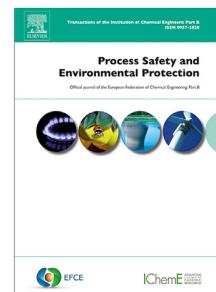


Journal Pre-proof

Masked autoencoder-based vision framework for robust fire detection in complex environments

Hanxiang Wang, Muhammad Fayaz, Awais Ahmad, Yanfen Li, Tan N. Nguyen, L. Minh Dang



PII: S0957-5820(25)01286-8

DOI: <https://doi.org/10.1016/j.psep.2025.108019>

Reference: PSEP 108019

To appear in: *Process Safety and Environmental Protection*

Received date: 13 July 2025

Revised date: 8 October 2025

Accepted date: 13 October 2025

Please cite this article as: H. Wang, M. Fayaz, A. Ahmad et al., Masked autoencoder-based vision framework for robust fire detection in complex environments. *Process Safety and Environmental Protection* (2025), doi: <https://doi.org/10.1016/j.psep.2025.108019>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2025 Published by Elsevier Ltd on behalf of Institution of Chemical Engineers.

1 Masked Autoencoder-Based Vision Framework for Robust Fire Detection in Complex 2 Environments

10 Abstract

11 Vision-based fire detection has become an increasingly important focus in computer vision, driven by the growing need for early
12 warning systems and public safety in surveillance environments. While conventional models have primarily relied on color-based
13 features to distinguish fire from background, maintaining high detection accuracy while ensuring computational efficiency remains
14 a persistent challenge, particularly in real-time surveillance systems. To address this, we introduce a novel fire detection framework
15 grounded in masked autoencoding and Vision Transformers (ViT), designed to balance detection performance with scalable deploy-
16 ment. Our architecture leverages self-supervised learning to reconstruct masked visual regions, enhancing the encoder's ability to
17 capture fine-grained fire cues in complex scenarios. The integration of global attention and hierarchical context modeling enables
18 the system to distinguish between fire and visually similar non-fire patterns, such as reflections and artificial lighting, under diverse
19 environmental conditions. Unlike prior models that are sensitive to background noise or rely heavily on channel saliency, our ap-
20 proach learns robust representations through reconstruction objectives, eliminating the need for hand-crafted modules. Extensive
21 experiments conducted on five benchmark datasets: BWF, DQFF, LSFD, DSFD, FG and DFAN demonstrate consistent improve-
22 ments over existing methods, with notable gains of 2.5% on BWF, 2.2% on DQFF, 1.42% on LSFD, 1.8% on DSFD, 1.14% on FG
23 and 1.10% on DFAN. The proposed model also maintains computational efficiency and generalizes effectively across a wide range
24 of fire conditions, supporting its deployment in practical, real-time systems.

25 **Keywords:** Fire Detection, Masked Autoencoder, Vision Transformer, Self-Supervised Learning, Feature Reconstruction,
26 Attention Mechanism.

31 1. Introduction

32 Fires represent a severe threat to human life and property
33 due to their rapid and often uncontrollable spread, especially
34 in densely populated regions such as urban residential areas,
35 transportation hubs, and **forested environments** [1]. Ensuring
36 prompt detection of fire outbreaks is crucial to mitigate damage
37 and improve public safety in various domains, including resi-
38 dential, commercial, and industrial settings.

39 Conventional fire detection (FD) systems typically employ
40 environmental sensors like smoke, temperature, and **particle**
41 **detectors** [2]. These sensors are cost-effective and relatively
42 easy to deploy, particularly in confined indoor environments.
43 However, their effectiveness diminishes significantly in open or
44 large-scale outdoor scenarios. Moreover, these systems often
45 activate only when they directly detect fire by-products such as
46 heat or smoke, potentially delaying the response time and re-
47 ducing the chances of **early intervention** [3, 4].

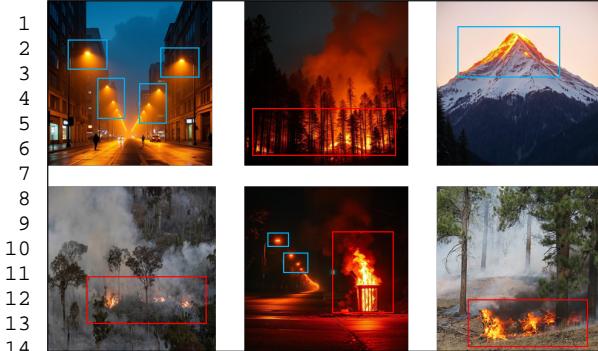
48 To overcome these limitations, there has been a growing in-
49 terest in vision-based FD technologies which utilize cameras as
50 sensors to monitor large areas in **real-time** [5, 6]. These sys-
51 tems offer broader coverage, quicker response, and better adapt-
52 ability to **various environmental conditions** [7]. Consequently,
53 numerous vision-based methods have been proposed, mainly
54 categorized as traditional machine learning (TML) and deep
55 learning (DL) approaches [8].

56 TML-based approaches are based on traditional feature ex-
57 traction techniques, such as flame texture, color, and **motion**

58 **patterns** [6]. The success of these methods depends heavily
59 on the quality and relevance of the features designed manually.
60 However, it is challenging to design a robust global feature ex-
61 traction due to the diverse and dynamic nature of fire. Variations
62 in flame color caused by different combustible materials,
63 lighting conditions, and environmental influences such as wind
64 or temperature fluctuations contribute to the unpredictability of
65 flame behavior. These factors can hinder consistent FD and in-
crease the likelihood of false positives. To effectively harness
TML-based methods, it is essential to master the challenge of
attaining a high true positive rate while minimizing the false
alarm rate.

66 In response to these challenges, deep learning has emerged as
67 a powerful alternative, offering end-to-end learning capabilities
68 and superior performance across a variety of computer vision
69 tasks [9, 10]. DL-based models automatically learn discrimi-
70 native features from large datasets, enabling them to generalize
71 effectively to new and unseen fire scenarios. This capacity to
72 extract complex patterns has led to substantial improvements
73 in detection accuracy and robustness. Notably, several DL-
74 based techniques have demonstrated enhanced reliability over
75 TML approaches, especially under variable and uncertain **con-**
76 **ditions** [11].

77 Despite these advancements, DL-based FD systems are not
78 without limitations. Their performance can degrade in visually
79 complex scenes, such as when fire-colored objects are present
80 or when the fire source is distant and small in the frame, as il-



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Figure 1: Sample images illustrating challenging scenarios for FD. The red bounding boxes indicate actual fire regions, while the blue bounding boxes highlight visually similar non-fire regions such as lighting artifacts, reflections, or sunlight patches. These examples demonstrate the difficulty in distinguishing real fires from fire-like patterns, particularly under varying environmental and lighting conditions.

1.1. Research Gap

Despite notable progress in FD technologies, significant challenges remain unresolved, particularly in achieving timely and reliable detection in complex and dynamic environments. Traditional sensor-based systems, while cost-effective and simple to deploy, are constrained by their limited range and delayed responsiveness, especially in outdoor or large-scale settings. These systems often rely on direct detection of smoke or heat, which restricts their ability to identify fires in their early stages. On the other hand, vision-based approaches, although more promising in terms of coverage and responsiveness, also encounter limitations.

Traditional Machine Learning (TML) techniques depend heavily on manually crafted features such as flame color, texture, and motion, which are highly susceptible to environmental variations including lighting conditions, background clutter, and fire appearance diversity. This dependence leads to inconsistent detection performance and a high rate of false alarms.

While Deep Learning (DL)-based models have demonstrated superior accuracy and generalization in FD tasks, their effectiveness declines in certain complex scenarios. Specifically, DL models struggle with false detections when visually similar objects (e.g., fire-colored materials) are present, or when fire is located at long distances, making it less visible in the frame. Additionally, many existing DL methods have not been rigorously evaluated across a wide range of real-world conditions, limiting their practical applicability.

Therefore, there is a clear need for an advanced FD framework that combines the strengths of deep learning with robust feature representation, capable of handling real-time detection in diverse and challenging environments while minimizing

false positives and improving detection accuracy across varying scales and contexts.

1.2. Main Contributions

To address the challenges and research gaps in vision-based FD, this study introduces a novel FD framework based on masked autoencoding and transformer architectures. The proposed method aims to improve early FD accuracy and robustness in complex environments by leveraging advanced feature learning and contextual reasoning. The main contributions of this work are as follows:

- We propose an Image Masked Autoencoder (ImageMAE) based FD framework that efficiently learns rich and discriminative fire-related features through a self-supervised reconstruction task. The framework uses an asymmetric encoder-decoder design, where the encoder processes only visible image patches, significantly reducing computational overhead.
- A Vision Transformer (ViT)-based encoder is utilized for feature extraction, capturing long-range dependencies and complex fire patterns such as varying flame shapes, colors, and textures. This enhances the model's ability to distinguish between real fires and fire-like objects in challenging visual conditions.
- The reconstruction module incorporates a novel pixel-level reconstruction loss with optional normalized pixel targets, improving feature invariance to environmental factors such as lighting changes and smoke occlusion. This leads to more robust representations suitable for early-stage FD.
- We conduct comprehensive experiments on a diverse fire dataset, demonstrating that the proposed method achieves superior FD accuracy compared to baseline approaches. Additionally, the model's lightweight decoder and efficient masking strategy enable scalability and real-time applicability in surveillance systems.

1.3. Study Outline

The remainder of this paper is organized as follows: Section II reviews recent literature on FD methods, highlighting existing challenges and advances. Section III details the proposed ImageMAE-based FD framework, including architecture and training strategies. Section IV presents the datasets used, experimental results, and a comprehensive analysis of the model's performance. The conclusion and future direction are presented in Section V.

2. Related Work

The development of early and accurate FD technology using vision sensors has been an active research area and can be broadly categorized into Traditional Machine Learning (TML) methods and Deep Learning (DL) methods. This section provides a detailed overview of both approaches.

1 2.1. Traditional Machine Learning Methods

2 Traditional machine learning (TML) approaches primarily
 3 focus on handcrafted features based on fire characteristics such
 4 as color, shape, motion, and texture [12]. Early color-based
 5 FD methods were proposed by [13, 14]; however, these meth-
 6 ods suffered from high false positive rates and limited accu-
 7 racy, restricting their practical application. [15] introduced an
 8 auto-adaptive edge detection technique to identify fire regions,
 9 while [16] applied methods like K-means clustering, optical
 10 flow, logistic regression, and temporal smoothing for FD. [17]
 11 employed multiple classification algorithms and combined their
 12 outputs for accurate real-time fire scene classification. [18] used
 13 covariance features with Support Vector Machines (SVM) for
 14 fire scene classification. Teng et al. [19] applied hidden Markov
 15 models to detect moving fire pixels extracted from pixel clus-
 16 ters. [20] developed BoWFire, a novel approach that fused
 17 color features with super-pixel texture discrimination to im-
 18 prove FD. While traditional machine learning methods have
 19 contributed to fire detection, they rely heavily on handcrafted
 20 features and standard classifiers, which are error-prone and
 21 labor-intensive. They often confuse fire with fire-like objects
 22 and perform poorly under varying lighting, smoke, or occlu-
 23 sion. These challenges highlight the need for approaches that
 24 can automatically learn robust features. Our ImageMAE-ViT
 25 framework addresses this by combining self-supervised learn-
 26 ing with global context modeling, enabling accurate and effi-
 27 cient fire detection across diverse conditions.

28 2.2. Deep Learning Methods

29 Deep learning (DL) techniques automatically learn features
 30 from raw image data and have demonstrated superior perfor-
 31 mance over traditional methods in FD. Sharma et al. (2017) in-
 32 troduced a challenging fire dataset and evaluated models such
 33 as ResNet50 [21] and VGG16 [22], with ResNet50 achieving
 34 promising results. [23] designed a convolutional neural net-
 35 work (CNN) for fire scene classification. [24] combined lo-
 36 cal binary patterns with AdaBoost to isolate fire regions, which
 37 were then input to a CNN for feature extraction and classifica-
 38 tion. [25] utilized pretrained GoogLeNet weights to balance
 39 efficiency and accuracy for FD. In another study, they fine-
 40 tuned AlexNet for indoor and outdoor fire scenarios, employ-
 41 ing adaptive camera prioritization strategies. [26] used a mod-
 42 ified InceptionV1 architecture for efficient FD. [27] compared
 43 various classifiers, including MLP, AdaBoost, AdaBoost-LBP,
 44 and CNN for FD. Zhang et al. [28] proposed a novel CNN
 45 model with three convolutional and three fully connected layers
 46 aimed at efficient FD. More recently, [29] modified VGG16 to
 47 reduce model size and training parameters while significantly
 48 improving accuracy. [5] presented a CNN architecture con-
 49 nected with an autoencoder for fire scene classification. Altaf
 50 et al. [1] proposed an optimized deep learning model with an
 51 attention module to improve detection efficiency, while Khan et
 52 al. [30] developed a multi-attention network with a new bench-
 53 mark dataset for real-time fire detection. Wang et al. [31] intro-
 54 duced FFD-YOLO, a modified YOLOv8 architecture tailored
 55 for forest fire scenarios, and Lv et al. [32] enhanced YOLOv8

56 with CARAFE and context-guided modules for robust perfor-
 57 mance. Complementing these works, Wang et al. [2] presented
 58 a multi-source data fusion framework using deep learning to
 59 improve detection accuracy across diverse conditions. Beyond
 60 deep learning, Maity et al. [6] proposed MLSFDD, a smart
 61 fire detection device for precision agriculture, demonstrating
 62 the growing interest in domain-specific fire detection applica-
 63 tions.

64 However, current deep learning models for FD often rely
 65 on relatively plain architectures, limiting their ability to ex-
 66 tract fine-grained details and accurately localize fire regions.
 67 To overcome these shortcomings, many researchers have in-
 68 corporated attention mechanisms into their models, utilizing
 69 various backbones such as attention-squeeze networks, deep
 70 CNNs with channel attention [33], spatial and channel atten-
 71 tion modules [34], InceptionV3 with CBAM [35], Vision Trans-
 72 former (ViT) with self-attention [36], non-local attention net-
 73 works [37], EfficientNetB0 with attention [38], and other self-
 74 attention frameworks [39, 40]. While these attention-based
 75 methods generally outperform plain CNN architectures by en-
 76 hancing the focus on salient features, they often process feature
 77 maps within a single receptive field, which may not fully cap-
 78 ture the complex spatial and contextual variations required for
 79 accurate FD.

80 Existing attention-based fire detection models employ a va-
 81 riety of mechanisms to improve feature learning. Common
 82 approaches include channel attention, which emphasizes im-
 83 portant feature maps, and spatial attention, which highlights
 84 relevant image regions. Beyond these, more advanced strate-
 85 gies such as combined channel-spatial attention modules (e.g.,
 86 CBAM), non-local attention networks, and transformer-based
 87 self-attention have been applied to capture long-range depen-
 88 dencies and contextual information. Despite their effective-
 89 ness, these methods still face challenges in complex scenar-
 90 ios with fire-like objects, reflections, smoke, or varied back-
 91 grounds, often resulting in higher false-positive rates. Many
 92 also struggle to robustly detect fires across varying scales, dis-
 93 tances, and environmental conditions such as occlusions and
 94 lighting changes, partly due to limited multiscale feature mod-
 95 eling.

96 To address the limitations of current deep learning models for
 97 fire detection, our method leverages a masked autoencoder [41]
 98 (ImageMAE) framework that inherently captures multi-scale
 99 and contextual information through self-supervised reconstruc-
 100 tion of masked image patches. Unlike conventional attention
 101 modules, which often operate on single-scale feature maps and
 102 struggle with complex environmental variations, our approach
 103 learns rich hierarchical representations by reconstructing miss-
 104 ing regions based on global context, enabling the encoder to fo-
 105 cus on discriminative fire-related patterns across multiple spa-
 106 tial scales. Current DL methods can be broadly classified into
 107 three categories: (i) plain CNN architectures, which achieve
 108 reasonable accuracy but often fail to capture fine-grained spa-
 109 tial and contextual details; (ii) attention-enhanced CNNs, which
 110 improve focus on salient features but still have limited multi-
 111 scale awareness; and (iii) transformer-based methods, which
 112 model long-range dependencies but can be computationally ex-
 113 pensive.

1 pensive and often struggle to distinguish fire from fire-like objects under challenging conditions. By integrating masked autoencoding with Vision Transformer-based feature extraction, 2 our framework addresses these challenges, reducing false positives, 3 improving robustness to lighting changes, occlusions, smoke, and varying scales, and enabling accurate, efficient, and 4 generalizable fire detection across diverse real-world scenarios.

5 3. Proposed Methodology

6 In this work, we propose a vision-based FD framework 7 grounded on the Masked Autoencoder (MAE) architecture [42], 8 originally introduced for general image representation learning. Our adaptation takes advantage of the strengths of MAE for 9 extracting discriminative and robust features from surveillance 10 imagery, which often includes challenges such as variable lighting, 11 occlusions caused by smoke or objects, and visually confusing 12 fire-like patterns. By reconstructing missing parts of the 13 input images, the model inherently learns the global context and 14 fine-grained fire-specific visual cues necessary for early and accurate 15 FD. The general flow of the proposed method is illustrated in Figure 2.

16 3.1. Overall Architecture and Motivation

17 The core idea behind our approach is to employ a self-supervised 18 learning scheme in which the model learns meaningful representations 19 by reconstructing masked portions of the input image. This paradigm 20 forces the network to reason about the overall scene and infer the 21 presence of fire even when only partial information is visible. Given 22 that early-stage fires are often small and visually subtle, such contextual 23 reasoning is critical. Our MAE-based architecture is composed of two main 24 components: an encoder that transforms visible parts of the input 25 into latent representations and a decoder that attempts to reconstruct 26 the original image from these latent codes supplemented by learned mask 27 tokens.

28 The asymmetric nature of the design, a relatively large encoder paired 29 with a lightweight decoder, allows us to efficiently train deep models with 30 reduced computational requirements while preserving or enhancing 31 performance. This is especially important for FD systems deployed on edge 32 devices or real-time monitoring setups, where computational resources and 33 latency constraints are significant factors.

34 3.1.1. Patch Tokenization and Masking Strategy

35 Given an input image $\mathbf{I} \in \mathbb{R}^{H \times W \times C}$, where H and W denote 36 spatial dimensions and C the number of color channels (usually 37 3 for RGB), we divide the image into fixed-size, nonoverlapping patches 38 of dimension $P \times P$ pixels. This process converts the input into a sequence 39 of discrete tokens, formalized as:

$$40 N = \frac{H \times W}{P^2} \quad (1)$$

41 where N is the total number of patches. Each patch $\mathbf{p}_i \in \mathbb{R}^{P^2 \times C}$ 42 is flattened and projected into a latent embedding vector $\mathbf{z}_i \in \mathbb{R}^D$ via a learnable linear transformation:

$$43 \mathbf{z}_i = \mathbf{W}_e \cdot \text{Flatten}(\mathbf{p}_i) + \mathbf{e}_i \quad (2)$$

44 Here, $\mathbf{W}_e \in \mathbb{R}^{D \times (P^2 \times C)}$ is the projection matrix, and $\mathbf{e}_i \in \mathbb{R}^D$ is the positional embedding that encodes spatial location information, vital for preserving the structure of the scene.

45 To promote the model's ability to infer global context and avoid relying on local neighboring information, we adopt a random masking procedure. A large fraction (commonly 75%) of the patch tokens is randomly removed from the input sequence. This is implemented by uniform random sampling without replacement, yielding a visible subset \mathcal{V} and a masked subset \mathcal{M} :

$$46 \mathcal{V}, \mathcal{M} = \text{RandomMasking}(\{\mathbf{z}_1, \dots, \mathbf{z}_N\}) \quad (3)$$

47 The removal of such a significant portion of the input forces the encoder to rely on partial observations to learn holistic and discriminative features, a property that is particularly advantageous for identifying subtle fire cues obscured by environmental factors such as smoke or lighting variation.

48 3.1.2. Encoder Design: Feature Extraction via Vision Transformer

49 The encoder is based on the Vision Transformer (ViT) architecture [43], modified to operate exclusively on the visible patch embeddings \mathcal{V} . Unlike traditional ViT models which process all tokens including special mask tokens, our encoder discards masked tokens entirely, thus saving computational and memory resources. This approach enables scaling to deeper architectures without prohibitive cost.

50 The ViT encoder consists of multiple stacked Transformer blocks, each containing a Multi-Head Self-Attention (MHSA) mechanism and a Feedforward Network (FFN). The self-attention mechanism enables the model to dynamically focus on different parts of the visible image, capturing both local and long-range dependencies crucial for recognizing complex fire patterns.

51 The self-attention operation for a single head is mathematically represented as:

$$52 \text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right)\mathbf{V} \quad (4)$$

53 In this mechanism, the input embeddings are transformed into three distinct matrices: queries (\mathbf{Q}), keys (\mathbf{K}), and values (\mathbf{V}). Attention computation involves comparing queries with keys to produce attention scores, which are then used to weight the corresponding values. A scaling factor, typically the inverse square root of the key dimension ($\sqrt{d_k}$), is applied to stabilize the output of the dot product. To enhance the model's ability to learn a variety of contextual features, multiple attention operations - known as heads - are performed in parallel. The results of these heads are concatenated to form a richer, more expressive representation.

54 The encoder outputs a sequence of latent representations \mathbf{Z}_{enc} corresponding to the visible patches:

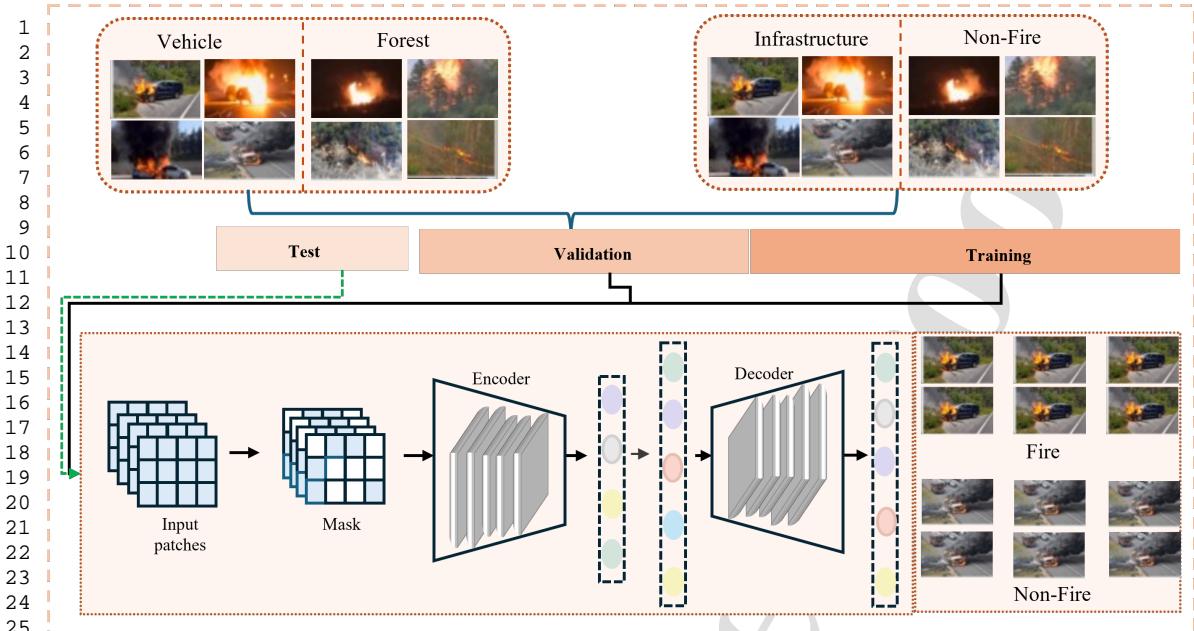


Figure 2: Fire recognition framework using ImageMAE.

$$\mathbf{Z}_{\text{enc}} = \text{Encoder}(\mathcal{V}) \quad (5)$$

These latent codes encode rich semantic and structural information about the scene, including the presence, shape, color, and dynamics of fire regions.

3.1.3. Decoder Design: Reconstructing Fire-Related Visual Patterns

The decoder is designed to reconstruct the original input image by predicting the pixel content of masked patches, leveraging the encoded visible tokens \mathbf{Z}_{enc} and a set of learned mask tokens $\mathbf{M} \in \mathbb{R}^{|\mathcal{M}| \times D}$. These mask tokens serve as placeholders for missing patches and are shared across all masked positions [44].

The combined sequence \mathbf{Z}_{full} is formed by concatenating encoded visible tokens with mask tokens, and then unshuffling to restore the original spatial order of patches:

$$\mathbf{Z}_{\text{full}} = \text{Unshuffle}(\mathbf{Z}_{\text{enc}} \cup \mathbf{M}) \quad (6)$$

The decoder consists of a smaller stack of Transformer blocks compared to the encoder, striking a balance between reconstruction capability and computational efficiency. The asymmetry of the architecture, a large encoder with a lightweight decoder, allows efficient pre-training without sacrificing representational power.

3.1.4. Reconstruction Loss and Optimization Objective

The primary training aim is to reduce the reconstruction error between the true pixels and predicted values of the masked patches. After processing by the decoder, each patch prediction $\hat{\mathbf{p}}_i$ is projected back to pixel space and reshaped to $P \times P \times C$. The mean squared error (MSE) loss is computed only over the masked patches:

$$\mathcal{L}_{\text{rec}} = \frac{1}{|\mathcal{M}|} \sum_{i \in \mathcal{M}} \|\hat{\mathbf{p}}_i - \mathbf{p}_i\|_2^2 \quad (7)$$

To further enhance robustness to variations in illumination and smoke density factors that can drastically change pixel intensities we explore a normalized reconstruction loss where pixel values of each patch are normalized by their mean and standard deviation:

$$\mathbf{p}_i^{\text{norm}} = \frac{\mathbf{p}_i - \mu_i}{\sigma_i}, \quad \mu_i = \text{mean}(\mathbf{p}_i), \quad \sigma_i = \text{std}(\mathbf{p}_i) \quad (8)$$

This normalised target increases the model's capacity to learn features that are unaffected by illumination variations, which is an important attribute in FD circumstances.

3.1.5. Fire Classification via Fine-tuning

Upon completion of the self-supervised pretraining phase, the decoder is discarded and the encoder is retained as a powerful feature extractor. For FD, we attach a classification head,

1 typically a Multi-Layer Perceptron (MLP), on top of the en-
 2 coder output corresponding to a special classification token or
 3 pooled features:

$$4 \\ 5 \\ 6 \quad \hat{y} = \text{Softmax}(\text{MLP}(\mathbf{Z}_{\text{enc}}^{[\text{CLS}]})) \quad (9)$$

7 The model is fine-tuned on labeled fire datasets by minimiz-
 8 ing the cross-entropy loss:

$$9 \\ 10 \\ 11 \\ 12 \quad \mathcal{L}_{\text{cls}} = - \sum_{c=1}^C y_c \log(\hat{y}_c) \quad (10)$$

13 where C is the number of classes (fire vs. no fire), y_c is
 14 the ground truth label, and \hat{y}_c is the predicted class probabili-
 15 ty. Fine-tuning adapts the pretrained features to the specific
 16 task of FD, enabling robust and accurate classification across
 17 diverse scenes.

21 3.1.6. Implementation and Computational Efficiency

22 Our implementation emphasizes efficiency and simplicity.
 23 Masking is performed by randomly shuffling patch tokens and
 24 removing a proportion based on the masking ratio, effectively
 25 simulating random patch sampling without replacement. After
 26 encoding, mask tokens are appended and the sequence is un-
 27 shuffled to maintain spatial consistency, ensuring the decoder
 28 can reconstruct masked patches in their original positions. Not-
 29 ably, this design does not require any specialized sparse tensor
 30 operations, making it practical for real-world deployment.

31 3.1.7. Benefits for Fire Detection

32 This ImageMAE-based approach offers several advantages
 33 for vision-based FD systems:

- 34 • **Robust Feature Learning:** Self-supervised reconstruc-
 35 tion encourages learning of generalized features that cap-
 36 ture essential fire characteristics under challenging con-
 37 ditions such as smoke, lighting variations, and complex
 38 backgrounds.
- 39 • **Computational Efficiency:** The asymmetric encoder-
 40 decoder design reduces training and inference costs, fa-
 41 cilitating real-time deployment.
- 42 • **Early Fire Recognition:** Contextual reasoning enabled by
 43 masking allows detection of fires in early, subtle stages
 44 when visual cues are partial or obscured.
- 45 • **Scalability:** The method is scalable to high-resolution im-
 46 ages and can be integrated with existing surveillance in-
 47 frastructure.

48 In summary, our ImageMAE-based FD method combines ef-
 49 ficient masked autoencoding with powerful transformer-based
 50 feature extraction to deliver accurate, timely, and scalable FD
 51 in real-world environments.

52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

3.2. Architecture Design

The proposed FD framework is constructed on a masked au-
 toencoder backbone and follows a carefully designed three-
 stage pipeline: (i) input preprocessing and patch embedding,
 (ii) masked autoencoder pretraining, and (iii) supervised fine-
 tuning for fire classification. Initially, images from surveillance
 streams are standardized in size and pixel values. Each image is
 partitioned into fixed-size non-overlapping patches, which are
 then flattened and projected into patch embeddings through a
 learnable linear mapping. Positional embeddings are added to
 preserve the spatial structure, ensuring that the model retains
 an understanding of the image layout. This transformation con-
 verts the 2D visual data into a 1D sequence of tokens, mak-
 ing it compatible with transformer-based processing. The cen-
 tral component of the system employs the masked autoencoder
 (MAE) strategy in a self-supervised manner. In this stage, a
 large proportion of image patches (around 75%) are randomly
 masked and excluded from the encoder input. This masking
 encourages the encoder to infer discriminative features from
 incomplete visual cues rather than relying on low-level pixel
 continuity. A Vision Transformer (ViT) serves as the encoder,
 processing the remaining visible patches and producing latent
 feature representations. These representations capture critical
 fire-related characteristics, such as the irregularity of flames,
 variations in color intensity, and the texture of smoke.

Subsequently, the decoder is provided with both the encoder
 outputs and a set of learned tokens representing the missing
 patches. Its role is to reconstruct the masked regions by pre-
 dicting the pixel values of the absent patches, guided by the
 contextual information derived from the visible ones. This re-
 construction objective drives the network to develop richer and
 more generalizable representations, particularly under adverse
 conditions such as fluctuating lighting, reflections, and partial
 occlusions. After pretraining, the decoder is discarded, and the
 pretrained encoder is retained as a feature extractor for fire de-
 tection. A lightweight classification head is then attached to
 the encoder's output, operating on the class token or pooled
 global features. During supervised fine-tuning, the combined
 encoder-classifier is trained on labeled fire datasets using cross-
 entropy loss, enabling precise discrimination between fire and
 non-fire categories. In deployment, the system bypasses mask-
 ing and processes complete images, ensuring efficient infer-
 ence. The encoder extracts high-level features, and the classifier
 outputs probability scores for fire presence.

The architectural choices are not arbitrary but motivated by
 the need for both robustness and efficiency. The MAE-based
 pretraining allows the encoder to learn from partially observed
 data, strengthening its ability to identify subtle fire cues. The
 ViT encoder contributes long-range attention, critical for dis-
 tinguishing flames from visually similar distractors such as ar-
 tificial lighting. Finally, the lightweight decoder ensures scal-
 ability for real-time surveillance. Together, these components
 form a unified and innovative framework that advances beyond
 prior approaches based on handcrafted color features or shallow
 CNNs, offering both interpretability and reliable performance
 in real-world fire detection scenarios.

1 **4. Experimental Results**

2 This section evaluates the performance of the proposed
 3 ImageMAE-based FD model on three publicly available FD
 4 benchmarks along with a self-curated fire dataset. The exper-
 5 imental analysis includes implementation details, dataset de-
 6 scriptions, evaluation metrics, ablation studies, cross-crop eval-
 7 uations, qualitative visualizations, and comparative results with
 8 state-of-the-art methods.

9 **4.1. Experimental Setup and Evaluation Metrics**

10 All experiments were carried out on a system equipped with
 11 an Intel Core i9 processor 3.60GHz, with NVIDIA GEFORCE
 12 RTX 3080 Ti GPU, and 64 GB of RAM. The model was im-
 13 plemented using the Keras deep learning API with TensorFlow
 14 as the backend. With a starting learning rate of 0.001 and a
 15 weight-decomposition approach that progressively lowers the
 16 learning rate to 0.0001 over the course of subsequent epochs,
 17 the model was optimised using the AdamW optimiser. The best
 18 results were obtained across all data sets when the training was
 19 conducted with a batch size of 32 for 50 epochs. The loss function
 20 used was sparse categorical cross-entropy.

21 To evaluate the performance, we utilized several standard
 22 metrics widely adopted in the FD domain: Accuracy (Acc),
 23 Precision (P), Recall (R), F1-score (F), False Positive Rate
 24 (FPR), and False Negative Rate (FNR). These metrics are com-
 25 puted using the following definitions:

- 26 • True Positive (TP): Number of correctly classified fire im-
 27 ages.
- 28 • True Negative (TN): Number of correctly classified non-
 29 fire images.
- 30 • False Positive (FP): Number of non-fire images incorrectly
 31 classified as fire.
- 32 • False Negative (FN): Number of fire images incorrectly
 33 classified as non-fire.

34 Based on these, the evaluation metrics are calculated as:

$$35 \quad \text{Acc} = \frac{TP + TN}{TP + FP + FN + TN} \quad (11)$$

$$36 \quad P = \frac{TP}{TP + FP} \quad (12)$$

$$37 \quad R = \frac{TP}{TP + FN} \quad (13)$$

$$38 \quad F = \frac{2 \times P \times R}{P + R} \quad (14)$$

$$39 \quad \text{FPR} = \frac{FP}{FP + TN} \quad (15)$$

$$40 \quad \text{FNR} = \frac{FN}{FN + TP} \quad (16)$$

41 Accuracy reflects the ratio of correct predictions over the to-
 42 tal number of instances. Precision measures the proportion of
 43 predicted fire instances that were actually fire. Recall evaluates
 44 the proportion of actual fire instances that were correctly iden-
 45 tified. The F1-score provides a harmonic mean of precision and
 46 recall. The FPR indicates the proportion of non-fire instances
 47 incorrectly predicted as fire, and FNR shows the proportion of
 48 fire instances incorrectly predicted as non-fire.

49 These evaluation metrics ensure a thorough and balanced as-
 50 sessment of the model's capability in both FD sensitivity and
 51 false alarm resistance.

52 **4.2. Benchmark Fire Detection Datasets**

53 To evaluate the effectiveness and robustness of the proposed
 54 FD framework, we employ a collection of publicly available
 55 benchmark datasets frequently used in FD research, shown in
 56 Table 1. These include Foggia's (FG) dataset [17], BoWFire
 57 (BWF) [20], DeepQuestAI Fire-Flame (DQFF) [40], the Large-
 58 Scale Fire Detection (LSFD) dataset [35], and the Drone Satel-
 59 lite Fire Dataset (DSFD) [50]. These datasets cover a wide
 60 range of challenging environments and are briefly described be-
 61 low.

62 **Foggia (FG):** The Foggia's dataset is composed of 31 video
 63 sequences recorded in both indoor and outdoor settings. It in-
 64 cludes various scenes captured under different lighting and en-
 65 vironmental conditions. The videos are converted into a total of
 66 62,690 frames, providing a rich source of temporal and visual
 67 diversity for training and evaluation purposes.

68 **BoWFire (BWF):** The BoWFire dataset is comparatively
 69 smaller in scale, containing a total of 226 still images. De-
 70 spite its limited size, it includes visually complex scenes such as
 71 sunsets, artificial lights, and flame-like regions, which present
 72 challenges to FD models and often contribute to false positives.

73 **DeepQuestAI Fire-Flame (DQFF):** The DeepQuestAI Fire-
 74 Flame dataset consists of 2,000 images obtained from diverse
 75 real-world environments, including urban streets, buildings,
 76 and natural landscapes. Its diversity in terms of background
 77 complexity and ambient lighting makes it valuable for assess-
 78 ing the generalization performance of FD models.

79 **Large-Scale Fire Detection (LSFD):** The Large-Scale FD
 80 dataset provides a significantly broader collection of fire-related
 81 scenes. With 50,000 high-resolution images, it was constructed
 82 by combining multiple datasets, such as FG and BWF, and aug-
 83 menting them with additional samples sourced from the inter-
 84 net. This dataset supports large-scale training and enables deep
 85 models to learn from a wide distribution of fire conditions.

86 **DFAN Dataset:** In the field of fire detection, most available
 87 datasets are limited to two categories: fire and normal. Such
 88 datasets mainly emphasize the identification of fire presence,
 89 without considering the specific objects affected by the flames.
 90 The DFAN dataset offers a notable advancement by incorporat-
 91 ing greater diversity and expanding the number of categories to

Table 1: Datasets description including number of samples, classes and environment.

Dataset	Classes	Samples	Environment
FG	2	62,690	The dataset comprises images captured in both indoor and outdoor settings, where red-colored objects are present near visible fire regions. .
BWF	2	226	A small-scale dataset featuring diverse and challenging indoor and outdoor environments.
DQFF	2	2,000	A medium-scale dataset containing both indoor and outdoor samples.
LSFD	2	50,000	A combined dataset consisting of the FG dataset and newly collected samples from the internet, including both indoor and outdoor scenes captured by CCTV and remote sensing devices.
DSFD	2	6,000	Outdoor samples from drones and satellites under varied angles, heights, times of day, and weather conditions.
DFAN	12	12000+	A large scale dataset with diverse range of challenges and classes.

Table 2: Comparative analysis with DL-based FD methods.

Dataset	Method	ACC	PR	RC	FS
ANetFire [25]	88.1	80.0	98.0	88.0	
GNetFire [45]	85.0	79.0	93.0	85.0	
CNNFire [46]	89.8	83.0	97.0	90.0	
BWF	EMNFire [47]	92.0	90.0	93.0	92.0
DQFF	DFAN [35]	95.0	95.0	94.0	95.0
	Ours	98.5	98.5	97.5	98.0
GNNetFire [45]	89.4	84.5	96.5	90.1	
EMNFire [47]	87.7	80.8	98.8	88.9	
DQFF	EfficientNetB0	95.4	91.8	97.6	94.8
	[38]	98.5	97.2	98.3	97.6
EMNFire [47]	92.8	88.3	98.7	93.2	
MIAPC [48]	96.0	94.9	97.3	96.1	
LSFD	MS-Net [49]	97.38	97.8	96.82	97.35
	Ours	99.0	99.6	98.2	98.9
EFDNet [33]	88.0	87.5	88.0	87.75	
ADFireNet [50]	90.86	90.9	90.86	90.88	
DSFD	M-SoftFireNet	93.50	93.57	93.51	93.53
	[51]	96.20	95.11	95.25	95.19
SE-EFFNet [5]	97.2	0.04 (FPR)	0.03 (FNR)	–	
STN-CNN [52]	96.2	3.68 (FPR)	2.46 (FNR)	–	
FG	Ours	98.80	0 (FPR)	0.02 (FNR)	–
LW [50]	90.0	90.43	90.49	89.99	
DFire [51]	91.20	90.63	91.17	90.36	
DFAN	DFAN [35]	89.36	86.1	94.00	89.84
	Ours	92.30	91.21	92.10	91.90

enhance recognition of objects in fire scenarios. Unlike conventional datasets that only classify fire versus non-fire, DFAN [51] introduces 12 distinct classes: fire on a boat, building, bus, car, cargo, electric pole, forest, non-fire, pick-up, SUV, train, and van, providing a more comprehensive representation of real-world fire events. This broader categorization is essential for capturing the complexity of fire situations and improving detection accuracy. The dataset is partitioned into training (70%), validation (20%), and testing (10%) subsets.

Drone Satellite Fire Dataset (DSFD): The Drone Satellite Fire Dataset incorporates a unique perspective by combining aerial views from drones and satellite images. Videos were captured using DJI drones at different heights, ranging from 10 to 70 meters, followed by a 60-frame skip mechanism to introduce diversity and eliminate redundancy. Images were also manually reviewed to ensure quality and relevance. The satellite component adds further variation in environmental settings such as day/night conditions and foggy weather, enhancing the

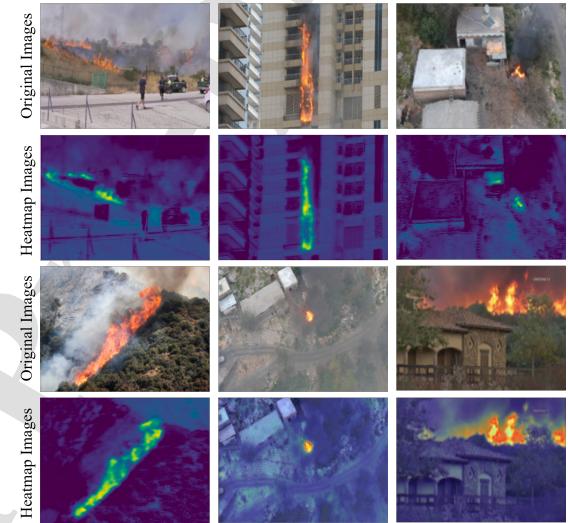


Figure 3: Results from our model visualized using Grad-CAM XAI methods to highlight important features and regions contributing to predictions.

complexity and utility of this dataset for high-altitude surveillance applications. To prepare these datasets for model training, we resized all images to a uniform resolution compatible with our architecture. A standard data split strategy is followed for all datasets, 20% and 10% of the data for training, validation, and testing. For smaller datasets like BWF, data augmentation techniques, including rotation, flipping, and zooming, were applied to enhance training diversity and mitigate overfitting. This comprehensive set of benchmark datasets provides a reliable foundation for evaluating the proposed model across multiple real-world FD scenarios.

4.3. Quantitative Analysis

The proposed method is comprehensively evaluated against TML and DL-based FD models across various datasets. The comparison uses six standard metrics: Accuracy (ACC), Precision (PR), Recall (RC), F1-score (FS), False Positive Rate

Table 3: Comparative analysis with TML-based FD methods.

Dataset	Method	ACC	PR	RC	FS
BWF [20]	FD-GCM [13]	–	55.0	54.0	54.0
	Bonfires[20]	–	51.0	65.0	67.0
	EFD-IP [53]	–	75.0	15.0	25.0
LSFD [33]	Ours	98.5	98.5	97.5	98.0
	FD-GCM [13]	69.6	63.9	90.0	74.7
	FFD-ANN [55]	71.7	71.1	73.2	72.1
FG [17]	FPC [54]	53.9	52.0	99.9	98.4
	Ours	99.0	99.6	98.3	97.6
	FD-CSM [17]	93.6	–	–	–
DSFD	FSD-YUV [35]	87.1	–	–	–
	FSD-RGB [35]	74.2	–	–	–
	FD-CV [35]	92.9	–	–	–
DSFD	CM-FDM [18]	90.3	5.9 (FPR)	14.2 (FNR)	–
	Ours	98.80	0 (FPR)	0.02 (FNR)	–

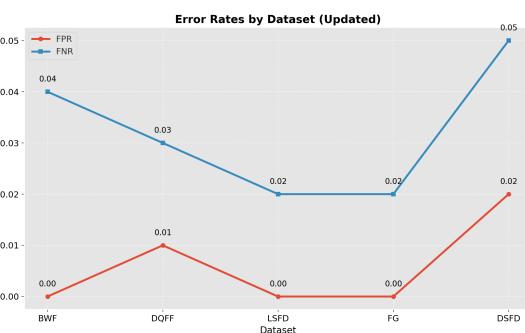


Figure 4: False Positive Rate (FPR) and False Negative Rate (FNR) across datasets showing low error rates for the proposed method.

(FPR), and False Negative Rate (FNR). The results, presented in Tables 3 and 2, demonstrate the superior performance of our method across all scenarios.

4.3.1. Comparison with TML-Based Methods

Table 3 reports the performance comparison with TML-based methods on BWF, LSFD, and FG datasets. Traditional methods such as FD-GCM [13] and Bonfires [20] performed poorly, especially in precision and recall. For example, on the BWF dataset, Bonfires achieved only 51.0% precision and 65.0% recall, while EFD-IP [53], despite a relatively high precision of 75.0%, failed drastically in recall (15.0%). Our method achieved 98.5% accuracy, 98.5% precision, 97.5% recall, and 98.0% F1-score, indicating strong performance in detecting fire across diverse conditions. On the LSFD dataset, our model reported 99.0% accuracy, substantially outperforming FD-GCM [54] (69.6%) and FFD-ANN [55] (71.7%). On the FG [17] dataset, where previous TML methods such as FD-CSM [17] and FSD-RGB [35] suffered from high FPR and FNR, our method achieved near-perfect performance with 0% FPR and only 0.02% FNR.

4.3.2. Comparison with DL-Based Methods

Our method consistently outperforms DL-based methods on BWF, DQFF, LSFD, DSFD, and FG datasets, as shown in Table 2. On BWF, it surpassed models like DFAN (95.0% ACC)

and EMNFire (92.0%) by achieving 98.5% accuracy and the highest scores across all other metrics. For the DQFF dataset, our model improved upon EfficientNetB0 (95.4%) and EMNFire (87.7%), reaching 97.6% accuracy and 98.5% recall. In the LSFD dataset, our approach outperformed advanced networks such as MS-Net and MIAPC, achieving 99.0% across all primary metrics. Our model also performed best on DSFD, a dataset collected from aerial views, reaching 96.20% accuracy. Finally, on FG, the proposed method achieved the highest accuracy of 98.80%, reducing FPR to 0% and FNR to 0.02%, even outperforming SE-EFFNet and EMNFire. These results validate the generalization and robustness of the proposed method across varying environments and fire characteristics. To further evaluate the performance of the proposed method, we used the DFAN [35] multiclass dataset, which is one of the more challenging datasets. The proposed method achieved the highest accuracy of 92.30%, with a precision of 91.21%, recall of 92.10%, and F1-score of 91.90%. These results demonstrate that the proposed method generalizes well to large-scale and challenging datasets and outperforms other state-of-the-art methods, as shown in Table 2. These results confirm that our method achieves strong generalization performance and maintains a lower false alarm rate across diverse and challenging FD scenarios. The higher performance across all metrics validates the effectiveness of the proposed attention-based multi-scale architecture.

4.4. Qualitative comparison

Figure 6 presents a qualitative comparison between the proposed method and two recent FD approaches. The examples include challenging fire and non-fire scenes, such as artificial lighting, reflected glows, mountain snow caps, and dense smoke without visible flames. The proposed method correctly identifies all scenarios, while the baseline methods from Khan et al. and Yar et al. show multiple misclassifications, particularly in visually ambiguous cases. These outcomes demonstrate the effectiveness of our transformer-based model in learning nuanced fire features and reducing false alarms under diverse conditions. The Confusion matrix of the proposed method on all five datasets is shown in Figure 5. Figure 4 shows the error rate of the model over all datasets, whereas the ROC curve is shown in Figure 7. Furthermore, Figure 3 illustrates the Grad-CAM visualizations of our model’s predictions on fire images. The highlighted regions indicate the areas that the model considers most important when detecting fire. As shown, the model consistently focuses on the flames and smoke regions, demonstrating its ability to accurately localize key features associated with fire. This confirms that our model not only achieves high classification performance but also provides interpretable insights into the decision-making process, reinforcing its reliability for practical fire detection applications.

4.5. Ablation Study

To assess the influence of the encoder backbone within our FD framework, we performed a model-based ablation study by replacing the original Vision Transformer (ViT) encoder with

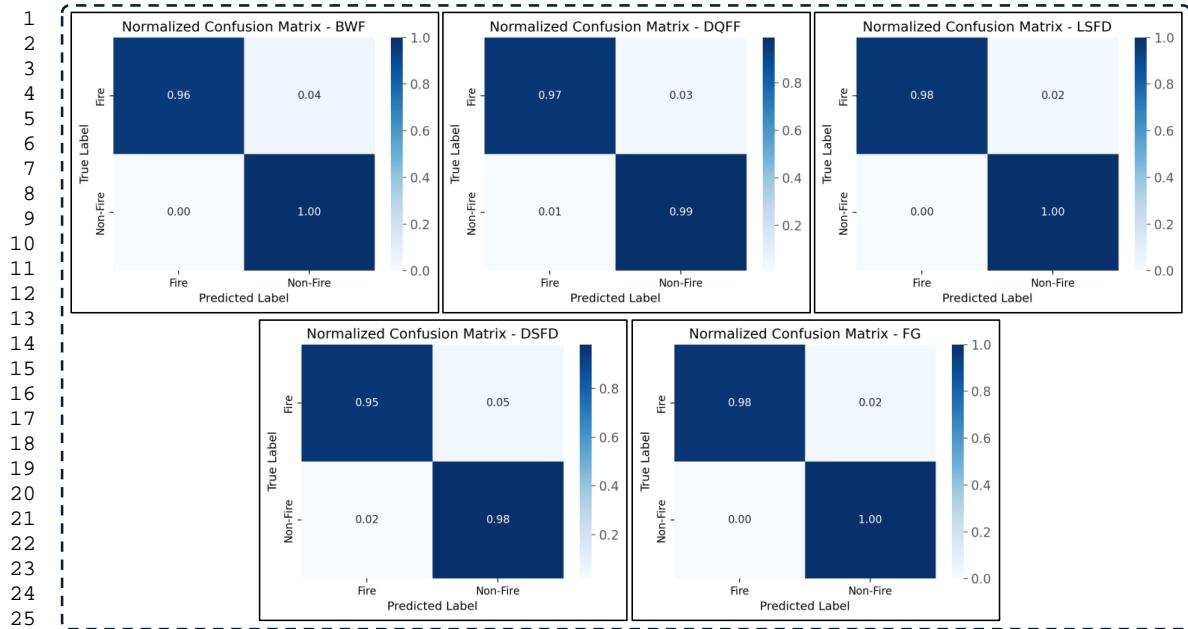


Figure 5: Confusion matrices for the proposed method across five datasets (BWF, DQFF, LSFD, DSFD, and FG), showing strong classification performance with high true positive and true negative rates.

Table 4: Ablation study on different masking ratios across five benchmark fire detection datasets.

Masking Ratio	BWF (%)	DQFF (%)	LSFD (%)	DSFD (%)	FG (%)
25%	91.2	89.5	88.7	90.1	87.9
50%	92.8	91.0	90.1	91.7	89.3
75%	98.5	98.5	99	96.20	98.80

several widely-used alternatives: Swin Transformer, ResNet50, ConvNeXt-Tiny, EfficientNetB0, and DenseNet121. All variants were embedded within the same masked autoencoder (ImageMAE) pipeline, using consistent training hyperparameters and loss design. The results, summarized in Table 5, reflect accuracy scores across five benchmark datasets. The original ViT encoder consistently outperformed all other backbones, achieving 98.80% on FG and 99.0% on LSFD, with similar superiority across the remaining datasets. Swin Transformer, which shares architectural similarities with ViT but incorporates local window attention, followed closely, showing promising performance with minimal drop in accuracy. In contrast, CNN-based encoders such as ResNet50, DenseNet121, and EfficientNetB0 demonstrated lower accuracy, particularly on complex datasets like DSFD and LSFD. This suggests a limitation in their ability to capture global dependencies and contextual semantics, which are critical for accurate fire localization and classification. ConvNeXt, a modernized CNN incorporating transformer-like features, performed better than traditional CNNs but still lagged behind transformer-based architectures. These findings affirm that transformer-based models, particularly ViT and Swin, are

more adept at modeling the complex spatial patterns necessary for robust FD in varied real-world environments.

Moreover, we further conducted an ablation study to evaluate the impact of different masking ratios on fire detection performance across five benchmark datasets (BWF, DQFF, LSFD, DSFD, and FG). As shown in Table 4, increasing the masking ratio from 25% to 75% consistently improves FD accuracy on all datasets. Lower masking ratios (25%) result in easier reconstruction but less robust feature learning, while a moderate ratio (50%) provides a better balance. The 75% masking ratio achieves the best performance across all datasets, encouraging the encoder to focus on salient fire-related features and learn more discriminative representations. These results justify the use of 75% masking for all experiments reported in this work.

4.6. Complexity Analysis

The inference time, model size, and computational complexity of CANet are compared with several state-of-the-art fire detection methods. The processing speed of AI-based models is largely influenced by their computational cost (FLOPS) and model size. Table 6 presents a comparison of proposed method with SE-EFFNet [5], EMNFE [47], EFDNet [33], GNetFire [45], ResNetFire [22], CNNFire [46], and DFAN [35]. Among these methods, EFDNet [33] and CNNFire [46] have the smallest model sizes of 4.8 MB and 3 MB, respectively, but their FLOPS are relatively higher (1130 M and 720 M), which limits their inference speed on low-computation devices. In contrast, DFAN [35] exhibit lower FLOPS (73.05 M), resulting in

1					
2					
3	Input				
4					
5	Ground Truth	Fire	Non-Fire	Fire	Fire
6	Yar et al. [36]	Fire	Fire	Fire	Fire
7	Khan et al. [40]	Fire	Non-Fire	Fire	Non-Fire
8	Our	Fire	Non-Fire	Fire	Fire
9					
10					
11	Input				
12					
13	Ground Truth	Fire	Fire	Fire	Fire
14	Yar et al. [36]	Fire			Fire
15	Khan et al. [40]	Non-Fire		Fire	Fire
16	Our	Fire		Fire	Fire
17					
18	Input				
19					
20	Ground Truth	Non-Fire	Fire	Non-Fire	Fire
21	Yar et al. [36]	Non-Fire	Fire	Non-Fire	Fire
22	Khan et al. [40]	Non-Fire	Fire	Non-Fire	Non-Fire
23	Our	Non-Fire	Fire	Fire	Fire
24					
25					
26					
27					
28					
29	Input				
30					
31					
32	Ground Truth	Non-Fire	Fire	Non-Fire	Fire
33	Yar et al. [36]	Non-Fire	Fire	Non-Fire	Fire
34	Khan et al. [40]	Non-Fire	Fire	Non-Fire	Non-Fire
35	Our	Non-Fire	Fire	Fire	Fire
36					

Figure 6: Qualitative comparison between the proposed method and two state-of-the-art models: Khan et al. and Yar et al. Each row shows predictions from a different method across various fire and non-fire scenarios. Black text indicates correct classification, while red text denotes misclassification. The proposed method (last row) demonstrates superior robustness in distinguishing fire from challenging non-fire cases, such as sunsets, lights, and smoke-like patterns.

higher FPS across all tested hardware, including Raspberry Pi (RPI), CPU, and GPU. The proposed method achieves the lowest FLOPS of 55.40 M and delivers the fastest inference speed, with 9 FPS on RPI, 50 FPS on CPU, and 130 FPS on GPU, outperforming all other methods in the comparison. This demonstrates that proposed method is highly efficient and suitable for real-time deployment on edge devices.

5. Conclusion

This study presents a transformer-based FD framework that addresses key limitations in prior deep learning methods, particularly those relying on shallow architectures or simplistic feature representations. While existing models often struggle with fine-grained fire discrimination and are limited by binary-class datasets, our method leverages a masked autoencoding mechanism to enhance visual understanding from partial observations. By incorporating a Vision Transformer encoder and reconstructive learning, the model learns robust, context-aware features

1 corporating transformer-based segmentation or detection modules. This would enable the system to highlight precise fire
2 zones rather than merely classifying the presence of fire. Additionally, attention will be given to developing modules that can
3 better interpret low-visibility scenarios, such as scenes dominated by smoke, haze, or indirect fire indicators.

4 We also plan to enhance the scalability of the framework by evaluating its performance across a broader range of datasets representing indoor, outdoor, urban, and rural fire conditions.
5 To this end, we will extend our dataset to include instances where fire is partially or entirely obscured, and where only smoke or glow is visible. These additions aim to improve the model's resilience in challenging environments.

6 Finally, future research may explore hybrid learning strategies that combine self-supervised reconstruction with supervised fire-type classification or region regression. By advancing the model's granularity, adaptability, and interpretability, we aim to bring the system closer to deployment in practical, high-stakes fire monitoring applications.

22 CRediT authorship contribution statement

23 **Hanxiang Wang:** Conceptualization, Methodology, implementation, Writing – original draft.

24 **Muhammad Fayaz:** Conceptualization, Methodology, Implementation, Visualization, Writing – original draft.

25 **Awais Ahmad:** Data curation, Validation, Writing – review & editing. **Yanfen Li:** Resources, formal analysis, Writing – review & editing. **Tan N. Nguyen:** Investigation, Writing – review & editing. **L. Minh Dang:** Supervision, Writing – review, and editing, Funding Aquisition.

35 Declaration of competing interest

36 The authors claim that they have no known conflicting financial interests or personal ties that may have seemed to affect the work presented in this study.

43 Acknowledgments

44 This work was supported by the Young Scientists Fund of the National Natural Science Foundation of China (No. 62502271), the Young Scientists Fund of the Natural Science Foundation of Shandong Province (No. ZR2025QC630), the Natural Science Foundation of Rizhao City (Nos. RZ2024ZR33 and RZ2024ZR34).

53 References

- [1] M. Altaf, M. Yasir, N. Dilshad, W. Kim, An optimized deep-learning-based network with an attention module for efficient fire detection., *Fire* (2571-6255) 8 (1) (2025).
- [2] S. Wang, M. Wu, X. Wei, X. Song, Q. Wang, Y. Jiang, J. Gao, L. Meng, Z. Chen, Q. Zhang, et al., An advanced multi-source data fusion method utilizing deep learning techniques for fire detection, *Engineering Applications of Artificial Intelligence* 142 (2025) 109902.
- [3] M. Wang, P. Yue, L. Jiang, D. Yu, T. Tuo, J. Li, An open flame and smoke detection dataset for deep learning in remote sensing based fire detection, *Geo-spatial Information Science* 28 (2) (2025) 511–526.
- [4] D. Gragnaniello, A. Greco, C. Sansone, B. Vento, Flame: fire detection in videos combining a deep neural network with a model-based motion analysis, *Neural Computing and Applications* 37 (8) (2025) 6181–6197.
- [5] Z. A. Khan, T. Hussain, F. U. M. Ullah, S. K. Gupta, M. Y. Lee, S. W. Baik, Randomly initialized cnn with densely connected stacked autoencoder for efficient fire detection, *Engineering Applications of Artificial Intelligence* 116 (2022) 105403.
- [6] T. Maity, A. N. Bhawani, J. Samanta, P. Saha, S. Majumdar, G. Srivastava, Mlsfdd: Machine learning-based smart fire detection device for precision agriculture, *IEEE Sensors Journal* (2025).
- [7] Y. Li, Y. Wang, X. Shao, A. Zheng, An efficient fire detection algorithm based on mamba space state linear attention, *Scientific Reports* 15 (1) (2025) 11289.
- [8] A. Hussain, H. Yar, N. Khan, Z. A. Khan, M. J. Kim, S. W. Baik, Dual stream deep attention networks for annual population projection, *Pattern Analysis and Applications* 28 (2) (2025) 71.
- [9] A. Hussain, W. Ullah, N. Khan, Z. A. Khan, M. J. Kim, S. W. Baik, Tds-net: Transformer enhanced dual-stream network for video anomaly detection, *Expert Systems with Applications* 256 (2024) 124846.
- [10] A. Hussain, W. Ullah, N. Khan, Z. A. Khan, H. Yar, S. W. Baik, Class-incremental learning network for real-time anomaly recognition in surveillance environments, *Pattern Recognition* (2025) 112064doi:
<https://doi.org/10.1016/j.patcog.2025.112064>.
- [11] H. Harkat, H. F. T. Ahmed, J. M. Nascimento, A. Bernardino, Early fire detection using wavelet based features, *Measurement* 242 (2025) 115881.
- [12] H. Harkat, J. M. Nascimento, A. Bernardino, H. F. T. Ahmed, Fire images classification based on a handcraft approach, *Expert Systems with Applications* 212 (2023) 118594.
- [13] T. Çelik, H. Özkarahanlı, H. Demirel, Fire and smoke detection without sensors: Image processing based approach, in: 2007 15th European signal processing conference, IEEE, 2007, pp. 1794–1798.
- [14] A. Rafiee, R. Dianat, M. Jamshidi, R. Tavakoli, S. Abbaspour, Fire and smoke detection using wavelet analysis and disorder characteristics, in: 2011 3rd International conference on computer research and development, Vol. 3, IEEE, 2011, pp. 262–265.
- [15] T. Qiu, Y. Yan, G. Lu, An autoadaptive edge-detection algorithm for flame and fire image processing, *IEEE Transactions on instrumentation and measurement* 61 (5) (2011) 1486–1493.
- [16] S. G. Kong, D. Jin, S. Li, H. Kim, Fast fire flame detection in surveillance video using logistic regression and temporal smoothing, *Fire Safety Journal* 79 (2016) 37–43.
- [17] P. Foggia, A. Saggese, M. Vento, Real-time fire detection for video-surveillance applications using a combination of experts based on color, shape, and motion, *IEEE TRANSACTIONS on circuits and systems for video technology* 25 (9) (2015) 1545–1556.
- [18] Y. H. Habiboglu, O. Güney, A. E. Çetin, Covariance matrix-based fire and flame detection method in video, *Machine Vision and Applications* 23 (2012) 1103–1113.
- [19] Z. Teng, J.-H. Kim, D.-J. Kang, Fire detection based on hidden markov models, *International Journal of Control, Automation and Systems* 8 (2010) 822–830.
- [20] D. Y. Chino, L. P. Avalhais, J. F. Rodrigues, A. J. Traina, Bowfire: detection of fire in still images by integrating pixel color and texture analysis, in: 2015 28th SIBGRAPI conference on graphics, patterns and images, IEEE, 2015, pp. 95–102.
- [21] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.
- [22] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, *arXiv preprint arXiv:1409.1556* (2014).
- [23] S. Frizzi, R. Kaabi, M. Bouchouicha, J.-M. Ginoux, E. Moreau, F. Fnaiech, Convolutional neural network for video fire and smoke detection, in: IECON 2016-42nd Annual Conference of the IEEE Industrial Electronics Society, IEEE, 2016, pp. 877–882.
- [24] O. Maksymiv, T. Rak, D. Peleshko, Real-time fire detection method combining adaboost, lbp and convolutional neural network in video sequence, in: 2017 14th international conference the experience of designing and application of CAD Systems in microelectronics (CADSM), IEEE, 2017,

- 1 pp. 351–353.
- 2 [25] K. Muhammad, J. Ahmad, S. W. Baik, Early fire detection using convolutional neural networks during surveillance for effective disaster management, *Neurocomputing* 288 (2018) 30–42.
- 3 [26] A. J. Dunnings, T. P. Breckon, Experimentally defined convolutional neural network architecture variants for non-temporal real-time fire detection, in: 2018 25th IEEE international conference on image processing (ICIP), IEEE, 2018, pp. 1558–1562.
- 4 [27] F. Saeed, A. Paul, P. Karthigaikumar, A. Nayyar, Convolutional neural network based early fire detection, *Multimedia Tools and Applications* 79 (13) (2020) 9083–9099.
- 5 [28] G. Zhang, M. Wang, K. Liu, Forest fire susceptibility modeling using a convolutional neural network for yunnan province of china, *International Journal of Disaster Risk Science* 10 (3) (2019) 386–403.
- 6 [29] N. Dilshad, T. Khan, J. Song, Efficient deep learning framework for fire detection in complex surveillance environment., *Comput. Syst. Sci. Eng.* 46 (1) (2023) 749–764.
- 7 [30] T. Khan, Z. A. Khan, C. Choi, Enhancing real-time fire detection: An effective multi-attention network and a fire benchmark, *Neural Computing and Applications* 37 (18) (2025) 11693–11707.
- 8 [31] Z. Wang, L. Xu, Z. Chen, Ffd-yolo: A modified yolov8 architecture for forest fire detection, *Signal, Image and Video Processing* 19 (3) (2025) 265.
- 9 [32] K. Lv, R. Wu, S. Chen, P. Lan, Cci-yolov8n: Enhanced fire detection with carafe and context-guided modules, in: International Conference on Intelligent Computing, Springer, 2025, pp. 128–140.
- 10 [33] S. Li, Q. Yan, P. Liu, An efficient fire detection method based on multi-scale feature extraction, implicit deep supervision and channel attention mechanism, *IEEE Transactions on Image Processing* 29 (2020) 8467–8475.
- 11 [34] F. Yuan, K. Li, C. Wang, Z. Fang, A lightweight network for smoke semantic segmentation, *Pattern Recognition* 137 (2023) 109289.
- 12 [35] H. Yar, T. Hussain, M. Agarwal, Z. A. Khan, S. K. Gupta, S. W. Baik, Optimized dual fire attention network and medium-scale fire classification benchmark, *IEEE Transactions on Image Processing* 31 (2022) 6331–6343.
- 13 [36] H. Yar, Z. A. Khan, T. Hussain, S. W. Baik, A modified vision transformer architecture with scratch learning capabilities for effective fire detection, *Expert Systems with Applications* 252 (2024) 123935.
- 14 [37] W. Wang, L. Zhang, T. Yang, S. Ma, Q. Zhang, P. Shi, F. Ding, Combined serum free light chain predicts prognosis in acute kidney injury following cardiovascular surgery, *Renal Failure* 44 (1) (2022) 1–10.
- 15 [38] S. Majid, F. Alenezi, S. Masood, M. Ahmad, E. S. Gündüz, K. Polat, Attention based cnn model for fire detection and localization in real-world images, *Expert Systems with Applications* 189 (2022) 116114.
- 16 [39] M. Jiang, Y. Zhao, F. Yu, C. Zhou, T. Peng, A self-attention network for smoke detection, *Fire safety journal* 129 (2022) 103547.
- 17 [40] Z. A. Khan, F. U. M. Ullah, H. Yar, W. Ullah, N. Khan, M. J. Kim, S. W. Baik, Optimized cross-module attention network and medium-scale dataset for effective fire detection, *Pattern Recognition* 161 (2025) 111273.
- 18 [41] X. Yu, C.-H. Chen, A robust operators' cognitive workload recognition method based on denoising masked autoencoder, *Knowledge-Based Systems* 301 (2024) 112370.
- 19 [42] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, R. Girshick, Masked autoencoders are scalable vision learners, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2022, pp. 16000–16009.
- 20 [43] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, et al., An image is worth 16x16 words: Transformers for image recognition at scale, *arXiv preprint arXiv:2010.11929* (2020).
- 21 [44] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, Bert: Pre-training of deep bidirectional transformers for language understanding, in: Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers), 2019, pp. 4171–4186.
- 22 [45] K. Muhammad, J. Ahmad, I. Mehmood, S. Rho, S. W. Baik, Convolutional neural networks based fire detection in surveillance videos, *Ieee Access* 6 (2018) 18174–18183.
- 23 [46] K. Muhammad, J. Ahmad, Z. Lv, P. Bellavista, P. Yang, S. W. Baik, Efficient deep cnn-based fire detection and localization in video surveillance applications, *IEEE Transactions on Systems, Man, and Cybernetics: Systems* 49 (7) (2018) 1419–1434.
- 24 [47] K. Muhammad, S. Khan, M. Elhoseny, S. H. Ahmed, S. W. Baik, Efficient fire detection for uncertain surveillance environment, *IEEE Transactions on Industrial Informatics* 15 (5) (2019) 3113–3122.
- 25 [48] Z. Deng, S. Hu, S. Yin, Y. Wang, A. Basu, I. Cheng, Multi-step implicit adams predictor-corrector network for fire detection, *IET Image Processing* 16 (9) (2022) 2338–2350.
- 26 [49] J. Feng, Y. Sun, Multiscale network based on feature fusion for fire disaster detection in complex scenes, *Expert Systems with Applications* 240 (2024) 122494.
- 27 [50] H. Yar, W. Ullah, Z. A. Khan, S. W. Baik, An effective attention-based cnn model for fire detection in adverse weather conditions, *ISPRS Journal of Photogrammetry and Remote Sensing* 206 (2023) 335–346.
- 28 [51] H. Yar, Z. A. Khan, I. Rida, W. Ullah, M. J. Kim, S. W. Baik, An efficient deep learning architecture for effective fire detection in smart surveillance, *Image and Vision Computing* 145 (2024) 104989.
- 29 [52] S. B. Avula, S. J. Badri, G. Reddy, A novel forest fire detection system using fuzzy entropy optimized thresholding and stn-based cnn, in: 2020 international conference on COMmunication Systems & NETworks (COMSNETS), IEEE, 2020, pp. 750–755.
- 30 [53] T.-H. Chen, P.-H. Wu, Y.-C. Chiou, An early fire-detection method based on image processing, in: 2004 International Conference on Image Processing, 2004, ICIP'04., Vol. 3, IEEE, 2004, pp. 1707–1710.
- 31 [54] T. Celik, H. Ozkaramanli, H. Demirel, Fire pixel classification using fuzzy logic and statistical color model, in: 2007 IEEE International Conference on Acoustics, Speech and Signal Processing-ICASSP'07, Vol. 1, IEEE, 2007, pp. I-1205.
- 32 [55] D. Zhang, S. Han, J. Zhao, Z. Zhang, C. Qu, Y. Ke, X. Chen, Image based forest fire detection using dynamic characteristics with artificial neural networks, in: 2009 international joint conference on artificial intelligence, IEEE, 2009, pp. 290–293.

Title Page

Masked Autoencoder-Based Vision Framework for Robust Fire Detection in Complex Environments
Authors:

1. Hanxiang Wang [†]
Research Assistant
School of Computer Science, Qufu Normal University, Rizhao, China
Email: hanxiang@qfnu.edu.cn
2. Muhammad Fayaz
Research Assistant
Department of Computer Science and Engineering, Sejong University, Seoul 05006, South Korea
Email: muhammadfayaz@sju.ac.kr
3. Awais Ahmad
Research Assistant
Department of Computer Science, Islamia college university, Peshawar 25000, Pakistan
Email: islamian1@gmail.com
4. Yanfen Li *
Research Assistant
School of Computer Science, Qufu Normal University, Rizhao, China
Email: hanxiang@qfnu.edu.cn
5. Tan N. Nguyen [†]
Assistant Professor
Department of Architectural Engineering, Sejong University, Seoul, 05006, South Korea,
Email: tnnguyen@sejong.ac.kr

Corresponding Author:

Prof: L. Minh Dang *
Professor
The Institute of Research and Development, Duy Tan University, Da Nang, 550000, Viet Nam,
Faculty of Information Technology, Duy Tan University, Da Nang, 550000, Viet Nam,
Email: danglienminh@duytan.edu.vn

Keywords:

- Fire Detection
- Masked Autoencoder
- Vision Transformer
- Self-Supervised Learning
- Feature Reconstruction
- Attention Mechanism

[†] Hanxiang Wang and Tan N. Nguyen are co-first authors and contributed equally to the work.

* These authors share corresponding authorship

Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

