transformer (such as 1 encoder layer each). Additionally, RT-DETR's decoder itself supports dynamic layer number adjustment for speed control, so our model can still flexibly balance speed and accuracy during deployment.

In summary, our improved RT-DETR model fuses the advantages of convolution and Transformer in the encoder part: convolution provides local perception and enhances multi-scale representation through ASFF, ECA, etc., while Transformer introduces global dependencies and gated expert mechanisms to enhance modeling capability. The decoder part continues RT-DETR's design, using multi-layer multi-head attention to iteratively optimize queries and output detection results, with each layer having auxiliary detection heads for training. The model's training loss includes detection loss (classification, bounding box regression) and auxiliary balance loss for MoE routing, with total objective function  $L = L_{det} + \lambda L_{balance}$ , where  $\lambda$  is the weight. Through the above improvements, we expect the model to more accurately detect fire and smoke targets in drone imagery, with specific performance improvements to be verified in experiments in the next section.

#### **4 EXPERIMENTAL DESIGN AND EVALUATION**

#### 4.1 Experimental Setup

**Dataset:** We evaluate our model using a self-collected and annotated UAV smoke fire dataset. This dataset contains wildfire flame and smoke images from various scenarios, totaling approximately 12,551 images. 70% are used for training, 15% for validation, and 15% for testing. The images are extracted from UAV aerial video frames with 1080p resolution, covering environments such as forests, grasslands, and mountainous areas, with fire conditions ranging from initial smoke to large-scale open flames. Annotations follow the COCO format, with each flame or smoke target marked by bounding boxes and categorized into two classes (fire or smoke). During training, we treat both classes as positive samples for detection (without distinguishing categories for mAP evaluation), while calculating individual class AP separately during evaluation for reference. Prior to model training, images undergo data augmentation including random scaling, cropping, and color jittering to improve the model's adaptability to fire conditions of different scales.

**Training Details:** We train all models under the PyTorch framework using the AdamW optimizer with an initial learning rate set to 1e-4. We first perform 2000 steps of linear warmup, followed by a linear decay strategy consistent with DINO to gradually reduce the learning rate from the initial value to the minimum value. Due to the relatively small dataset size, training employs pre-trained weight initialization: the ResNet50 backbone loads ImageNet pre-trained parameters, while the Transformer encoder components use Xavier random initialization. The Mixture of Experts (MoE) parameters are initialized with uniform distribution, and router biases are appropriately adjusted to encourage balance. Training is conducted for 70 epochs with a batch size of 64 (distributed data parallel training on two NVIDIA V100

GPUs). For the loss function, the detection branch uses Focal Loss (classification) and CIoU loss (bounding box), along with denoising training techniques from DN-DETR to stabilize convergence. The MoE routing balance loss coefficient  $\lambda$  is set to 0.01, which has been experimentally verified to achieve good results. During training, we observed that the auxiliary branch loss stabilizes after approximately 40 epochs, with overall convergence reaching optimal performance at epochs 50-60.

**Evaluation Metrics:** We adopt the standard COCO object detection metrics, specifically Average Precision (AP). The report primarily focuses on: mAP (mean AP) under IoU threshold 0.5:0.95 and  $mAP_{50}$  under IoU=0.5. Additionally, to more intuitively reflect detection performance, we provide Precision and Recall metrics (using IoU=0.5 to determine true positives). Inference speed is measured by frames per second (FPS) on a single NVIDIA V100 GPU with batch size=1, tested at 640×640 scaled resolution. Model parameters (Million) and computational complexity (GFLOPs) are also provided as references. For dual-category (fire and smoke) detection, we calculate AP for each class but primarily evaluate overall model capability using comprehensive AP. All experiments are run multiple times and averaged to reduce random fluctuations.

**Comparison Methods:** We select several mainstream object detection models as baselines: (1) Two-stage representative: Faster R-CNN (ResNet50); (2) Single-stage representatives: YOLOv7-min (official version) and its standard version YOLOv7; (3) Transformer representative: original RT-DETR (Res18), as well as our implemented versions with various improvement modules removed for ablation studies. All aforementioned models are fine-tuned on the same dataset with identical training configurations to ensure fair comparison.

# 4.2 Overall Performance Comparison

Table 1 presents the performance comparison between our proposed model and mainstream detection models on the UAV smoke fire dataset test set. The results demonstrate that our improved RT-DETR achieves optimal performance across all metrics. Specifically, under the IoU=0.5 standard, our model achieves an  $mAP_{50}$  of 88.8%, representing approximately a 2 percentage point improvement over the original RT-DETR and surpassing the YOLOv7-min model by about 5 percentage points. For detection recall, our model achieves 87.9%, showing significant improvement compared to the original RT-DETR's approximately 86.7%. This indicates that our model reduces missed detections while not introducing additional false positives. The two-stage Faster R-CNN performs worst due to its insensitivity to small targets, achieving only about 80% mAP with recall below 80%, making it difficult to meet practical requirements.

Table I Detection Performance of Models on the UAV Smoke-Fire Dataset							
Item	#Epochs	#Params (M)	GFLOPs	${\rm FPS}_{\rm bs=1}$	$mAP_{50}$	$mAP_{0.5:0.95}$	Recall
Faster R-CNN	70	41.30M	134.38	21.27	0.804	0.507	0.792
YOLOv7-min	70	6.0M	6.5	171.0	0.832	0.575	0.828
YOLOv7	70	36.5M	51.6	62.7	0.894	0.643	0.861
RT-DETR	70	21.9M	29.7	86.9	0.868	0.612	0.867
Improved-RT-DETR	70	27.4M	37.1	71.5	0.888	0.638	0.879

Note:  $mAP_{50}/mAP_{0.5:0.95}$  at IoU 0.50/0.50–0.95

In terms of speed, our model achieves approximately 71.5 FPS for single-frame inference on NVIDIA V100, far exceeding real-time requirements (30 FPS), though slightly lower than the original RT-DETR. This is mainly due to the introduction of additional convolutional fusion and expert parameters, which increase computational overhead. However, our model's speed remains significantly higher than the two-stage Faster R-CNN (only around 21 FPS). YOLOv7-min has the fastest inference speed, reaching 171 FPS, outperforming our model. This is because Transformer self-attention and MoE computations are more time-consuming on high-resolution feature maps. Compared to the standard YOLOv7, while it has higher accuracy than our model, it also increases corresponding parameters and computational speed are acceptable for model deployment on small UAVs in fire monitoring scenarios where accuracy is prioritized. If TensorRT acceleration is used for Transformer computations, there is further room for speed improvement.



Figure 4 AP Convergence Curves (IoU = 0.50/0.50–0.95) for RT-DETR Variants

In summary, our model comprehensively outperforms the baseline RT-DETR in accuracy, particularly in detecting small flames and distant smoke columns, which is also demonstrated in the case analysis figures discussed later. Achieving such performance improvements while maintaining near real-time speed proves the effectiveness of our proposed improvement strategies (multi-scale fusion, attention enhancement, MoE expansion).

### 4.3 Ablation Studies

To quantify each improvement module's contribution to model performance, we designed a series of ablation experiments, with training results summarized in Figure 4. We conduct comparative analysis by progressively removing modules:

The results show that by gradually adding these modules, the model's detection accuracy steadily improves. Among them, ASFF multi-scale fusion brings the largest gain: after removing ASFF and ECA channel attention, mAP drops from 63.8% to 60.7%, a decrease of 3 percentage points, indicating that without ASFF, the model struggles to fully utilize multi-scale features, significantly degrading small target detection performance. Although ECA's contribution is less significant than ASFF, it remains non-negligible. After removing the MoE expert layer, mAP decreases by approximately 1.9 percentage points. This demonstrates that the MoE mechanism indeed provides performance improvement, validating that expert routing can enhance the model's ability to characterize different fire patterns.

Notably, the original baseline model achieves only 61.2% mAP, significantly lower than the complete model's 63.8%. This indicates that various improvements work synergistically to create the final significant enhancement. Without any component, model performance degrades to varying degrees. Particularly, ASFF fusion is crucial for information integration in small targets and complex backgrounds, serving as the key factor for our model's breakthrough over the baseline.

To intuitively demonstrate each module's role, we further compare detection results under different configurations for typical scenarios. As shown in Figure 5: in an image containing multiple distant smoke columns and multiple nearby open fire, the original model misses some smoke columns and incompletely boxes the open fire; after adding ASFF and ECA, small flames are correctly localized, proving that multi-scale fusion effectively enhances small-scale target signals; with the addition of MoE, fire boxes become more compact and accurate, and smoke is detected be slightly cause multiple experts collaborate to enhance feature response in fire regions; finally, the complete model (with MoE) has almost no missed detections in complex areas like smoke column edges, and no false detection of clouds as smoke, indicating that MoE experts further improve the model's ability to distinguish different fire appearances.





(d) RT-DETR+ASFF+ECA+MoE

**Figure 5** Comparison of Detection Effects under Different Improvement Module Configurations. (a) Original RT-DETR, Missing Some Tiny Flames; (b) +ASFF+ECA, Some Tiny Flames are Detected; (c) +MoE, Fire Target Boxes are more Accurate, and smoke is Detected; (d) ASFF+ECA+MoE full Model, all Fire Targets are Correctly Detected

# **5 RESULES AND DISCUSSION**

# 5.1 Analysis of Model Improvement Effects

Based on the comprehensive experimental results above, we can clarify each improvement component's contribution to model performance enhancement:

**ASFF Multi-scale Fusion:** Greatly improves the model's detection capability for fire targets of different scales. Particularly in detecting distant small smoke columns, ASFF's upsampling fusion enables the model to utilize high-resolution features, significantly improving recall rate. Meanwhile, since ASFF adaptively selects feature sources at each spatial location, it reduces interference from irrelevant scale features, lowering false detection rate (Precision also improves). This is validated in both ablation experiments and visualizations. ASFF can be said to solve the insufficient cross-scale fusion problem of the original RT-DETR, and its importance aligns with conclusions from previous research on small object detection.

**ECA Attention Mechanism:** Helps the model better focus on discriminative features of flames and smoke. Through combined use of ECA and ASFF, the model can automatically increase channel weights for fire source highlight regions while suppressing background noise channels, playing a subtle but important role in improving detection accuracy. Although ECA's removal with ASFF only slightly decreases mAP by 1.9% in ablation studies, the localization accuracy improvement brought by ECA when used with other modules is visually apparent. This indicates ECA improves the signal-to-noise ratio of features, making the model's confidence judgment for targets more precise. Compared to SE modules, ECA requires no explicit dimensionality reduction and expansion operations, offering higher computational efficiency, making it very suitable for our real-time model.

**MoE Mixture of Experts:** Enhances the model's adaptability to diverse fire patterns. Since fire forms are highly varied, a single network struggles to handle all situations well, while MoE allows multiple experts to learn separately, for example, some experts focus on learning dense smoke scenarios while others specialize in open flame burning patterns. When actual input arrives, the routing network automatically selects appropriate expert combinations for processing. This mechanism effectively improves detection robustness in complex scenarios. Our model can detect partially occluded fire sources even in extreme cases (such as dense smoke obscuring open flames), which is nearly impossible with the original model. Although MoE's overall mAP improvement is less obvious than ASFF, in several difficult samples we tested, detection results with MoE enabled show significant improvement compared to when MoE is disabled. This indicates MoE's advantages mainly manifest in difficult cases—it provides the model with more capacity to characterize special situations, thereby improving the overall performance lower bound.

**Multi-layer Transformer Integration:** This paper primarily uses the highest-layer Transformer encoding. We attempted to simultaneously apply encoders to mid-layer features and fuse them, resulting in approximately 0.4 percentage point mAP improvement, but considering the computational cost increase of about 15%, we ultimately did not include it as a main result. However, this phenomenon merits discussion: multi-level encoding indeed further

improves performance, indicating Transformer also helps mid-layer features, but possibly due to high resolution of midlayer features causing time-consuming attention with limited gains. Under stronger hardware or more optimized implementations, this strategy can serve as an option for balancing accuracy. Our framework design already supports flexible selection of encoding feature layers, which can be enabled as needed in the future.

## 5.2 Comparison with YOLO Series Methods

Although our improved model is based on RT-DETR, it's necessary to compare and discuss it with current state-of-theart YOLO series methods. From Table 1, our model significantly outperforms YOLOv7-min in accuracy, particularly advantageous in recall rate, indicating Transformer's benefits in capturing global information and discovering hidden targets. YOLO, due to anchor mechanisms and receptive field limitations, may miss some inconspicuous smoke points. On the other hand, YOLO remains faster, mainly attributed to efficient pure CNN architecture implementation on GPUs. Therefore, in actual deployment, if pursuing ultimate speed while accepting certain missed detections, YOLOv7min/YOLOv8-min remain good choices. However, in accuracy-prioritized scenarios (such as wildfire early warning requiring extremely low false negatives), our model provides more confident detection results.

Notably, new models like YOLOv11 also incorporate Transformer concepts (such as Decoupled Head, self-attention modules), continuously improving performance. If real-time visual Transformers further optimize speed in the future, Transformer detectors have potential to comprehensively surpass YOLO series. This research also demonstrates that by introducing excellent YOLO modules like ASFF and attention mechanisms, Transformer models' shortcomings (multi-scale and local features) can be addressed, thereby leveraging Transformer's strength in modeling global dependencies. This provides insights for future detection model design combining CNN and Transformer advantages.

## **5.3Model Limitations and Improvement Directions**

Despite our model achieving good performance on our dataset, some limitations remain: (1) High model complexity with nearly 27.4 million parameters is oversized for some embedded platforms, hindering real-time deployment on UAV terminals. Future work could consider model pruning, distillation, or lighter backbones (such as MobileNet series) to reduce model size. (2) Our model currently only utilizes visible light image features, not yet addressing fire point detection in nighttime infrared imaging. Introducing multi-spectral data (infrared + visible light) for multi-modal fusion detection could significantly improve all-weather applicability. (3) MoE routing mechanism increases training instability; we observed that routing tends to favor certain experts in early training, requiring loss weight adjustment for convergence. Future work could explore more stable expert selection algorithms or introduce online hard example mining to make different experts' roles more distinct.

Additionally, due to our limited dataset scale, model potential may not be fully exploited. If larger-scale, more diverse UAV fire data could be collected and pre-training or semi-supervised learning employed, model performance could further improve. Some latest research directions such as video temporal information utilization, 3D convolution modeling of fire dynamics, and generative adversarial networks for synthesizing training samples are also worth trying to compensate for insufficient real data.

Overall, this research provides an effective solution for fire target detection model improvement. By combining multiscale fusion, attention enhancement, and expert routing, we significantly improved detection accuracy while maintaining real-time performance. Looking forward, applying these strategies to more scenarios (such as urban fire monitoring, industrial accident warning) and combining with other sensor information, intelligent fire detection systems will become more robust and reliable.

# 6 CONCLUSION AND OUTLOOK

This paper designs an improved RT-DETR-based detection model for UAV fire detection tasks and conducts systematic experimental research on a self-built dataset. We introduce ASFF multi-scale feature fusion modules, ECA efficient channel attention mechanisms, and gated MoE mixture of experts structures into the RT-DETR model's hybrid encoder, while adopting multi-layer Transformer feature integration strategies, significantly improving the model's detection performance for flame and smoke targets of different scales. Experimental results show that compared to original RT-DETR and classic methods like YOLO and Faster R-CNN, our model has obvious advantages in detection accuracy and recall rate, particularly more accurate and reliable identification of small fire targets. Under IoU=0.5 metrics, our model achieves 88.8% mAP, improving approximately 2 percentage points over baseline with significantly reduced missed detection rate. Through ablation experiments, we quantified each improvement component's contribution, with ASFF multi-scale fusion contributing most, while ECA attention and MoE expert mechanisms also provide positive gains. Although model parameters increase somewhat, inference speed remains near real-time, meeting most UAV inspection application requirements.

Research proves that combining multi-scale fusion, attention mechanisms, and MoE expert routing can effectively enhance Transformer detectors' performance in fire monitoring domains. This provides useful reference for future development of high-precision intelligent fire monitoring systems. Looking ahead, we will further improve from the following directions: (1) Explore model lightweighting techniques such as knowledge distillation and network pruning for deployment on computation-constrained UAV platforms; (2) Expand training data including nighttime infrared fire imagery and simulated data augmentation to improve model adaptability to various conditions; (3) Extend the model to

tasks like fire spread prediction by combining video temporal information and multi-modal sensor data, achieving functionality from "seeing fire" to "predicting fire development." In the near future, with continued development of deep learning and edge computing, we have reason to expect more intelligent and efficient aerial fire monitoring systems to play key roles in forest fire prevention.

### **COMPETING INTERESTS**

The authors have no relevant financial or non-financial interests to disclose.

# REFERENCES

- [1] Chen Y, Zhang Y, Xin J, et al. A UAV-based forest fire detection algorithm using convolutional neural network. 2018 37th Chinese Control Conference (CCC). IEEE, 2018: 10305-10310.
- [2] Haucap J, Rasch A, Stiebale J. How mergers affect innovation: theory and evidence. International Journal of Industrial Organization, 2019, 63: 283-325.
- [3] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollar. Focal loss for dense object detection. In Proceedings of the IEEE/CVF International Conference on Computer Vision, 2017, 2: 2980–2988.
- [4] Jocher G, Stoken A, Borovec J, et al. ultralytics/yolov5: v3. 0. Zenodo, 2020.
- [5] Chien-Yao Wang, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. Yolov7: Trainable bag-of-freebies sets new state-ofthe-art for real-time object detectors. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2023: 7464–7475.
- [6] Mukhiddinov M, Abdusalomov A B, Cho J. A wildfire smoke detection system using unmanned aerial vehicle images based on the optimized YOLOv5. Sensors, 2022, 22(23): 9384.
- [7] Zhao Y, Lv W, Xu S, et al. Detrs beat yolos on real-time object detection. Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2024: 16965-16974.
- [8] Xizhou Zhu, Weijie Su, Lewei Lu, et al. Deformable detr: Deformable transformers for end-to-end object detection. In International Conference on Learning Representations, 2020.
- [9] Shilong Liu, Feng Li, Hao Zhang, et al. Dab-detr: Dynamic anchor boxes are better queries for detr. In International Conference on Learning Representations, 2021.
- [10] Lv W, Zhao Y, Chang Q, et al. Rt-detrv2: Improved baseline with bag-of-freebies for real-time detection transformer. arXiv preprint arXiv:2407.17140, 2024.
- [11] Liu Z, Zhang K, Wang C, et al. Research on the identification method for the forest fire based on deep learning. Optik, 2020, 223: 165491.
- [12] Jiaqi Shi, Jinhu Wang, Junhui Xu, et al. Research on forest fire monitoring technology based on UAV and convolutional neural network. Advances in Applied Mathematics, 2022, 11: 3200.
- [13] Jie Li, Xuanbing Qiu, Enhua Zhang, et al. Fire recognition algorithm based on convolutional neural network. Journal of Computer Applications, 2020, 40(S2): 173-177.
- [14] Qiang Chen, Jian Wang, Chuchu Han, et al. Group detr v2: Strong object detector with encoder-decoder pretraining. arXiv preprint arXiv:2211.03594, 2022.
- [15] Shazeer N, Mirhoseini A, Maziarz K, et al. Outrageously large neural networks: The sparsely-gated mixture-ofexperts layer. arXiv preprint arXiv:1701.06538, 2017.
- [16] Fedus W, Zoph B, Shazeer N. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. Journal of Machine Learning Research, 2022, 23(120): 1-39.
- [17] Riquelme C, Puigcerver J, Mustafa B, et al. Scaling vision with sparse mixture of experts. Advances in Neural Information Processing Systems, 2021, 34: 8583-8595.
- [18] Yuan, Jinghuil. A Margin-Maximizing Fine-Grained Ensemble Method. arXiv preprint arXiv:2409.12849, 2024.