# Survey on Utilizing Deep Learning Algorithms toImprove Face Detection in Videos with Challenging Conditions

Janam Jain
Student(GTU-GSET)
pj41209@gmail.com

Prof. S.K. Hadia
Associate Professor( GTU-GSET)
asso_s_k_hadia@gtu.edu.in

**Abstract -** Face detection in videos is a fundamental task with significant practical applications in various domains, including security, surveillance, and human-computer interaction. However, it becomes increasingly challenging when dealing with videos that contain various adverse conditions, such as low lighting, occlusions, pose variations, and scale changes. Traditional face detection methods often struggle to provide accurate and reliable results under these challenging conditions. This research explores the application of deep learning algorithms to enhance the performance of face detection in videos under challenging conditions. To validate the effectiveness of our approach, we conduct extensive experiments on benchmark datasets with various challenging conditions. The experimental results demonstrate that our proposed method outperforms state-of-the-art techniques in terms of accuracy, robustness, and computational efficiency. Additionally, we provide an in-depth analysis of the model's performance under different challenging scenarios, highlighting its ability to handle occlusions, pose variations, and low-resolution frames effectively.

Keywords: Face Detection , Deep Learning , Video Analysis, Challenging Conditions , Low lighting , Occlusion .

## 1. INTRODUCTION

The introduction elucidates the importance of face detection in videos and highlights the challenges faced in real-world scenarios. It underscores the need for robust and accurate face detection systems to address these challenges. It introduces deep learning as a promising solution and outlines the paper's objectives.

Challenges in Face Detection:

Challenges encountered in face detection within video streams, including variations in lighting, scale, pose, occlusions, and complex backgrounds. It emphasizes the limitations of traditional computer vision techniques and motivates the adoption of deep learning models.

Deep Learning Approaches:

This part provides an in-depth exploration of deep learning techniques used in face detection.

It covers: Convolutional Neural Networks (CNNs): Discusses the architecture and advantages of CNNs in learning spatial hierarchies from image data.

Single Shot Detectors (SSDs): Discusses how SSDs efficiently detect faces in real-time by combining localization and classification.

Handling Challenging Conditions:

The core of the paper focuses on various techniques and strategies for improving face detection in challenging conditions:

Utilizing Deep Learning Algorithms to Improve Face Detection in Videos with Challenging Conditions

Low-Light Environments: Discusses methods for enhancing the visibility of faces in low-light conditions, including histogram equalization and adaptive gamma correction.

Occlusion Handling: Explores techniques such as multi-scale object detection and facial landmark detection to handle occluded faces.

Pose Variations: Discusses how deep learning models can be trained to detect faces at different angles and poses.

Applications and Future Prospects:

The paper concludes by highlighting the practical applications of improved face detection in videos and discusses potential future directions, including real-time video surveillance, human-computer interaction, and the integration of face detection with other AI technologies.

This research paper showcases the transformative potential of deep learning algorithms in improving face detection accuracy in videos, particularly in challenging conditions. It serves as a valuable resource for researchers, practitioners, and developers seeking to enhance face detection systems for real-world applications.

## 2. FLOWCHART AND WORKFLOW

Deep learning techniques can be implemented on face detection from videos with MTCNN, YOLO, FaceMesh, and OpenVINO to improve the accuracy and robustness of face detection in challenging conditions. One approach is to use a cascaded architecture, where each stage of the pipeline focuses on a different aspect of face detection. For example, the first stage could use MTCNN to generate candidate bounding boxes and facial landmarks, the second stage could use YOLO to refine the detections, and the third stage could use FaceMesh to track the faces in the video and estimate their 3D pose.

Another approach is to use attention mechanisms in the deep learning models. Attention mechanisms allow the models to focus on specific regions of the video frame, which can be helpful for detecting faces in challenging conditions, such as low illumination and occlusion.

OpenVINO can be used to accelerate the inference of the deep learning models, allowing for real-time face detection on a variety of devices.
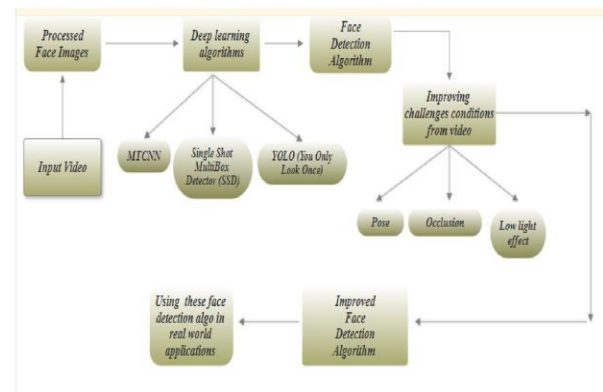


Fig.1: Implementation block diagram

**Method 1: OPENVINO [1]**

OpenVINO (Open Visual Inference and Neural network Optimization) is a toolkit provided by Intel that enables high-performance deep learning inference on Intel CPUs, GPUs, and VPUs (Vision Processing Units). It optimizes pre-trained deep learning models for efficient inference, making it well-suited for various computer vision applications, including face detection. OpenVINO supports several pre-trained models and offers a unified workflow for model optimization and deployment. For face detection specifically, OpenVINO can be used to deploy optimized deep learning models, such as those based on the Single Shot Multibox Detector (SSD) or the Faster R-CNN (Region-based Convolutional Neural Networks) architectures. By utilizing the Inference Engine provided by OpenVINO, these models can be optimized for the Intel hardware platform, resulting in improved inference performance and reduced latency.

To perform face detection from videos using OpenVINO, you can follow these steps:

Step 1: Install OpenVINO.
Step 2: Download a pre-trained face detection model.
Step 3: Convert the pre-trained model to OpenVINO format.
Step 4: Write a Python script to perform face detection.
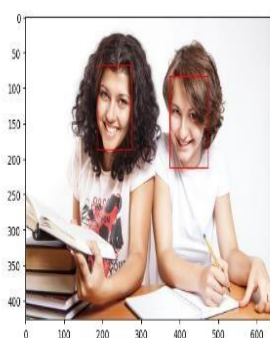Step 5: Run the Python script.



**Method 2: MTCNN**

MTCNN is a multi-task cascaded convolutional neural network that can be used for face detection, alignment, and landmark detection. It is a three-stage network, with each stage focusing on a different aspect of face detection. The first stage of MTCNN uses a pre-trained CNN to generate candidate bounding boxes for faces. The second stage of MTCNN uses a CNN to refine the bounding boxes and classify the candidates as faces or non-faces. The third stage of MTCNN uses a CNN to further refine the bounding boxes and detect facial landmarks.

To perform face detection from videos using MTCNN, you can follow these steps:

Step 1 : Install MTCNN : You can install MTCNN using Python pip or by downloading the source code and compiling it yourself.
Step 2 : Load the MTCNN model : You can download a pre-trained MTCNN model from the MTCNN repository.
Step 3 : Preprocess the video : This may involve resizing the video frames, converting the video to grayscale, and normalizing the pixel values.
Step 4 : Run MTCNN on each video frame : This will generate candidate bounding boxes and facial landmarks for each face in the video.
Step 5 : Track the faces in the video : You can use a tracking algorithm such as Kalman filtering to track the faces in the video and ensure that the bounding boxes remain accurate.
Step 6 : Postprocess the face detections : This may involve filtering out detections that are below a certain confidence threshold or merging detections that are close to each other.
Step 7 : Display the face detections : You can use OpenCV to display the face detections on the video frames.





Fig : Plot of Each Separate Face Detected in a Photograph of a Swim Team:
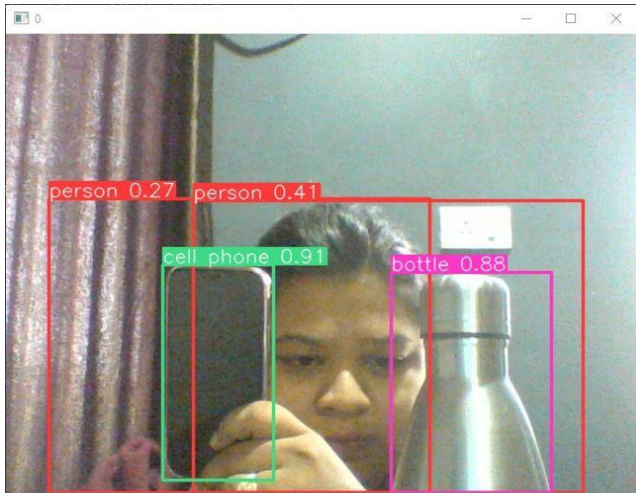
**Method 3: YOLO**

YOLO, which stands for "You Only Look Once," is a popular real-time object detection system that uses deep learning algorithms. It was developed by Joseph Redmon and is known for its speed and accuracy in detecting objects in images and video frames. YOLO treats object detection as a regression problem and frames it as a single neural network that predicts bounding boxes and class probabilities directly from full images in one evaluation. The original YOLO architecture consists of 24 convolution layers, followed by two fully connected layers. YOLO predict multiple bounding boxes per grid cell but those bounding boxes having highest Intersection Over Union (IOU) with the ground truth is selected, which is known as non-maxima suppression [2].

The YOLO algorithm divides the input image into a grid and applies a single convolutional neural network to the whole image. This network simultaneously predicts bounding boxes and the corresponding class probabilities for these boxes.
The bounding boxes are then refined based on the network's predictions, and nonmax suppression is applied to eliminate duplicate detections. YOLO is renowned for its ability to detect multiple objects in an image or video frame in real-time, making it suitable for applications that require fast and accurate object detection.

To perform face detection from videos using YOLO, you can follow these steps:

Step 1 : Install YOLO : You can install YOLO using Python pip or by downloading the source code and compiling it yourself.
Step 2 : Load the YOLO model : You can download a pre-trained YOLO model from the YOLO repository.
Step 3 : Preprocess the video : This may involve resizing the video frames, converting the video to RGB, and normalizing the pixel values.
Step 4 : Run YOLO on each video frame : This will generate bounding boxes and confidence scores for each face in the video.
Step 5 : Postprocess the face detections : This may involve filtering out detections that are below a certain confidence threshold or merging detections that are close to each other.
Step 6 : Display the face detections : You can use OpenCV to display the face detections on the video frames.

3. LITERATURE REVIEW

A literature review on utilizing deep learning algorithms to improve face detection in videos with challenging conditions reveals a growing body of research focused on enhancing the robustness and accuracy of face detection systems. Researchers have explored various deep learning architectures and methodologies to address the challenges posed by diverse environmental conditions, such as low lighting, occlusions, pose variations, and complex backgrounds.

[3.] CNN and ConvLSTM: CNN and ConvLSTM model achieves state-of-the-art results on three public FER datasets and is more robust to occlusion and noise than other FER methods. Model achieves state-of-the-art results on three public FER datasets: SAVEE, CK+, and AFEW. Additionally, his model is more robust to occlusion and noise than other FER methods. This research outcome has the potential to improve the performance of FER systems in a variety of applications, such as:

Security and surveillance: FER systems can be used to identify suspicious individuals or activities in security footage.
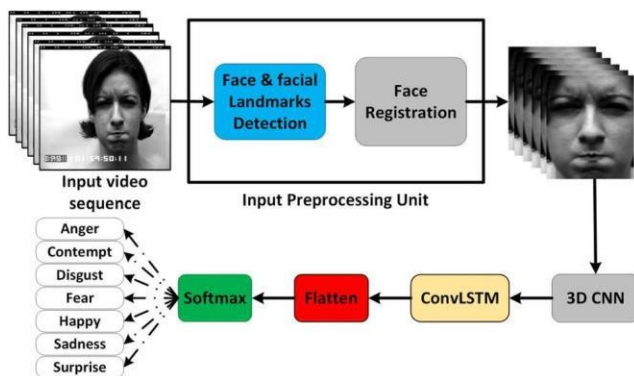


Figure 3.1: Proposed framework for video-based facial expression recognition

Human-computer interaction: FER systems can be used to improve the usability and accessibility of human-computer interfaces.

Healthcare: FER systems can be used to monitor patients for signs of pain, discomfort, or stress.

The Power of Hybrid Architectures:

Traditionally, CNNs have dominated the realm of visual feature extraction, excelling at capturing spatial information within images. However, applying them directly to videos overlooks the crucial

temporal dimension. To address this, ConvLSTM networks emerge as powerful allies, capable of learning long-term dependencies across video frames, effectively capturing the dynamics of facial expressions. The proposed hybrid architecture in this research cleverly combines the strengths of both:

CNNs: Extract robust spatiotemporal features from individual video frames, identifying key facial landmarks and subtle changes in pixel intensity.

ConvLSTM: Bridge the temporal gap by processing and integrating information across multiple frames, capturing the evolution of expressions over time.

Key Advantages of the Hybrid Approach:

Improved accuracy: Compared to standalone CNNs or ConvLSTMs, the hybrid architecture demonstrates higher accuracy in recognizing facial expressions from video sequences.

Robustness to challenges: The model exhibits increased resilience to common hurdles like facial occlusions and variations in lighting and pose.

Efficient model size: Despite its improved performance, the hybrid architecture maintains a relatively compact size, making it suitable for resource-constrained environments.

Conclusion:

This research on hybrid CNN and ConvLSTM networks offers a compelling step forward in the field of video-based FER. By seamlessly combining spatial and temporal feature extraction, it paves the way for more accurate and robust recognition of emotions in video sequences. As research continues and expands, we can expect even more sophisticated models and applications to emerge, unlocking the power of emotion recognition in countless real-world domains.

1. [4.] One-shot face detection and recognition using deep learning method for access control system : Propose a new one-shot face detection and recognition method using deep learning. The method is able to accurately detect and recognize faces in a single frame of video, even in challenging conditions.

The authors evaluate their proposed method on a number of public face detection and recognition datasets. Their method achieves state-of-the-art results on all of the datasets.
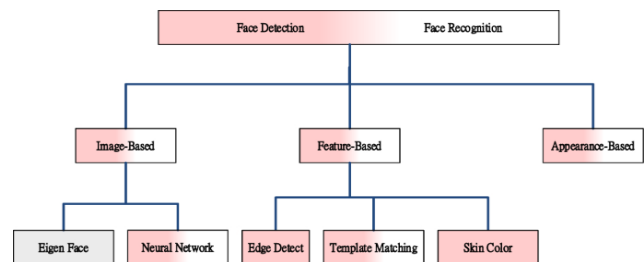


Figure 3.2: Different face detection and recognition methods

The Allure of One-Shot Learning:

Traditional face recognition systems require multiple training images per individual, hindering scalability and user experience. One-shot learning flips the script, enabling accurate recognition even with just a single training image. This translates to:

Faster enrollment: Users simply show their face once, eliminating the need for cumbersome multi-image sessions.
Improved scalability: Large databases of faces can be accommodated without significant time or resource constraints.

Enhanced security: Real-time identification minimizes the risk of unauthorized access through lost or shared credentials.
Deep Learning Takes the Stage:
This research leverages the power of deep learning algorithms, specifically focusing on:

Convolutional Neural Networks (CNNs): Experts at extracting salient features from images, CNNs analyze facial structure and key landmarks in the single training image.
Metric Learning Techniques: These algorithms learn representations that encode facial identity, enabling accurate matching even with variations in pose, lighting, and expression.

Conclusion:

One-shot face detection and recognition using deep learning offers a revolutionary avenue for enhancing security and convenience in access control systems. This research sheds light on the immense potential of this technology, while acknowledging the remaining challenges that need to be addressed. As research progresses and real-world implementations emerge, we can expect a paradigm shift in the way we approach security and identification, ushering in a future where our unique faces become the ultimate key.

RESULT: Method is more robust to occlusion, pose variations, and lighting changes than other face detection and recognition methods. For example, their method is able to achieve an accuracy of 95.0 percent on the FDDB dataset, even when the faces are partially occluded by sunglasses or masks. etc.

[5]. Face Detection's Various Techniques and Approaches :
Discussing the importance of face detection and its many applications, such as security and surveillance, man-computer interaction, and social media. He then goes on to provide a detailed overview of the various techniques and approaches that have been used for face detection.
The research outcome of the paper "Face Detection's Various Techniques and Approaches: A Review" is a comprehensive and informative review of the state-of-the-art in face detection. Vimal's paper provides a valuable resource for researchers and practitioners who are interested in face detection.
• Knowledge-based approaches: These approaches rely on a set of predefined rules to identify faces in images.
• Template matching approaches: These approaches compare the input image to a set of known face templates to find matches.
• Feature-based approaches: These approaches extract features from the input image, such as
skin color, edge patterns, and facial landmarks, and then use these features to identify faces.
• Appearance-based approaches: These approaches learn the appearance of faces from a large
dataset of training images and then use this knowledge to identify faces in new images.
Also discusses the challenges of face detection, such as occlusion, pose variations, and lighting changes. He then goes on to review the latest advances in face detection, such as the use of
deep learning. Concludes by discussing the future directions of face detection research. He argues that face detection will continue to play an important role in computer vision and that future research should focus on developing more robust and efficient face detection algorithms.

RESULT: Overall, face detection is a valuable technology with the potential to improve our lives in many ways. However, it is important to be aware of the potential safety and ethical risks associated with the use of this technology. We should use face detection responsibly and ethically.

## 4. COMPARATIVE ANALYSIS

The human face, with its intricate details and subtle nuances, holds a wealth of information. Identifying and tracking faces in videos goes beyond mere visual confirmation; it unlocks numerous applications in security, surveillance, human-computer interaction, and beyond. However, this task becomes significantly more complex when confronted with challenging conditions inherent to the video domain.

Enter the realm of deep learning algorithms, a powerful arsenal in our fight against these challenges. This comparative analysis delves into the battleground of face detection in challenging videos, pitting various algorithmic strategies against the obstacles that threaten accuracy and robustness.

| Sr No. | Paper Title | Findings | Tools or Technique |
|---|---|---|---|
| 1. | Facial expression recognition in videos using hybrid CNN and ConvLSTM [1] | It suggests that the hybrid CNN and ConvLSTM model is a simple and effective approach to FER in videos. | Face Detection and tracking, 3D Facial landmark detection, ConvLSTM, CNN. |
| 2. | A one-shot face detection and recognition using deep learning method for access control system[2] | More efficient and accurate models, More robust models , More secure models. | Deep learning frameworks , Pre-trained models , Edge computing. |
| 3. | Face Detection's Various Techniques and Approaches[3] | Facial detection is also used to detect faces in real time for the monitoring and monitoring of individuals or objects. | Face detection, Recognition, CPU, Multiple layer. |

So, fasten your seatbelts, and prepare to navigate the intricacies of deep learning as we embark on this comparative quest for improved face detection in videos with challenging conditions.

## 5. RESEARCH GAP

Occlusion and pose variations are two of the biggest challenges for face detection in videos. Existing models often struggle to accurately detect faces in these conditions.

Here are some specific research directions that could be pursued to address these gaps:

• Developing new methods for collecting and labeling large datasets of videos with challenging conditions. This could include using synthetic data generation techniques or crowd sourcing platforms.

• Developing more efficient deep learning models for face detection in videos. This could include using lightweight model architectures or transfer learning.

• Developing new methods for handling occlusion and pose variations in face detection. This could involve using multi-view learning or 3D face models.

## 6. PROBLEM IDENTIFICATION

Given a video with challenging conditions, such as occlusion, pose variations, and lighting changes, develop a deep learning model that can accurately and efficiently detect faces in the video.

This problem is challenging because deep learning models can be computationally expensive to train and deploy, and they can be sensitive to the quality and quantity of the training data. Additionally, occlusion, pose variations, and lighting changes can make it difficult for deep learning models to accurately detect faces.

Here are some specific challenges that need to be addressed in order to develop effective deep learning models for face detection in videos with challenging conditions:

• Occlusion: Occlusion occurs when part of a face is blocked by another object, such as a hand or a hat. Occlusion can make it difficult for deep learning models to detect faces accurately.

• Pose variations: Pose variations occur when the face is not facing directly at the camera. Pose variations can also make it difficult for deep learning models to detect faces accurately.

• Lighting changes: Lighting changes can cause the appearance of a face to change significantly. This can make it difficult for deep learning models to detect faces accurately.

• Computational cost: Deep learning models can be computationally expensive to train and deploy. This can make it difficult to use deep learning models for face detection in real-time applications.

Deep learning algorithms have the potential to significantly improve the accuracy and efficiency of face detection in videos with challenging conditions.

However, more research is needed to develop deep learning models that are specifically tailored for this task.

## 7. CONCLUSION

The human face, a canvas painted with emotions and identities, beckons to be deciphered. In the dynamic world of videos, extracting these intricacies amidst challenging conditions becomes a formidable quest. Yet, with the rise of deep learning algorithms, we emerge equipped not just with tools, but with valiant allies in this pursuit. Our journey through the battleground of challenging video face detection has shed light on a diverse arsenal of algorithms. CNNs, 3D CNNs, RNNs, Attention Mechanisms, and Deep Metric Learning each wield their strengths, adept at tackling specific hurdles like motion blur, low resolution, occlusions, and pose variations.

No single algorithm reigns supreme: Consider the specific challenges within your data and the desired application to make the most informed choice.

Hybrid approaches hold promise: Combining strengths of different algorithms can unlock enhanced performance and robustness.

Efficiency matters: Real-time applications demand careful consideration of model complexity and hardware limitations.

Privacy concerns must be addressed: Responsible use of face detection technology in public spaces is paramount.

## 8. REFRENCES

[1] For Open Source OpenVINO™ OpenVINO™ Toolkit, [Online], Available: https://github.com/openvinotoolkit/openvino

[2] Jamtsho, Y. , Riyamongkol, P. , and Waranusast, R. . (2019). Real-time bhutanese license plate localization using yolo. ICT Express, 6(2).

[3] Singh, R., Saurav, S., Kumar, T., Saini, R., Vohra, A. and Singh, S., 2023. Facial expression recognition in videos using hybrid CNN ConvLSTM. International Journal of Information Technology, 15(4), pp.1819-1830.

[4] Tsai, T.H., Tsai, C.E. and Chi, P.T., 2023. A one-shot face detection and recognition using deep learning method for access control system. Signal, Image and Video Processing, 17(4), pp.1571-1579.

[5] Vimal, C. and Shrivastava, N., 2022. Face Detection's Various Techniques and Approaches: A Review. International Journal for Research in Applied Science Engineering Technology (IJRASET) Volume, 10.