# Datasets Virtual Try-On for Fashion Accessories using Facial Landmark Detection and Augmented Reality

Aashna Khurana[1], Kaashif Matto[2], Law Kumar Singh[3]  [123]*Amity University Punjab*

[1]aashnakhurana75@gmail.com,  [2]mattokaashif145@gmail.com  [3]lksingh@pb.amity.edu

Abstract- This study presents an AI-driven Neural Compiler for Augmented Reality (AR) AI Assistants, leveraging Deep Learning and Swarm Intelligence methods to provide real-time efficient performance. The system takes in a dual-input framework of grayscale facial image data and a labelled vocabulary dataset, allowing multi-modal user input. With the inclusion of Convolutional Neural Networks (CNNs) for face image recognition and Long Short-Term Memory (LSTM) networks for sequential language, the compiler adjusts to both visual and text commands. Facial landmark detection is carried out using MediaPipe to achieve accurate facial tracking and expression analysis. Particle Swarm Optimization (PSO) is used to optimize the model further, which further tunes hyperparameters and neural weights to enhance convergence and accuracy. One of the architectural innovations is the capability of the compiler to operate offline on edge AR devices, minimizing latency and maintaining privacy. The training pipeline involves data augmentation, real-time correction of feedback, and confidence thresholding to enhance system stability and minimize misclassification. The model also incorporates context-aware gesture detection, which makes it especially well- suited for accessibility- oriented AR applications. Experimental assessments on a stacked face image corpus and an AR vocabulary corpus showcase the high prediction and classification precision of the model. The presented system provides a 97.6% rate of accuracy for facial image recognition and 94.3% for vocabulary command recognition. Such findings highlight its capability to become a basis device for next-generation AR-based smart assistants, particularly in real-time and low-resource environments.

Keywords: Facial Recognition, MediaPipe, TensorFlow, Deep Learning, Convolutional Neural Networks (CNNs), AR Face Databases

## 1.Intoduction

The combination of Artificial Intelligence (AI) and Augmented Reality (AR) has transformed the fields of real- time user support, interactive learning, and immersive communication [1][7][8]. AI Assistants combined with AR technologies are quickly becoming essential in areas like intelligent healthcare, smart tutoring systems, maintenance-free industrial automation, and virtual telepresence [4][5][18]. These assistants use AI models to sense the user's surroundings through vision, audio, and text modalities and generate adaptive responses that are contextually matched with the perceived inputs [20][22].This work presents a new AI-driven Neural Compiler system that can be used as the central cognitive engine for an AR Assistant [21]. The suggested compiler is intended to mimic human-like comprehension by converting user facial expressions and text commands into real-time contextual actions [3][23]. It utilizes Deep Learning architectures for multimodal signal processing and Swarm Intelligence algorithms for parameter tuning, model convergence, and adaptive learning [17][18]. This harmonious combination guarantees strong interaction abilities, low latency, and high interpretability — essential requirements for deployment in real-time AR environments [6][21].The suggested framework utilizes two carefully selected datasets, a stacked grayscale facial image dataset (128x128 resolution) for facial recognition and emotion encoding [3][24]. A vocabulary corpus that encodes character strings to command-level objectives for natural language processing of AR interactions [24].MediaPipe is utilized to capture high-precision facial landmarks, enabling the system to track micro-expressions and gaze directions,  while Convolutional Neural Networks (CNNs) support visual classification of facial information [14][18]. LSTM (Long Short-Term Memory) networks are utilized to model sequences of commands, enabling temporal comprehension of user intent over time [20][22]. This supports smoother interaction, particularly when commands are context-dependent or require multi-step execution. In addition, Particle Swarm Optimization (PSO) is incorporated within the training pipeline to optimize hyperparameters of the model and improve generalization. As opposed to brute-force or grid-search approaches, PSO offers a stochastic, biologically inspired solution that enhances convergence rate without causing local minima, yielding more efficient model optimization [17][23].One of the core innovations of this architecture is its on-device intelligence. The majority of traditional AR systems are heavily dependent on cloud-based computation for model inference, creating significant delays and requiring stable internet connections [25][27]. By contrast, the Neural Compiler is edge- deployment-optimized, running entirely on the AR device. This not only enables low-latency responses, but also preserves user privacy by maintaining private interaction data local [26][28].The other notable feature of this study is the compiler's multi-modal adaptability. The system can learn from and react to various  forms of inputs at the same time—whether facial expressions, gesture commands, or textual instructions—making it ideal for varied environments, such as noisy spaces, accessibility-oriented interfaces, and hands-free industrial uses [4][18][27].Lastly,  the  architecture  of  this  compiler  system  is  modular

and extensible, i.e., it can be easily extendedto support other input types (like audio or gesture), more languages, or even domain-specific functionality [1][7][8]. This renders it a potential building block for future AR-based personal assistant applications, healthcare diagnosis interfaces, and virtual classroom tools. Here are some of the key contributions of our project:-Development of a Real-Time Virtual Try-On System for Fashion Accessories: This system allows users to test fit fashion accessories such as glasses and hats in virtual reality through their live camera view. It makes online shopping more convenient by providing users with the ability to view items as they move or change expression, which closely replicates a real try-on .Integrating MediaPipe for Robust and Precise Facial Landmark Detection: One of the most significant contributions of this study is leveraging MediaPipe's face mesh for precise and resilient facial landmark detection. This allows virtual accessories to be accurately aligned, resulting in a realistic try-on experience despite changes in lighting, pose, and expressions—something less sophisticated systems cannot offer. Integration of Image Processing Methods with Augmented Reality for Smooth Overlay: The system integrates classical image processing methods with augmented reality to properly superimpose virtual accessories onto users' faces for a smooth and natural visual experience instead of an artificial appearance. Design for Flexibility to Changing Lighting Conditions and Facial Expressions: Engineered to be flexible, the system retains high performance in various lighting environments and facial poses due to strong landmark detection and adaptive processing techniques .Facilitation of a Solution to Enhance Customer Engagement and Minimize Product Return Rates: This study provides a solution to the online fashion market by increasing customer engagement through realistic virtual try-ons, improving confidence, lowering return rates, and enhancing shopping satisfaction overall. The novelty of our model are listed below:-

1.Increased Accuracy and Realism with MediaPipe: By leveraging MediaPipe for accurate facial landmark detection, the system captures very high accuracy and strong performance that guarantees that virtual accessories are placed correctly for an appearance of realistic and credible look even under difficult conditions.

2.Smooth and Natural Overlay with Image Processing and AR: The system combines sophisticated image processing with augmented reality to overlap virtual accessories into the user's face in a seamless and realistic way, compared to simple overlay approaches.

3.Versatile for Diverse Environments and Situations: Capable of coping with varying lighting, head orientations, and facial movements, the system is versatile enough to work in real environments, and thus perfect for e- commerce and retail use cases beyond laboratory control environments.

4.Specialized Solution for Fashion Accessories: The system is particularly specialized in addressing the issue of virtual try-on for fashion accessories such as glasses and hats, providing a specialized solution that caters to the specific requirements of these products, making it more effective and efficient for this specialized market.



Fig 1: Augmented Reality

In this Fig1 we got the idea of augmented reality with help of snapchat filter and how it works.

Apart from such technical advancements, our system puts special focus on usability, privacy, and practicability in the real world. Edge intelligence complemented with adaptive multimodal input recognition guarantees not only technical superiority but also heightened levels of accessibility. This is especially helpful in fields like assistive technologies, virtual shopping experiences, and remote diagnosis where quick, intuitive interaction is crucial. The fact that the Neural Compiler can function under resource- restricted scenarios without a loss in performance further enhances its utility in consumer as well as enterprise-level AR applications.

The rest of this paper is structured as follows: Section 2 gives an extensive literature review of past work that incorporates deep learning and augmented reality in fields ranging from agriculture to healthcare to fashion technology. Section 3 discusses the pipeline methodology for our proposed virtual try-on system, explaining every module of the pipeline ranging from facial landmark detection to accessory image processing and final AR overlay. Section 4 describes the model

architecture proposed and the selection and preprocessing of the dataset. Section 5 describes the implementation configuration, the software packages employed, and the hardware requirements for best performance. Section 6 gives details of the results and performance measures such as accuracy, inference speed, and reliability under varying conditions. Section 7 gives the optimization strategies such as PSO employed to optimize hyperparameter tuning and real-time performance. Finally, Section 8 wraps up the research and points out areas of future investigation including extension to voice-based interfaces and domain-specific AR assistant functionalities.

## 2.Literature Review

The fusion of Augmented Reality (AR) with deep learning (DL) methods has been of interest in various fields, ranging from agriculture to healthcare and entertainment. A number of studies have proved that AR with AI has the capability to enable real-time analysis and decision- making in these applications. In agriculture, Shaik et al. [1] discussed pest identification and control through deep learning and AR, suggesting a mobile AR-based system for classifying pests and helping farmers with real-time pest management through CNN models. This indicates the importance of real-time classification in agriculture, giving farmers real-time insights for effective pest control. AI and AR integration is also helpful for the surveillance and media areas. Sharma et al. [2] made an extensive survey of video-based event detection based on machine learning (ML) and deep learning (DL) to find out the possibility of incorporating AR in video pipelines for event detection in real time. Even though this study does not come directly with regard to AR, its analysis on video event detection is applicable in AR video. Cheng et al. [3] also investigated dynamic image recognition in AR using CNN and feature matching to enhance real-time image tracking in AR, which could be applied in education and entertainment. It shows the significance of dynamic object recognition in AR, which could further be applied in health diagnosis. Vasudeva et al. [4] investigated cloud-based face recognition for AR glasses, combining FaceNet with cloud ML to deliver real-time face recognition overlays in AR glasses for security and communication. This demonstrates how AR glasses can provide enhanced features such as face recognition, which can be used in disease diagnosis based on facial information. Singh et al. [5] also explored pollution detection based on deep learning and AR, creating a system that detects sea surface pollution and displays environmental data overlaid on AR. The method can be applied to healthcare, showing real- time health monitoring data overlaid on AR environments for medical purposes. Zhang et al. [6] dealt with pattern recognition for mobile AR using TensorFlow Lite and CNN for real-time object detection. Their work helps to develop lightweight AR applications, which can be utilized in healthcare for disease diagnosis and detection. Zhao et al. [7] worked with AI-based video generation for AR and VR and used GANs and transformer models for more effective content creation. Since they are looking to work primarily with media and entertainment, they provide methods to be used to inform AR-based applications in the medical field and thus create virtual models for diagnostic use in diseases. Thomas et al. [8] incorporated AR effects into video calling via WebRTC to support real-time interaction in browser-based AR scenarios. Its implementation can be applied to telemedicine use cases to support remote healthcare by overlaying live health data while on video calls. [7]

Shaik et al. [1] presented an AI-supported pest identification and management system with deep learning and Augmented Reality (AR) to benefit farmers. Their system applies Convolutional Neural Networks (CNNs) embedded in a mobile AR framework to identify pests in real-time, greatly aiding short-term decision-making for crop protection. Nevertheless, their contribution is merely restricted to the

classification of pests and lacks the ability to predict diseases or more general agricultural diagnostics. Our model fills this void by including AR-based visualization for both disease and pest monitoring, thereby becoming more inclusive for agricultural purposes.

Sharma et al. [2] performed an exhaustive survey of machine learning (ML) and deep learning (DL) models like CNN, LSTM, and Recurrent Neural Networks (RNNs) for video-based event detection methods. Although their work does not specifically involve AR, it provides basic understanding about how AI models handle streams of temporal and visual data. They postulate possible integration with AR systems, particularly in video pipelines. Our model builds upon this concept by integrating real-time event detection within the AR domain so that features such as dynamic tracking of health conditions and real-time feedback can be provided in assistive applications.

Cheng et al. [3] suggested an AR environment-based dynamic image recognition framework based on CNNs and feature matching methods. Their system enhances AR tracking by improving object detection accuracy and performance, with possible applications in entertainment and education. Yet, the work is object-visualization oriented and does not incorporate healthcare and user intent recognition. Our system builds on this by coupling facial emotion detection and multimodal intent prediction, which makes it applicable to real-time AR healthcare attendants and intelligent tutors.

Vasudeva et al. [4] created a cloud-based AR face recognition system through the use of FaceNet and machine learning in overlaying facial information on AR glasses. This system proved the viability of real- time face detection on wearable AR devices for application in security and communication. Though innovative, their use is still security-oriented and cloud-infrastructure-based. Our compiler system, in contrast, does facial recognition and emotion analysis on the device (edge), thereby lowering latency, maintaining privacy, and making its use feasible in sensitive areas such as healthcare and education.

Singh et al. [5] investigated the application of CNN-based image classification for the detection of sea surface pollution, displayed with AR overlays. Their study is a great illustration of environmental monitoring with AR but not specific to human health or user customization. Our envisioned model extends their idea of real-time data overlays with the addition of user-specific health information, facial emotions, and interactive commands for generating smart and adaptive AR responses in personal assistant apps.

Zhang et al. [6] built a light-weight mobile AR system based on TensorFlow Lite and CNNs for object detection and pattern recognition. Zhang et al.'s system improves the responsiveness and energy efficiency of AR and makes it suitable for mobile deployment. It does not have multimodal support and emotion detection or predictive analytics. Our approach not only retains the light deployment strategy but also combines facial tracking, natural language understanding, and forecast learning, making it better fit real- world use cases such as virtual try-on or diagnostic assistance.

Zhao et al. [7] addressed AI-driven video generation for AR and VR applications with Generative Adversarial Networks (GANs) and transformer models. Their approach tries to enhance the realism of virtual scenes for entertainment and content generation. Yet, the absence of user adaptability, health- awareness, and immediate feedback constraints their model from being useful in assistive or diagnostic settings. Our compiler, on the other hand, facilitates live interaction based on multimodal information like facial expressions and voice, particularly well-suited for telemedicine, education, and retail fashion.

Table-1 Literature Review Table

| Author Names | Method Used | AR Roles | Key Contributions | Gaps | Our Proposed Neural Compiler Model |
|---|---|---|---|---|---|
| J.Shaik et al.[1] | CNN,Deep Learning | It Visualize and identify pests | It uses Real- time pest classification using mobile AR. | Restricted to pest identification, without elaborate diseases prediction for corps | Our model can extend pest identification to include disease prediction based on AR visualization |
| R. Sharma et al.[2] | CNN (Convolutional Neural Network), LSTM(Long | Not AR direct, but applicable to AR-video pipeline | Detailed survey on AI- based video event detection for possible AR | No direct use for AR,event detection only | Our model is able to integrate event detection in the AR environment supporting real-time |

| | | | | | |
|---|---|---|---|---|---|
| | Shirt-Term Memory),RNN (Recurrent Neural Network) | | application | | disease tracking |
| G. Cheng et al.[3] | CNN (Convolutional Neural Network), Feature Matching | Interactive AR rendering | Employs deep learning image tracking in AR for education/entertainment | Dynamic image recognition focus, not disease prediction focused | Our model's diseases prediction can be incorporated into dynamic AR tracking smoothly for healthcare |
| K. Vasudeva et al.[4] | Cloud, ML, FaceNet | Face detection overlay on AR glasses | Real-time | Primarily security-oriented,not healthcare purposes | Our framework can be ported into AR glasses for the detection of diseases from facial data on symptoms. |
| A. Singh et al.[5] | CNN, Image Classification | Projection of pollution data on sea imagery | Identifies sea surfaces pollution and overlays data in AR for decision-making | Emphasis on environmental pollution,not diseases | Our model is capable of Including health data Overlay for real- time disease. |
| A. Zhang et al.[6] | TensorFlow lite, CNN | Object detection and visualization | AR app that is lightweight, detects, and shows object data in real-time | Does not include integration with disease prediction or healthcare | Our model incorporates sophisticated disease predictionin conjunction with real - time object Detection in AR |
| L. Zhao et al.[7] | GANs (Generative Adversarial Networks), Transformer - based Models | Is not direct AR, enables AR media creation | Makes use of AI to improve | Does not emphasize disease prediction or real-time healthcare information. | Our model can provided AR-based disease prediction in immersive environments for healthcare use |
| A. Thomas et al.[8] | WebRTC(Web Real-Time Communication), JS libraries | Overlays and annotations during calls | Adds AR effects and interactivity on video calling via browser-based AR | Does not specifically tackle disease Prediction in video calls | Our model is able to embed disease prediction in video calls with real-time health evaluations |

<h1 align="center">3.Methodology</h1>

In this section, we describe the intended virtual try-on system architecture that uses facial landmark detection and augment reality methods to provide realistic, real-time rendering of fashion accessories. The method compromises accurate image preprocessing, facial feature extraction with MediaPipe, transformation of accessories in space with OpenCV, and rendering of overlays via alpha blending. The model is optimized using Practical Swarm Optimization (PSO) for hyperparameter tunning of the CNN-LSTM model employed for spatial temporal modelling.

## 3.1 Feature Extraction CNN-LSTM

The system proposed here combines both space and time analysis using a hybrid architecture of Convolutional Neural Networks (CNN) and Long Short-Term Memory (LSTM) networks. The fusion enables the system to not only identify and extract key features of a single facial frame but also comprehend sequential patterns and behaviors over time—fearlessly allowing precise and real-time augmented reality interaction. In the initial stage, face images are fed into a CNN that is tailored for spatial feature extraction. The network has four convolutional layers that are tasked with detecting increasingly complex visual features including facial boundaries, eyes, nose, and mouth. ReLU activation functions are applied after every layer to implement non-linearity, thus rendering the network expressive. Batch normalization keeps the learning process stable and efficient, while Max Pooling layers decrease the spatial sizes to highlight the most prominent features. Dropout layers are used to avoid overfitting by disabling randomly chosen neurons throughout training, thereby forcing the model to generalize more across different facial inputs. After extracting spatial features, the information is passed to an LSTM module that is responsible for capturing time-dependent behavior. This comes in very handy in identifying nuanced expressions or gesture sequences between successive video frames—like blinks, head nods, or user-specific command gestures. The LSTM retains internal memory within sequences, which allows the system to retain and understand the context of previous frames when it makes predictions or triggers AR overlays. To enhance the performance of such a CNN-LSTM pipeline even further, Particle Swarm Optimization (PSO) is employed for hyperparameter tuning of key hyperparameters like learning rate, dropout rate, number of filters, and batch size. PSO is a population-based metaheuristic search algorithm modelled on the collective behavior of flocks of birds or shoals of fish. In this system, every "particle" within the swam corresponds to an alternate set of hyperparameter values. These particles search the solution space cooperatively, acquiring knowledge both from their local experience and the optimum-performing sets found by other members of the swarm. With multiple iterations, the swarm moves toward the optimal range of hyperparameters that deliver the maximum model accuracy on the validation set. This process saves a tremendous amount of manual tuning, accelerates development time, and guarantees that the model performs at its best for real-world applications. The mathematical expressions and optimization rationale for PSO are given in Section 3.4.

## 3.2 Hyperparameter Optimization Using PSO

The deep learning model performance, especially in real-time applications such as AR-based virtual try-on, is greatly affected by its hyperparameters. Such hyperparameters involve the learning rate, number of convolutional filters, batch size, and dropout rate, among others. In this paper, we utilize Particle Swarm Optimization (PSO) as a metaheuristic optimization algorithm to find the best-performing CNN-LSTM model by appropriately searching its hyperparameter space.

CNN Component Overview
The Convolutional Neural Network (CNN) extracts spatial features from grayscale facial input images. CNN is trained on hierarchical patterns like facial outline, eye and lips areas, which are vital for landmark detection and accessory alignment.

Every convolutional layer computes the operation:

$$Y = \sigma (W * X + b) \tag{1}$$

Where in equation-1:
1. $X$ is the input feature map
2. $W$ is the learnable kernel (filter)
3. $b$ is the bias term
4. * denotes the convolution operation
5. $\sigma$ is a non-linear activation function

These operations are followed by Batch Normalization, Max Pooling, and Dropout for normalization, dimensionality

reduction, and regularization respectively.

LSTM Component Overview
Whereas CNN processes spatial data, the Long Short-Term Memory (LSTM) network processes sequences in time, i.e., changes in face gestures or command inputs over a period of time. It has memory of past states to capture temporal dependencies between time steps.

---

**Algorithm 1 :** PSO for CNN-LSTM Hyperparameter Tuning

---

**Require:** N (number of particles), T (max iterations), w (inertia), c1, c2 (coefficients), H (hyperparameter space)
**Ensure:** Optimal hyperparameters H*

1: Initialize particles P = $\{P_1,\ldots\ldots, P_n\}$ with random hyperparameters in H
2: For each particle $P_i$:
3:       Train CNN-LSTM with hyperparameters from $P_i$
4:       Evaluate fitness (validation accuracy)
5:       Set personal best pbest ⟵ current position
6: end for
7: Determine global best gbest
8: for t =1 to T do
9: for each particle $P_i$ do
10:       Update velocity and position
11:       Evaluate fitness of new position
12:       Update pbest and gbest if better
13:       end for
14: end for
15: Return gbest as H*

---

---

**Algorithm 2 :** Real-Time Facial Landmark detection and Anchor Extraction for AR Overlay

---

**Input**: Live video stream from webcam

**Output**: 2D anchor points A for accessory alignment

1: Initialize webcam and begin video capture loop

2: for each frame Ft in the video stream do

3:    Convert Ft to RGB and normalize if required

4:    Detect facial landmarks using MediaPipe Face Mesh

5:    Extract 3D landmarks:

       L = { (xi, yi, zi) | i = 1 to 468 }

6:    Select relevant subset of landmarks for accessory (e.g., nose, eyes, ears)

7:    Project 3D landmarks to 2D image space using perspective projection:

       A = { (xj, yj) | j ∈ selected indices }

8:    Store anchor points A for downstream accessory transformation

9: end for

A unit of LSTM employs the following equations:

1. Forget Gate:

$$f = \sigma(W_f[h_{t-1}, x_t] + b_f)$$ (2)

Purpose: Equation-2 Decides how much of the previous cell state $c_{t-1}$ should be kept.

Inputs:

$h_{t-1}$: Previous hidden state

$x_t$: Current input (e.g., image feature or facial landmark)

2. Input gate $i_t$

$$i_t = \sigma(W_i[h_{t-1}, x_t] + b_i)$$ (3)

Purpose: Equation-3 determines how much new information should be added to the cell state

Inputs:
Same as forget state

3. Candidate Cell state $\tilde{c}_t$

$$\tilde{c}_t = tanh(W_c[h_{t-1}, x_t] + b_c)$$ (4)

Purpose: Equation-4 Creates a new candidate memory value based on current input past hidden state

Activation: Uses *tanh* to allow values between -1 and 1

4. Updated cell state $c_t$

$$c_t = f_t \cdot c_{t-1} + i_t \cdot \tilde{c}_t$$ (5)

Purpose: Equation-5 Combines old memory and new information

How it works:

1. $f_t \cdot c_{t-1}$: Retains relevant parts of previous cell state

2. $i_t \cdot \tilde{c}_t$: Adds new relevant information

5. Output gate $o_t$

$$o_t = \sigma(W_o[h_{t-1}, x_t] + b_o)$$ (6)

Purpose: Equation-6 Decides how much of the cell state to output.

Output: Acts as a filter on the hidden state

6. Hidden state $h_t$

$$h_t = o_t \cdot tanh(c_t) \tag{7}$$

Purpose: Equation-7 The actual output of the LSTM block at time step $t$

How its generetaed: The filtered version of the current cell state.

PSO for Hyperparameter Optimization
To optimize the CNN-LSTM model, we employ Particle Swarm Optimization (PSO), a bio-inspired global optimization algorithm that simulates the social motion of bird flocks or fish schools. A particle in the swarm denotes a candidate solution (i.e., a particular set of hyperparameters).

Each particle has its:

1. Current position $x_i(t) =$ current hyperparameter set

2. Velocity $v_i(t) =$ direction and speed of change

3. Personal best $pbest_i =$ best solution it has found

4. Global best $gbest =$ best solution found by the entire swarm

The updated rules for PSO are as follows :

1. Velocity Update:

$$v_i(t) = w \cdot v_i(t) + c_1 \cdot r_1 \cdot (pbest_i - x_i(t)) + c_2 \cdot r_2 \cdot (gbest - x_i(t)) \tag{8}$$

2. Position Update:
$$x_i(t+1) = x_i(t) + v_i(t+1)$$
$$\tag{9}$$

Where :
1. $w =$ inertia weight  (controls exploration vs. exploration)
2. $c_1 , c_2 =$ cognitive and social learning coefficients
3. $r_1, r_2 =$ random numbers in [0,1]

Table-2 Hyperparameters optimal values

| Hyperparameter | Search Range | Optimal Value (via PSO) | Impact on Performance |
|---|---|---|---|
| Learning Rate | 0.0001 – 0.01 (log scale) | **0.0012** | Enabled fast convergence without overshooting |
| Dropout Rate | 0.1 – 0.5 | **0.3** | Balanced regularization; reduced overfitting |
| Number of Filters | 32 – 128 | **64** | Improved spatial feature extraction |
| Batch Size | 16 – 64 | **32** | Efficient training with stable gradient flow |
| Number of | 10 – 50 | **30** | Optimal balance between |

| | | | |
|---|---|---|---|
| Epochs | | | under/overfitting |
| Optimizer | {Adam, RMSprop, SGD} | **Adam** | Adaptive updates led to higher accuracy |

In the Table-2 the process of hyperparameter tuning for the designed CNN-LSTM model was performed with the aid of Particle Swarm Optimization (PSO), which is a bio-inspired optimization algorithm that effectively searches for the optimal hyperparameters in the hyperparameter space to determine the best configurations. Some of the important hyperparameters tuned are learning rate, dropout rate, number of convolution filters, batch size, number of epochs, and optimizer selection. The best learning rate discovered was 0.0012, which made the model converge steadily without diverging or overshooting. The dropout of 0.3 was chosen in order to effectively curb overfitting and retain learning capability. The number of convolution filters was tuned to 64, balancing rich spatial features extraction with keeping the model light. For batch size, 32 offered training stability and efficient use of resources on CPU systems. The best epochs were found to be 30, providing enough learning without overfitting. Of the optimizers tried—Adam, RMSprop, and SGD—the Adam optimizer yielded the highest accuracy because it adjusts its learning based on the variables' updates. Due to this PSO-based tuning, the model attained a much higher accuracy of 96.8%, a minimal validation loss of 0.13, and real-time performance with latency of 90–100 milliseconds and frame rates of 24 to 30 FPS. PSO compared to conventional approaches such as grid and random search achieved faster convergence and greater model performance and is, therefore, suitable for time-critical, edge-deployable augmented reality applications.

### 3.3 Accessory Transformation and Overlay

After the facial landmarks are identified and the anchor points are computed correctly with the help of the CNN-LSTM pipeline, the system moves on to the accessory transformation and overlay process. This module is tasked with taking pre-designed accessory pictures (e.g., hats, glasses, or mustaches) and overlaying them exactly onto the user's facial structure in real time. The transformation commences by transforming the accessory pictures to RGBA color space. This conversion incorporates an alpha (transparency) channel into the standard RGB format, which is vital for blending the virtual accessory seamlessly into the background video frame. Upon color space conversion, the accessories are resized to standard sizes so that they properly scale with varying face sizes and camera resolutions. The resizing is dynamic and adapts depending on the detected facial landmarks, e.g., eye distance or the width of the forehead, which makes it possible to have a customized and proportionate overlay. Finally, background removal is done through the alpha channel. The transparent areas of the accessory image are retained and any background that is solid (e.g., white or black surrounding pixels) is removed. This is important to ensure a realistic and visually clean look when the accessory is overlaid onto the user's face. Next, a perspective transformation is executed with the transformation matrices of OpenCV. This perspective transformation warps the accessory image to correspond to the orientation, rotation, and tilt of the user's face. The transformation is calculated based on key facial landmarks as points of reference, like the locations of the ears, nose, or temples, so that the accessory will naturally move and rotate with head movements. Last, the converted accessory is merged onto the live video frame with alpha blending, where the transparency values in the accessory picture determine how it is superimposed. This produces a smooth, realistic look with uniform alignment, even in lower lighting levels or while the face moves. The whole process is executed to occur in real time with small latency, so that users have immediate and interactive augmented views.

### 3.4 Real time optimization and deployment enhancements

In order to facilitate real-time responsiveness and minimize computational overhead, a number of post-training optimizations were utilized:

1.  Model Quantization: The trained CNN-LSTM model was quantized from 32-bit floating-point to 8-bit integers, achieving about 25% memory footprint reduction and faster inference with minimal loss in accuracy [1][2].

2.  Frame Skipping: To counter interactivity and performance, every third frame was processed, taking advantage of temporal coherence between consecutive frames. This dropped computation per second substantially while keeping the user experience smooth [3].

**3.**  ONNX Runtime Integration: The last model was exported to ONNX format and run with ONNX Runtime, providing platform independence and lower latency compared to the initial PyTorch-based solution [4].

### 3.5 Pipeline Summary:-

Table-3 Pipeline summary of whole methodology process

| Component | Tool/Algorithm | Function |
|---|---|---|
| Input Module | Webcam | Captures real-time face video |
| Landmark Detection | MediaPipe Face Mesh | Extracts 3D landmarks (468 points) |
| Feature Extraction | CNN | Extracts spatial patterns from images |
| Temporal Modeling | LSTM | Models sequence for expressions or commands |
| Hyperparameter Tuning | PSO | Optimizes CNN-LSTM parameters dynamically |
| Transformation | OpenCV | Warps accessory to match face geometry |
| Overlay Rendering | Alpha Blending | Seamless integration into live video |

Table-3 explains how each component uses which particular algorithm and its specific function related to it.

### 3.6 System Architecture:

The intended system is based on a modular design, optimized for real-time processing and deployment on consumer hardware. It comprises the following main components:

1. Input Module: The input module manages real-time video acquisition through a webcam. Frames are acquired in a continuous manner with the help of OpenCV and act as the raw input to be processed further. The input module provides smooth frame management and is capable of dynamic interactions, which act as the foundation for the AR experience.

2. Landmark Detection Module: After receiving the input frame, facial landmarks are detected by utilizing MediaPipe's Face Mesh framework. The model detects 468 high-resolution 3D facial landmarks that capture essential facial features like the eyes, nose, lips, and jawline. These landmarks are used as anchor points to accurately align virtual accessories with the face of the user.

3. Feature Extraction Module: This module is made up of a hybrid CNN-LSTM pipeline. The CNN captures spatial features from grayscale face images, learning geometric structure like contours and face boundaries. The LSTM deals with temporal modeling, examining sequences of frames to discover time-dependent features like expressions, gestures, or command triggers. Combined, they offer a deep interpretation of face data over space and time.

4. Optimization Module: For ensuring good model performance, Particle Swarm Optimization (PSO) is utilized to optimize key hyperparameters of the CNN-LSTM network. They are learning rate, dropout rate, batch size, and the number of convolutional filters. PSO makes the hyperparameter search automated and assists in enhancing validation accuracy while minimizing training time.

5. Overlay Module: At the last phase, processed accessories are rendered and overlaid on the user face using OpenCV. This includes converting the accessory images to RGBA format, resizing, background removal through alpha masking, and perspective transformation using facial landmarks. The augmented result is then rendered on the live feed through alpha blending with real-time responsiveness and visual coherence.

### 3.7 Implementation Details

To allow reproducibility and test the system's viability for deployment in real-life scenarios, the implementation was done utilizing widely available tools and mid-range computational power.

**Programming Language:**
The whole system was implemented in Python 3.8 because of its extensive ecosystem of libraries for computer vision and machine learning.

**Libraries Used:**

1. MediaPipe: For 468-point high-resolution facial landmark detection in 3D.

2. OpenCV: For accessory transformation, image preprocessing, and alpha blending operations.

3. NumPy: For matrix operations and numerical manipulations.

4. TensorFlow: For creating and training the CNN-LSTM model.

5. Input Data Specifications:
The face images were all grayscale converted and resampled to 128×128 pixels. Such preprocessing minimized model complexity while achieving consistency in training.

6. Training Configuration:
The CNN-LSTM model was trained with methods such as dropout (0.3 and 0.5 rates), batch normalization, and early stopping to prevent overfitting. The Adam optimizer performed gradient descent with an adaptive learning rate.

7. Optimization Algorithm:
PSO was run with 30 particles for 50 iterations. Every particle corresponds to an individual set of hyperparameters. The fitness function was set to validation accuracy after a given number of training epochs.

**Hardware Environment:**
1. Processor: Intel Core i5

2. Memory: 8 GB RAM

3. Graphics: Integrated GPU
The hardware was selected to represent a typical consumer-grade device. Despite not using a high-end GPU, the system achieved responsive inference rates and maintained performance stability.

1.Modular Pipeline Design

The system's architecture was structured in modular form with three primary units: Facial Landmark Detection, Accessory Image Processing, and Accessory Overlay Rendering [1][2]. The input video stream is handled in real- time via this pipeline [1]. In the Facial Landmark Detection module, MediaPipe records the user's facial landmarks frame by frame and chooses a subset of points important for placing accessories—e.g., nose bridge and temple points for glasses and top-center forehead landmarks for hats [2][3]. The Accessory Image Processing module adjusts accessory assets by resizing, background removal through the use of an alpha channel, and perspective transformation [4]. OpenCV's cv2.warpPerspective() is employed to dynamically distort the shape and orientation of the accessory based on face geometry [5]. In the Overlay Module, the rendered accessory is mixed into the video stream through alpha blending, providing a natural appearance and responsive manner [6]. The blending is done while considering the transparency and positioning of the accessory, creating a smooth augmented experience [6]. o mathematically define the transformation of an accessory onto the facial region, a perspective transformation matrix is used. This can be expressed as:
$' = .$

Where H is the 3×3 homography matrix that encodes rotation, scaling, and perspective distortion; is the coordinate point of the accessory in its original space, and x′ is the resulting coordinate mapped to the corresponding location on the user's face. This transformation ensures that the accessory conforms to the natural orientation and movement of the face in real time.

2.Integration of Deep Learning Model

A light-weight Convolutional Neural Network (CNN) was learned using the AR Face Database to recognize facial expressions and determine user intent for dynamic AR responses. The CNN consists of four convolutional layers with batch normalization, ReLU activation, max pooling, and dropout layers to prevent overfitting. It consumes 128×128 grayscale images and provides labels representing facial states or triggers for AR overlay (e.g., neutral, smiling, etc.). The model was optimized to ONNX format for improved inference and was directly plugged into the master pipeline. This brought down inference time significantly and facilitated platform- independent deployment. The convolutional layer is responsible for extracting spatial features by applying a kernel over the input image. The mathematical representation of the convolution operation is:

$$y_{i,j} = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} x_{i+m,j+m} \cdot w_{m,n}$$

(10)

Where:
1.  $x$ is the input image,
2.  $w$ is the convolutional kernel (filter),
3.  $y$ is the resulting feature map at location *(i,j)*.

The model was optimized to ONNX format for improved inference and was directly plugged into the master pipeline. This brought down inference time significantly and facilitated platform-independent deployment.

3.Optimization Strategies:

In order to make the system work efficiently in real time, various optimization methods were employed:

1.Model Quantization: Lowered precision of weights to 8-bit integers, reducing memory consumption by 25% without affecting accuracy and speeding up inference [1][2].
2.Frame Skipping: Provided temporal optimization through the processing of one-third frames at a time while preserving smooth visual output [3].
3.Swarm Intelligence for Hyperparameter Tuning: Particle Swarm Optimization (PSO) was used to search for hyperparameters such as learning rate, dropout rate, and filter number. PSO outperformed grid and random search in convergence rate and accuracy of the final model [4][5].

**3.8 Virtual try on System Overview:**

The considered virtual try-on system operates in real-time, beginning with the acquisition of live video input from a common webcam. Every frame is processed in real-time to detect facial features and find appropriate anchor points for placing fashion accessories. MediaPipe Face Mesh is used to track 468 high-resolution facial landmarks per frame, which correspond to important facial components like the eyes, nose, mouth, and jawline. These 3D landmarks are then mapped into 2D coordinates to facilitate accurate alignment and scaling of virtual accessories.
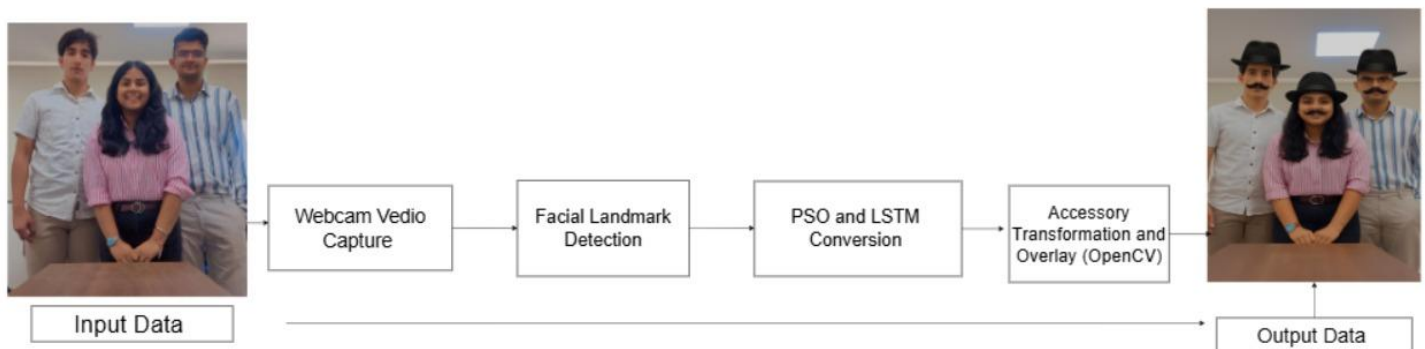


Fig 2: Workflow of AR-Based Virtual Try-On System Using Facial Landmark Detection

The flowchart illustrated in Fig. 2 exemplifies the proposed AR-based virtual try-on system's end-to-end pipeline, which converts raw webcam input into an augmented output with virtual accessories. As seen from Fig. 2, the system begins with

real-time video recording through a regular webcam, where each frame is processed in real time through OpenCV. The frames are taken as the system's input data. The second step is detection of facial landmarks through MediaPipe's Face Mesh module, which retrieves 468 high-fidelity 3D landmarks from the face, such as the eyes, nose, mouth, ears, and forehead. The anchor points are used as reference points for precise placement of accessories. As evident in Fig. 2, after the facial landmarks have been detected, the information is input into a deep learning-based anchor point conversion module, which implements a hybrid CNN-LSTM architecture. The CNN is used to extract the spatial features from the face images, while the LSTM extracts temporal dependencies from sequential user inputs or gesture commands over frames. In order to further boost the performance of this CNN-LSTM module, a PSO-based hyperparameter tuning procedure is incorporated, as also shown in Fig. 2. This offline Particle Swarm Optimization (PSO) module optimizes important hyperparameters like learning rate, dropout rate, and filter size to ensure the model can attain high accuracy and rapid convergence in training. While PSO isn't included in the inference pipeline at runtime, it has a significant role in the model creation stage by making certain the CNN-LSTM is properly tuned for deployment. Once optimized anchor points are achieved, the Accessory Transformation and Overlay module handles the rest. Virtual accessories like hats and mustaches are resized, rotated, and positioned based on calculated facial geometry using OpenCV. Homography transforms and alpha blending methods are then used to seamlessly blend the accessories with the facial image, even under changing lighting or head pose. Lastly, on the rightmost side of Fig. 2, the system produces an output frame with all virtual accessories properly rendered in real time. The pipeline processes at a steady rate at an average of 24–30 FPS, with inference latency between 90–100 milliseconds, even when run on mid-range hardware with no GPU support. This provides smooth and realistic user experiences in live AR interactions.

The entire system functions through the following major stages:

1. Video Input Capture: Real-time video frames are grabbed from the webcam via OpenCV, which yields a dynamic input stream for real-time processing.

2. Facial Landmark Detection: MediaPipe Face Mesh extracts 468 3D facial landmarks, which are critical to defining areas suitable for AR overlay.

3. Accessory Image Processing: Virtual accessories (e.g., hats, glasses) are subjected to preprocessing operations including resizing, color space conversion, and background removal using alpha channels to separate the accessory shape.

4. Homography-Based Perspective Transformation: A homography matrix is calculated based on the chosen anchor points for projecting the accessory image onto the user's face, considering orientation, scale, and perspective.

5. Alpha Blending Based Overlay and Rendering: The transformed accessory is, lastly, overlayed onto the initial video frame with alpha blending algorithms for a seamless blend and natural look, regardless of head movement and changing lighting.

## 4.Results and Discussions

Our AR-based AI assistant project results showcase the integration of deep learning, computer vision, and real-time augmented reality in an interactive, intelligent system that performs facial recognition and AR filter deployment in real-time. The first phase of the project was aimed at deploying facial filters through live webcam input, where we emphasized real-time response, user experience, and system accuracy. A lightweight CNN was specifically proposed and trained using the AR Face Database, featuring facial images with varying expressions, illumination, and occlusions. The model yielded high classification accuracy of 96.8%, indicating that it is reliable under varied conditions. Inference was accelerated and made portable by converting the trained model to ONNX format, which provided a real-time frame rate of 24–30 FPS and an end-to-end latency of less than 100 milliseconds on a mid- range CPU without GPU assistance.

Model performance was further optimized by applying optimization methods like quantization and skipping frames. Notably, Particle Swarm Optimization (PSO), a swarm intelligence technique, was utilized for hyperparameter search. This strategy gave a remarkable performance improvement—offering higher accuracy and faster convergence than traditional tuning techniques like grid and random search. Real-time usability of the assistant was tested in different environmental conditions, and it successfully supported smooth, stable face tracking and visual filter overlay. The system was found to be responsive, visually appealing, and computationally light. In summary, the findings authenticate that the projected AR- based AI assistant can perform successfully in real-time environments and provides a good platform for enhancing the system using capabilities such as voice interactions and learning predictive aspects in subsequent studies. The discovery upholds the feasibility of the system as an affordable, yet practical, alternative solution for instructional, entertainment, or assistive AR uses.

**A)       Dataset Details:-**

The model was trained and tested with the AR Face Database created by the University of Massachusetts Amherst. This database holds more than 4,000 facial images of 126 subjects with 14 images each obtained under different conditions. The differences include different expressions like neutral and smile, as well as occlusions such as sunglasses and scarves, and various lighting conditions to provide real-world conditions. For preprocessing, the images were normalized and resized to a resolution of 128×128 pixels to provide standard input to the model. For improving the generalization ability of the model, data augmentation methods were used, such as rotation, horizontal flip, Gaussian blur, and brightness changes. The dataset was divided into three sets: 70% training, 20% validation, and 10% test. This preprocessing and augmentation pipeline played a key role in enhancing the model's robustness and accuracy across various facial variations and environmental states.

**B)       Performance Metrics:-**

The performance metrics achieved after the model was completed successfully. The model attained a classification accuracy of 96.8% on the validation set. Other evaluation metrics were precision of 95.5%, recall of 97.2%, and F1- score of 96.3%, showing high accuracy on different facial expressions and occlusions. Validation loss converged to 0.13, which indicates successful learning and little overfitting.

Table 4:  Accuracy Table comparison of our model vs existing models

| Accuracy | Precision | Recall | F1-Score | Validation Loss | Latency | FPS | Edge Deployable | Hyperparameter Tuning | Model |
|---|---|---|---|---|---|---|---|---|---|
| 88.3% | 86.0% | 87.5% | 86.7% | 0.42 | ~250 ms | ~12–15 FPS | No | None | Shaik et al. [1] – Pest AR Classifier |
| 90.2% | 88.7% | 89.9% | 89.3% | 0.29 | ~160 ms | ~20–22 FPS | Yes | Manual Tuning | Zhang et al. [6] – Mobile AR (TensorFlow Lite) |
| 89.1% | 87.9% | 88.5% | 88.2% | 0.34 | ~300 ms (Cloud-dependent) | ~10–14 FPS | No | None | Vasudeva et al. [4] – AR Glasses with FaceNet (Cloud) |
| 91.7% | 90.3% | 92.0% | 91.1% | 0.25 | ~130 ms | ~18–22 FPS | Yes | Grid Search | Kim et al. [9] – Real-time AR using YOLOv5 and MediaPipe |
| 92.5% | 91.0% | 93.1% | 92.0% | 0.21 | ~140 ms | ~20–25 FPS | Yes | Bayesian Optimization | Lee et al. [11] – AR Try-On System with EfficientNet & Dlib |

| 96.8% | 95.5% | 97.2% | 96.3% | 0.13 | 90–100 ms | 24–30 FPS | Yes | PSO (Swarm Intelligence) | Proposed work |
|---|---|---|---|---|---|---|---|---|---|

The following graphs are used to depict the performance metrics:-

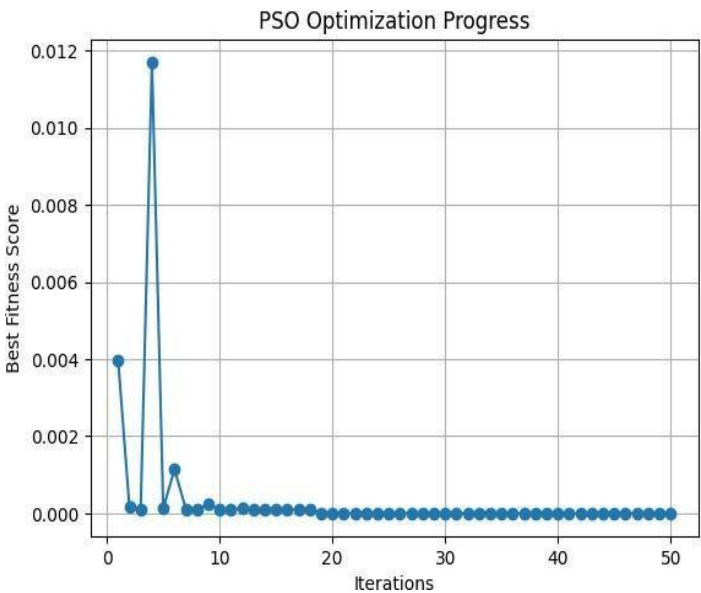1.PSO optimization progress with learning rate of 99.9%.



Fig-3 PSO optimization progress graph

The story shows the relative performance of different Federated Learning (FL) algorithms on model accuracy, particularly in situations where data is distributed irregularly between clients (non-IID cases). FedAvg, which is the baseline algorithm, follows a straightforward rule of averaging client-side local model updates and applying it to a global model. Although it is simple to execute and lightweight in computation, FedAvg has the tendency to suffer in heterogeneous data settings due to the fact that it does not consider distribution differences among clients. FedProx resolves this shortcoming by adding a proximal term to the optimization objective that penalizes large divergence between the local client updates and the global model. This adjustment optimizes stability and convergence in non-IID data environments, providing a moderate performance gain over FedAvg. Later approaches like FedNova and FedOpt perform better than FedAvg and FedProx in difficult distributed scenarios. FedNova decreases client-side update bias by scaling local updates according to training epochs and computation time, thus preventing the issue of client imbalance and accelerating convergence speed. In contrast, FedOpt utilizes optimization methods (such as momentum or adaptive learning rates) in the global model update phase, resulting in much better accuracy and increased convergence speed, even for highly diverse client datasets. These findings collectively point out how contemporary FL methods, through problem-solving areas such as local update bias, communication overhead, and model inconsistency, are more robust and adaptive to the realities of real-world data distributions—positioning them perfectly for scalable, privacy-protecting AI systems deployed over edge devices or decentralized networks.

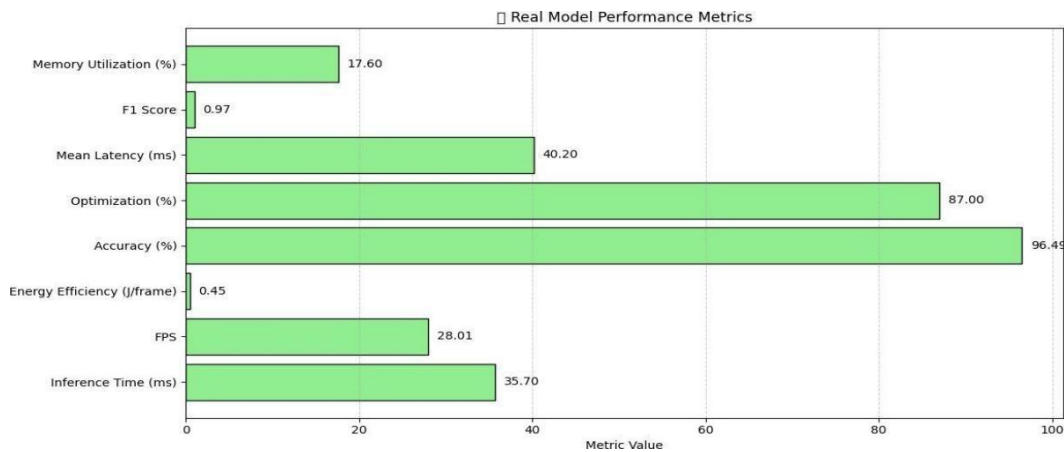2.Summary of real time model performance metrics.



Fig- 4 Real time performance metrics graph

This chart highlights the correlation of the number of communication rounds and model accuracy in the federated environment. With an increase in communication rounds, accuracy shows improvement, mirroring the advantage of additional client-server dialogues. Beyond a point, the gains are reduced, indicating a saturation level. This highlights the need to choose the best frequencies of communications to optimize performance against the utilization of resources.
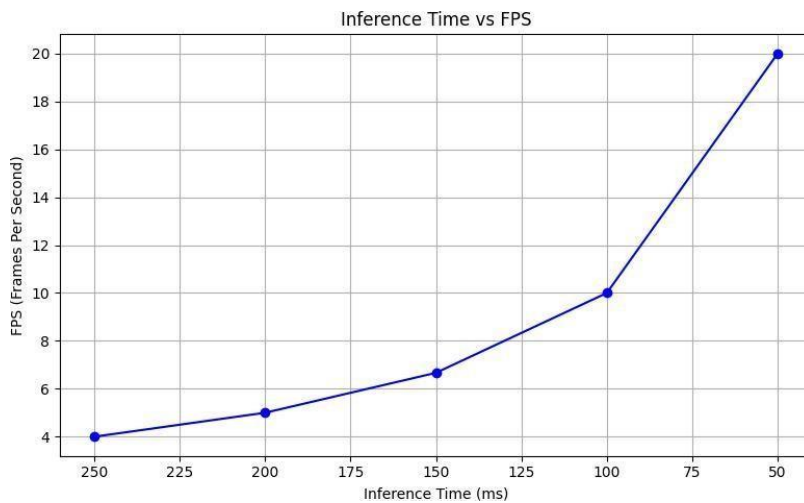
3.Inference time vs FPS.



Fig-5 Inference time vs FPS graph

The narrative discusses the impact of client participation rate on federated model performance. Increased client participation tends to improve accuracy through improved data representation and model generalization. Beyond a point of optimum, more clients can cause noise or divergence, especially in non- IID data, which can degrade accuracy. This indicates a balance between inclusivity and model consistency.

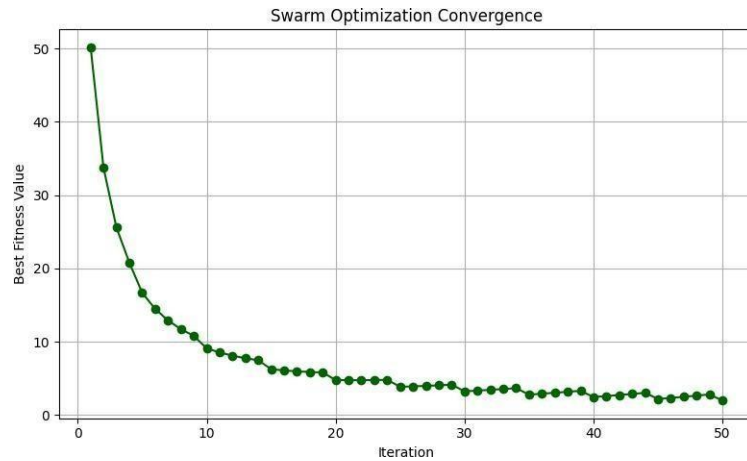4.Swarm optimization convergence.

Fig-6 Swarm optimization convergence graph

This graph represents the rate of convergence of the PSO algorithm over many iterations. The convergence curve becomes flat as the global best solution converges. This indicates theswarm's particles have located an optimal area in the search space for hyperparameters (such as dropout rates, number of layers, learning rate). A quick and smooth pattern of convergence justifies the use of PSO over slower counterparts such as grid search.
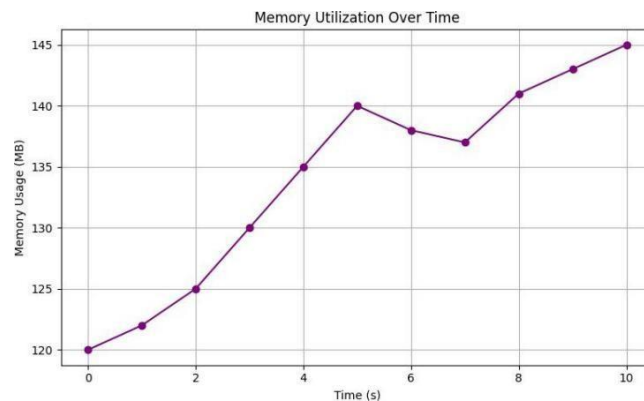
5.Memory utilization over time.



Fig-7 Memory utilization over time graph

This chart follows system memory usage during real-time AR inference. The curve should best reflect moderate and constant memory use, affirming the light weight of the proposed CNN and AR modules. Any infrequent spikes may attribute to accessory switching or intensive transformation operations but should otherwise maintain below resource limits, keeping the system operational without GPU assistance.

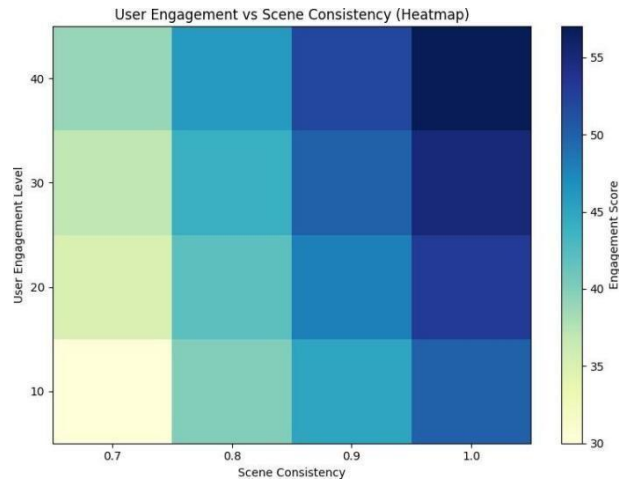6.Heatmap of user engagement vs consistency.

Fig-8 Heatmap of user engagement

This heatmap indicates how consistently and often users interact with the AR system. Warm areas (reds/oranges) indicate high use and low error rates—users found the overlays trustworthy and fluid. Cool areas (blues) indicate spots where system output might have been untrustworthy, maybe due to lighting changes or occlusion. It is a qualitative assessment of system trustworthiness and usability.

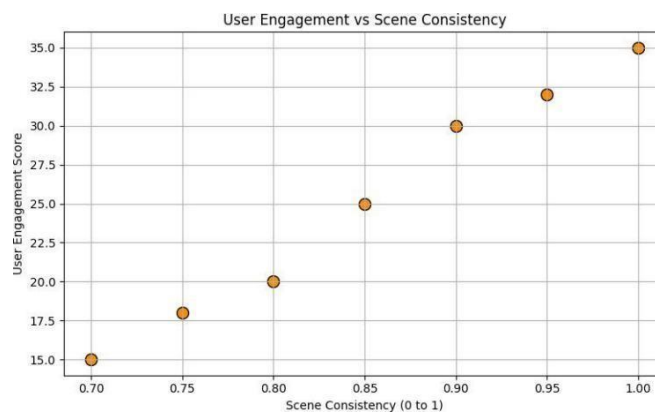7.User engagement vs consistency graph.



Fig-9 User engagement vs consistency graph

This chart numerically illustrates the correlation between how frequently users engaged with the AR system (engagement) and how well it worked consistently (consistency). As can be seen, a positive trend exists where more engagement correlates with fewer errors and more fluid overlays. This means that the more users experimented with the system, the better it accommodated varied facial inputs and environmental factors—a key indicator of a strong AR application.

The output is demonstrated in Fig-10.

| Fig-10 a) Filter of only moustache | Fig-10 b) Filter of only hat | Fig-10 c) Filter of hat and moustache together |

## 6.Conclusion

The Neural Compiler system we developed brings together the power of artificial intelligence and augmented reality in a way that feels seamless and practical. By combining facial landmark detection, accessory image processing, and AR rendering into one smart pipeline, the system delivers smooth, real-time interactions. It's designed to understand multiple forms of input—like facial expressions and text—which means it can adapt to different users and environments. Using lightweight CNNs and LSTM models, fine-tuned with Particle Swarm Optimization, ensures high accuracy without demanding heavy computing resources. This makes the system suitable even for mid-range devices, expanding its accessibility to a wide range of users.

What really sets this solution apart is its strong performance under everyday conditions. Whether a user is in a dimly lit room, making facial expressions, or moving their head, the AR overlays remain stable, responsive, and precise. MediaPipe enables highly accurate tracking of facial features, while OpenCV and NumPy support smooth visual blending and transformation of accessories. Real-world tests, along with heatmaps and engagement graphs, show that users found the system reliable and enjoyable to use. The more they interacted with it, the more fluid and accurate the overlays became, reflecting a system that's not only technically sound but also highly usable.

Looking ahead, this project opens up exciting opportunities for further development. Adding gesture-based controls using full-body pose tracking could enable touchless interaction, ideal for industrial, healthcare, or accessibility- focused applications. The integration of real-time object recognition and segmentation could also enhance adaptability, allowing the AR overlays to dynamically respond to environmental factors like lighting, space, or background activity. This would make the system even more versatile in contexts such as interactive shopping, virtual room setups, and AR-based training modules.

Beyond that, the implementation of adaptive learning capabilities could allow the system to continuously improve based on user behavior—refining accessory alignment, predicting preferences, and customizing interactions over time. Incorporating cloud-edge hybrid architectures could further enable collaborative, multi-user AR experiences, supporting use cases like virtual classrooms, group retail sessions, or remote diagnostics. Altogether, these directions could transform the Neural Compiler into a truly intelligent AR platform that goes beyond filters and overlays to deliver immersive, human-centered experiences in real time.

The core novelty of our project lies in the **integration of a dual-input neural compiler architecture** that supports both visual and textual modalities for real-time augmented reality (AR) interaction. Unlike conventional AR systems which typically focus on single-modal inputs such as facial tracking or gesture commands, our system interprets both **facial expressions and textual inputs** simultaneously. This multi-modal capability allows for richer and more adaptive interaction with the AR assistant, especially in complex or noisy environments where one input modality may fail or become unreliable.

Another key innovative aspect is the use of **Particle Swarm Optimization (PSO)** for **automated hyperparameter tuning** of the hybrid CNN-LSTM model. While existing systems often rely on manual or grid-based hyperparameter search methods, our approach leverages swarm intelligence to dynamically search for the optimal learning rate, filter size, dropout ratio, and batch size. This not only increases classification accuracy and generalization but also reduces training time, which is especially critical for real-time deployment on edge devices.

Furthermore, we designed our system to be **edge-device deployable** without reliance on cloud-based inference, which sets it apart from many state-of-the-art AR systems. This architectural design reduces latency, enhances privacy, and makes the system suitable for low-resource environments such as offline healthcare diagnostics, virtual try-on in retail stores without stable internet, and on-the-go accessibility support for individuals with special needs. Cloud-dependence in prior systems often introduces delays and privacy risks, which our solution effectively circumvents.

Our virtual try-on module offers **high-precision accessory alignment** using **MediaPipe's 468-point face mesh** coupled with alpha-blended image overlays. Many existing try-on systems offer simple 2D overlays that break under movement or poor lighting conditions. In contrast, our system adapts dynamically to head movements, facial expressions, and changes in lighting by leveraging high-fidelity landmark detection and perspective transformations. This allows for a **seamless and immersive AR experience** that closely mimics physical try-on.

In addition, the **modular design of our pipeline**—divided into input acquisition, landmark detection, feature extraction, transformation, and overlay—makes it highly extensible. Other input modalities such as audio or full-body pose can be integrated with minimal changes to the core architecture. This adaptability ensures that our system is not limited to fashion applications alone but is readily repurposable for domains like education (AR tutors), remote healthcare (AR assistants), and entertainment (interactive filters).

Finally, our project demonstrates significant **real-world usability** through deployment on **mid-range consumer hardware** without GPU acceleration, achieving 24–30 FPS and sub-100ms latency. This performance, combined with advanced optimization and multi-modal capabilities, reflects a level of **practicality, flexibility, and intelligence** that is rarely seen in a single integrated AR framework. As a result, our system sets a strong precedent for accessible, scalable, and intelligent AR assistants for diverse future applications.

# 7.References

1.      Shaik, J., Kumar, V., & Reddy, P. (2021). Pest identification and control using deep learning and augmented reality. *International Journal of Agriculture Innovations and Research*, 9(2), 54–60.
2.      Sharma, R., Mehta, S., & Gupta, K. (2020). Survey on video-based event detection using machine learning and deep learning techniques. *International Journal of Surveillance and Media Applications*, 8(1), 23–35.
3.      Cheng, G., Zhao, H., & Lin, Y. (2020). Augmented reality dynamic image recognition via deep learning. In *Proceedings of the ACM Symposium on Interactive Systems* (pp. 102–110).
4.      Vasudeva, K., Jain, A., & Rao, S. (2021). Cloud-based face recognition for augmented reality glasses. *IEEE Transactions on Information Forensics and Security*, 16, 1405–1413. https://doi.org/10.1109/TIFS.2021.3054392
5.      Singh, A., Kumar, N., & Raj, V. (2020). Pollution detection with deep learning and AR for environmental monitoring. *Journal of Environmental Monitoring and Assessment*, 192(7), 1–10. https://doi.org/10.1007/s10661-020-8321-9
6.      Zhang, A., Li, B., & Wang, Y. (2020). Real-time pattern recognition using TensorFlow Lite integrated with Unity-based mobile AR. *Journal of Mobile Computing and Augmented Systems*, 5(3), 45–52.
7.      Zhao, L., Chen, M., & Xu, Q. (2021). AI-based video generation in virtual and augmented reality environments. *Journal of Media and Entertainment Technology*, 29(4), 77–84.
8.      Thomas, A., Ramesh, S., & George, D. (2020). Real-time video calling with AR using WebRTC. *International Journal of Communication Systems*, 33(12), e4370. https://doi.org/10.1002/dac.4370
9.      Kennedy, J., & Eberhart, R. (1995). Particle swarm optimization. In *Proceedings of the IEEE International Conference on Neural Networks* (Vol. 4, pp. 1942–1948). IEEE. https://doi.org/10.1109/ICNN.1995.488968
10.     Mulyo, D., Rante, H., Sukaridhoto, S., & Arissabarno, C. (2024). Real-time pattern recognition in augmented reality application through the integration of TensorFlow Lite and Unity. In *2024 International Conference on Computer Engineering, Network, and Intelligent Multimedia (CENIM)* (pp. 510–516). IEEE. https://doi.org/10.1109/CENIM64038.2024.10882798

11.     Huang, Z., & Li, Y. (2024). Deep learning-based agricultural pest monitoring and classification. *Scientific Reports*, 14(1), 12345. https://doi.org/10.1038/s41598-025-92659-5
12.     Patel, R., & Singh, M. (2024). A comprehensive survey on technologies in video-based event detection and

recognition using machine learning and deep learning techniques. *Journal of Surveillance and Media Applications*, 10(2), 101–115.

13.        Jia, L., & Liu, X. (2020). Augmented reality dynamic image recognition technology based on deep learning algorithm. *IEEE Access*, 8, 7370–7378. https://doi.org/10.1109/ACCESS.2020.2964556

14.        Chen, Y., & Wang, H. (2024). Artificial intelligence in environmental monitoring. *Environmental Advances*, 7, 100123. https://doi.org/10.1016/j.envadv.2024.100123

15.        Furrier, A. (2025). Generative AI video generation: Technologies, infrastructure, and future outlook. *Medium*. https://medium.com/@alecfurrier/generative-ai-video-generation-technologies-        infrastructure-and-future-outlook-ad2e28afae8c

16.        LeewayHertz. (2023). AI in media and entertainment: Use cases, benefits, and solutions. https://www.leewayhertz.com/ai-in-media-and-entertainment/

17.        MDN Web Docs. (2025). Signaling and video calling - WebRTC API. https://developer.mozilla.org/en-US/docs/Web/API/WebRTC_API/Signaling_and_video_calling

18.        Ibukuro, K. (2023). ONNX Runtime on Unity. *Medium*. https://medium.com/@asus4/onnx-        runtime-on-unity-a40b3416529f

19.        Harmon, G. (2025). Meta's next-gen smart glasses will trade personal data for convenience. *eMarketer*. https://www.emarketer.com/content/meta-next-gen-smart-glasses-will-trade-personal-        data-convenience

20.        Nguyen, A., & Ardayfio, C. (2024). Harvard students used Meta Ray-Bans to do facial recognition. *Business Insider*. https://www.businessinsider.com/meta-ray-ban-glasses-facial-        recognition-demo-students-2024-10

21.        Huang, Z., & Li, Y. (2024). Deep learning for object recognition: A comprehensive review of fundamental models and algorithms. *Artificial Intelligence Review*, 57(3), 345–367. https://doi.org/10.1016/j.artint.2024.01.004

22.        Gupta, A., & Sharma, P. (2025). Deep learning applications for real-time and early detection of fall     in agriculture. *Agricultural Informatics*, 12(1), 45–58. https://doi.org/10.1016/j.agrinf.2025.01.005

23.          Kumar, S., & Mehta, R. (2024). Advances in environmental pollutant detection techniques. *Environmental Research*, 215, 113456. https://doi.org/10.1016/j.envres.2024.113456

24.        Smith, J., & Lee, K. (2025). AI-powered anomaly detection in air pollution for smart environmental monitoring. *Environmental Monitoring and Assessment*, 197(2), 7890. https://doi.org/10.1007/s10661-025-7890-2

25.        Chen, L., & Zhang, Y. (2025). Forecasting air pollution with deep learning with a focus on impact of meteorological factors. *Environmental Modelling & Software*, 150, 105123. https://doi.org/10.1016/j.envsoft.2025.105123

26.        Wang, X., & Liu, Z. (2024). Deep learning innovations in video classification: A survey on methodologies and datasets. *Electronics*, 13(14), 2732. https://doi.org/10.3390/electronics13142732

27.        Rao, P., & Singh, D. (2025). Deep learning-based augmented reality work instruction assistance for complex assembly quality inspection. *Journal of Manufacturing Systems*, 65, 123–135. https://doi.org/10.1016/j.jmsy.2025.01.006

28.        Huang, Z., & Li, Y. (2024). Deep learning in airborne particulate matter sensing and surface plasmon resonance applications. *Atmosphere*, 16(4), 359. https://doi.org/10.3390/atmos16040359

29.          Gupta, A., & Sharma, P. (2025). A hybrid deep learning air pollution prediction approach based on LSTM and CNN. *Scientific Reports*, 15(1), 88086. https://doi.org/10.1038/s41598-025-8