# Assignment -1 Report

Advanced Computer Vision (CS-673)

**Indian Institute of Technology Mandi**

# Submitted by

-Nikhil Kumar (S24118)

# Q1. Data source: Any standard single Face image-based dataset

**Dataset Overview**

- **Face Images:** Collected from public sources, including celebrity portraits (e.g., Narendra Modi, Bill Gates, Donald Trump) and standard face-image datasets. Organized into separate folders for training, validation, and testing. Each CSV (Train.csv, Val.csv, Test.csv) lists the file paths and labels.
- **Non-Face Images:** Sourced from object-centric datasets (e.g., cars, aeroplanes, ships). The number of non-face samples matches the face set. These images serve two purposes: training the face vs. non-face classifier and negative mining in detection.

| Split | # Face Samples | # Non-Face Samples | Total Images |
|---|---|---|---|
| Train | 690 | 690 | 1380 |
| Validation | 148 | 148 | 296 |
| Test | 148 | 148 | 296 |

## I. Implement the Histogram of Local Binary Paterns (LBP) from screech using python. Read images from a directory Train/Val/Test and generate Train.csv/Val.csv/Test.csv

- Local Binary Patterns (LBP) is a texture descriptor that encodes local spatial patterns by thresholding neighborhood pixels around each pixel.
- we uses Uniform LBP (with radius = 3 and 24 sampling points) to compute compact and rotation-invariant texture features.
- For each grayscale image, we compute the LBP map and extract a 26-bin normalized histogram representing the frequency of each uniform pattern.
- Images are organized into Train, Validation, and Test folders, each containing face and non-face subfolders.

- The script reads all images from these folders, computes LBP features, and assigns the corresponding class label.
- Extracted features are stored in CSV files (lbp_train.csv, lbp_val.csv, lbp_test.csv) with one row per image, containing 26 LBP features and a class label.

**II. Use KNN classifier or One-vs-Rest SVM classifier (choose best distance metrics and K in KNN and best kernel and hyperparameters in SVM case using Val.csv) report the accuracy on Test.csv for best case of each.**
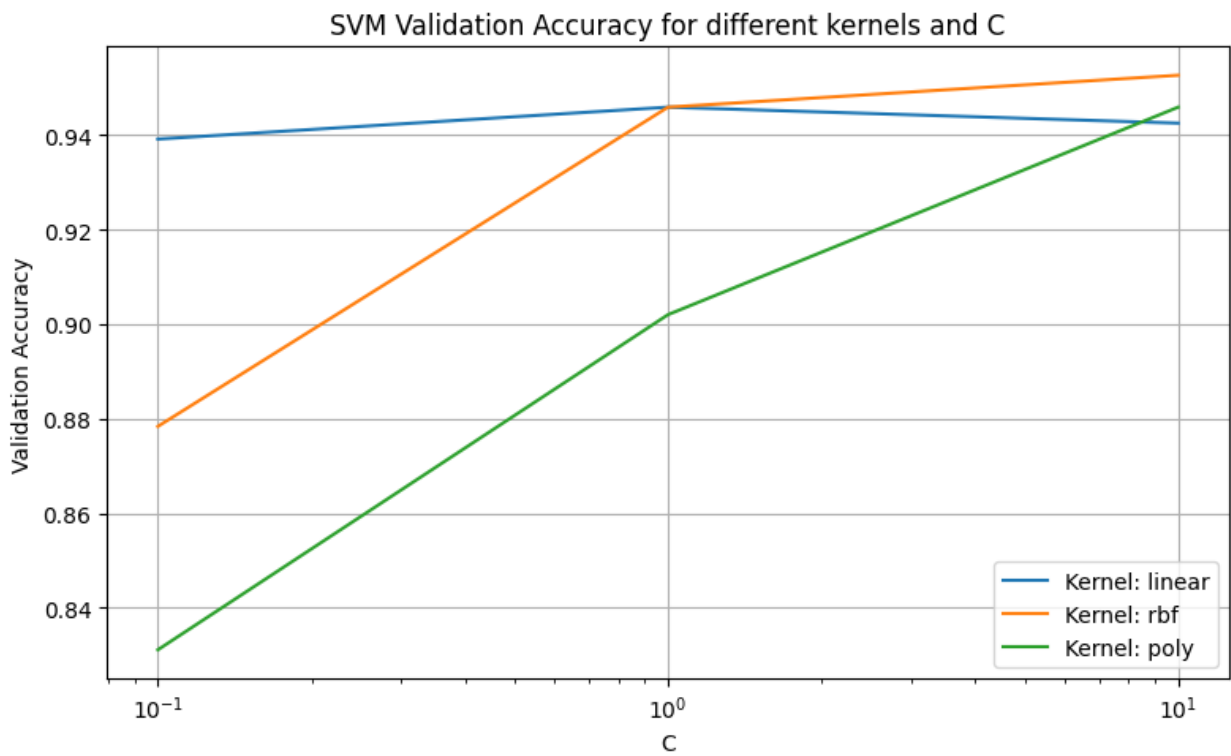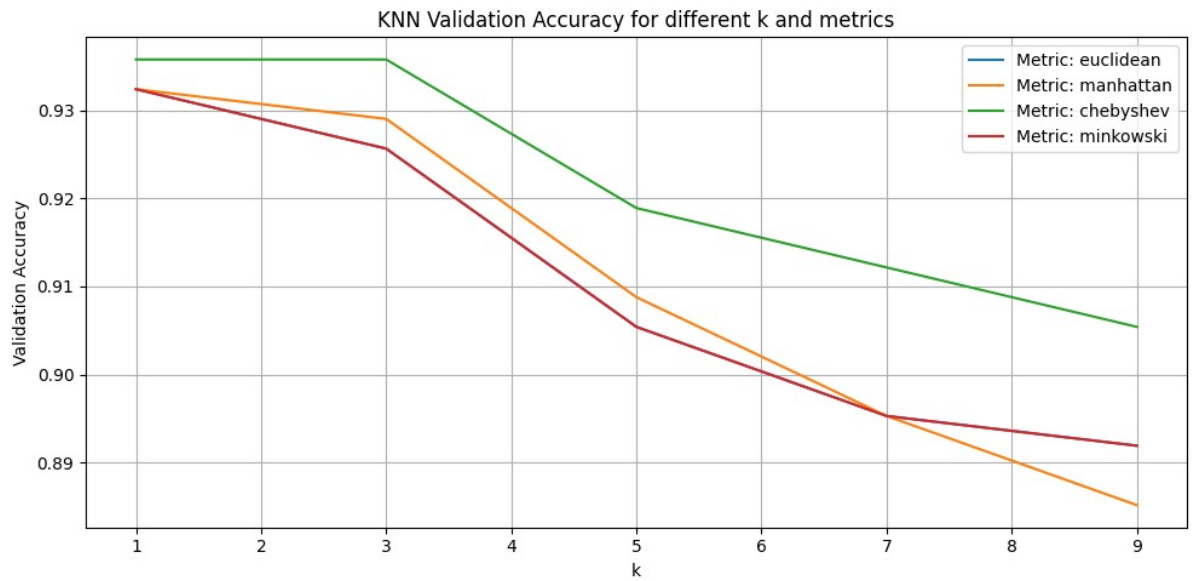
=> KNN Classifier:

- Hyperparameters tuned: number of neighbors (k ∈ {1, 3, 5, 7, 9}) and distance metric (Euclidean, Manhattan, Chebyshev, Minkowski).
- Best validation accuracy achieved using:
    1. k = 1
    2. metric = Chebyshev
- Final KNN model retrained on Train+Val and tested on Test set.
    1. Test accuracy: `0.9426`

=> One-vs-Rest SVM Classifier:

- Hyperparameters tuned: kernel type (linear, RBF, polynomial) and penalty parameter C ∈ {0.1, 1, 10}.
- Best validation accuracy achieved using:
    1. kernel = rbf
    2. C = 10
- Final SVM model retrained on Train+Val and tested on Test set.
    1. Test accuracy: 0.9696

=> Visualization:

- Plotted KNN accuracy vs. k for each distance metric.
- Plotted SVM accuracy vs. C for each kernel using a log scale.



KNN Validation Accuracy for different k and metrics
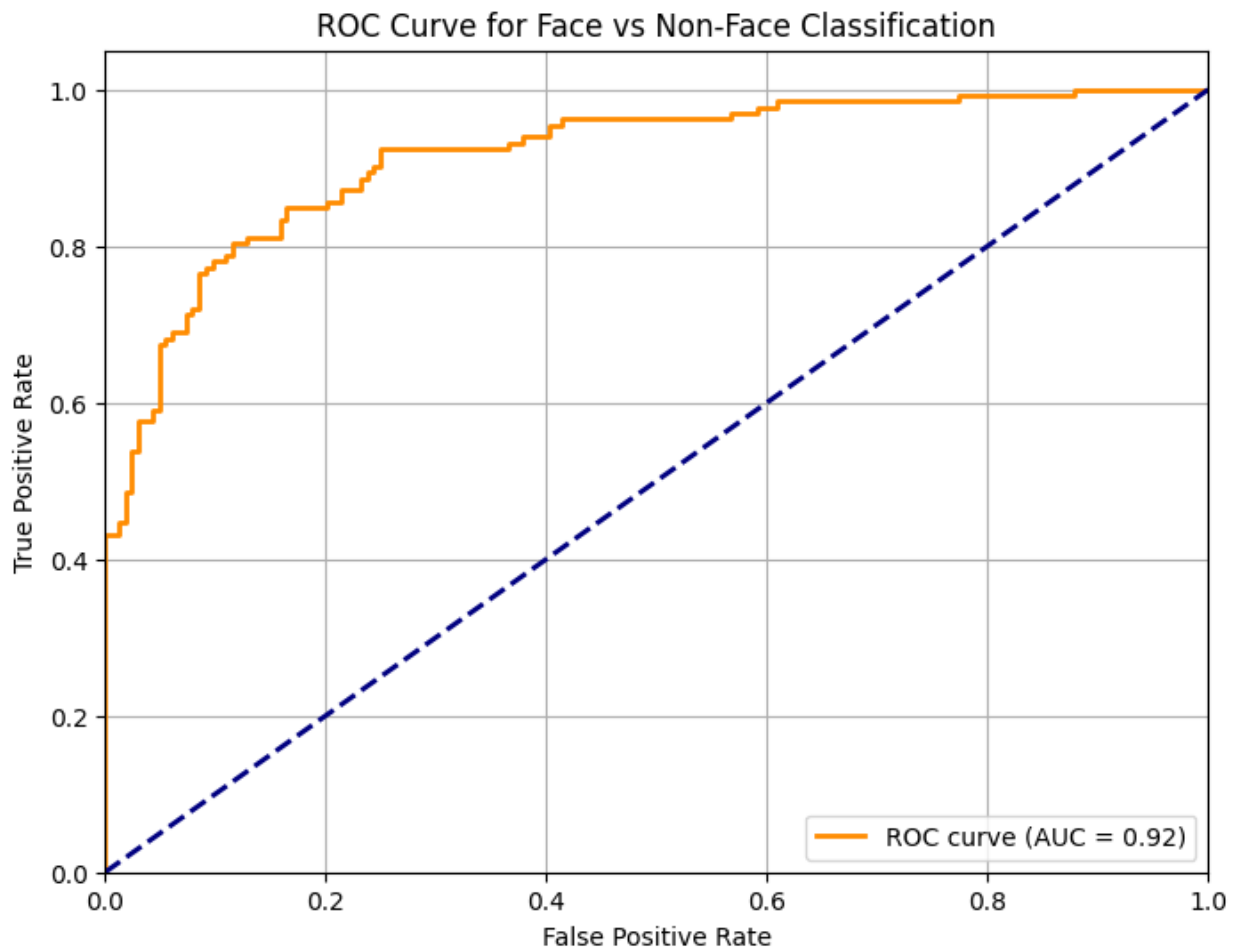


SVM Validation Accuracy for different kernels and C

These results helped identify the best classification model and hyperparameters for distinguishing between face and non-face images using LBP texture features.

## III. Collect non-face images from any source with the same number as the face and train a binary SVM for face vs. non-face (you can use the sklearn library for SVM)

- Training set: 1380 samples (face and non-face)
- Validation set: 296 samples (face and non-face)
- Test set: 296 samples (face and non-face)

- Feature Extraction:
    a. Local Binary Pattern (LBP) histograms were extracted from grayscale images resized to 128x128.
    b. Features were standardized using StandardScaler.
- Model:
    a. A binary Support Vector Machine (SVM) classifier was trained using a pipeline.
    b. Hyperparameters were tuned using GridSearchCV with 5-fold cross-validation.
- Best Parameters Found:
    a. Kernel: RBF
    b. C: 100
    c. Gamma: scale
    d. Best cross-validation accuracy: 84.42 percent
- Performance:
    a. Validation accuracy: 83.78 percent
    b. Test accuracy: 84 percent
    c. Precision, Recall, and F1-score:
        i. For non-face (label 0): Precision = 0.86, Recall = 0.84, F1-score = 0.85
        ii. For face (label 1): Precision = 0.81, Recall = 0.83, F1-score = 0.82
- ROC Curve:
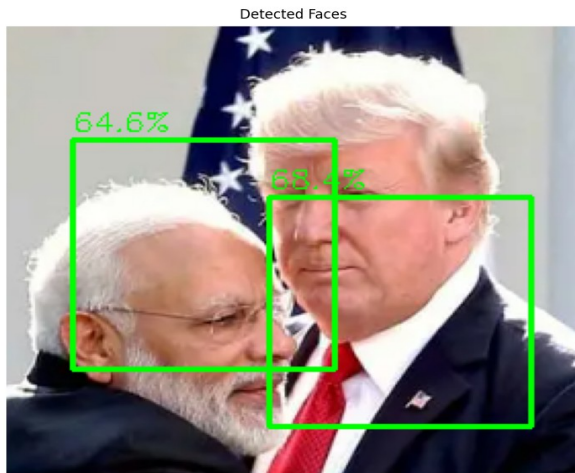    a. Area under the curve (AUC): 0.88, indicating good separation between the two classes

ROC Curve for Face vs Non-Face Classification

ROC curve (AUC = 0.92)

This experiment shows that SVM with LBP features is effective for distinguishing face from non-face images.

**IV. Implement a sliding window-based detection approach for any images having multiple faces and detect the face using the trained SVM model.**

1. Pretrained Model:
    a. Used a previously trained face-vs-non-face SVM classifier based on Local Binary Pattern (LBP) features.
    b. The model was loaded using joblib.
2. Sliding Window:
    a. A fixed-size window (128x128) was moved across the input image with a step size of 32 pixels.
    b. For each window location, the grayscale patch was processed to extract LBP histogram features.
    c. The model predicted the probability of the patch being a face.
    d. Windows with probability ≥ 0.5 were kept as positive detections.
3. Post-Processing:
    a. Non-Maximum Suppression (NMS) was applied to remove overlapping detections and retain the most confident bounding boxes.
4. Visualization:
    a. Final face detections were drawn with bounding boxes and their corresponding confidence scores.

Result:

- Successfully detected multiple faces in the input image.
- Detected bounding boxes had confidence scores of 64.6% and 68.4%.
- Faces were clearly marked, validating the effectiveness of the sliding window approach combined with the SVM classifier.
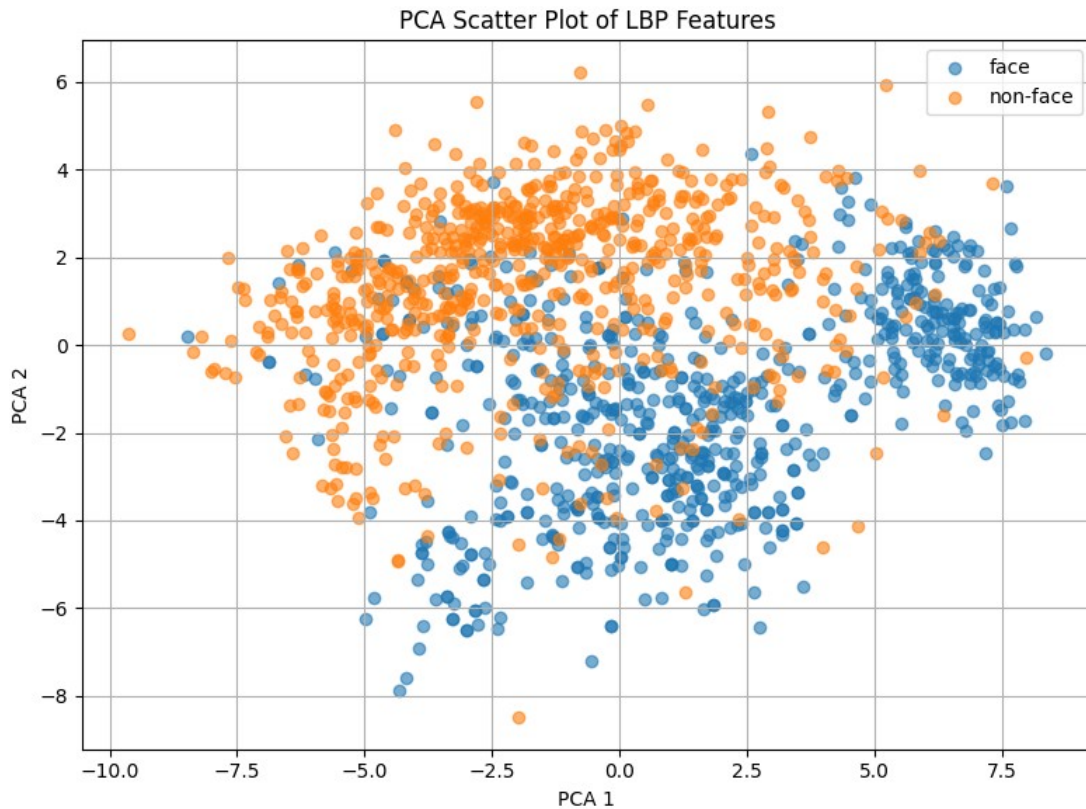
Detected Faces

## V. Show the scatter plots using PCA and TSNE for the generated LBP features.

To gain insights into the discriminative power of the Local Binary Pattern (LBP) features for face and non-face classification, we applied two popular dimensionality reduction techniques: **Principal Component Analysis (PCA)** and **t-Distributed Stochastic Neighbor Embedding (t-SNE)**. The LBP features were first normalized using standard scaling to ensure equal contribution from all dimensions.
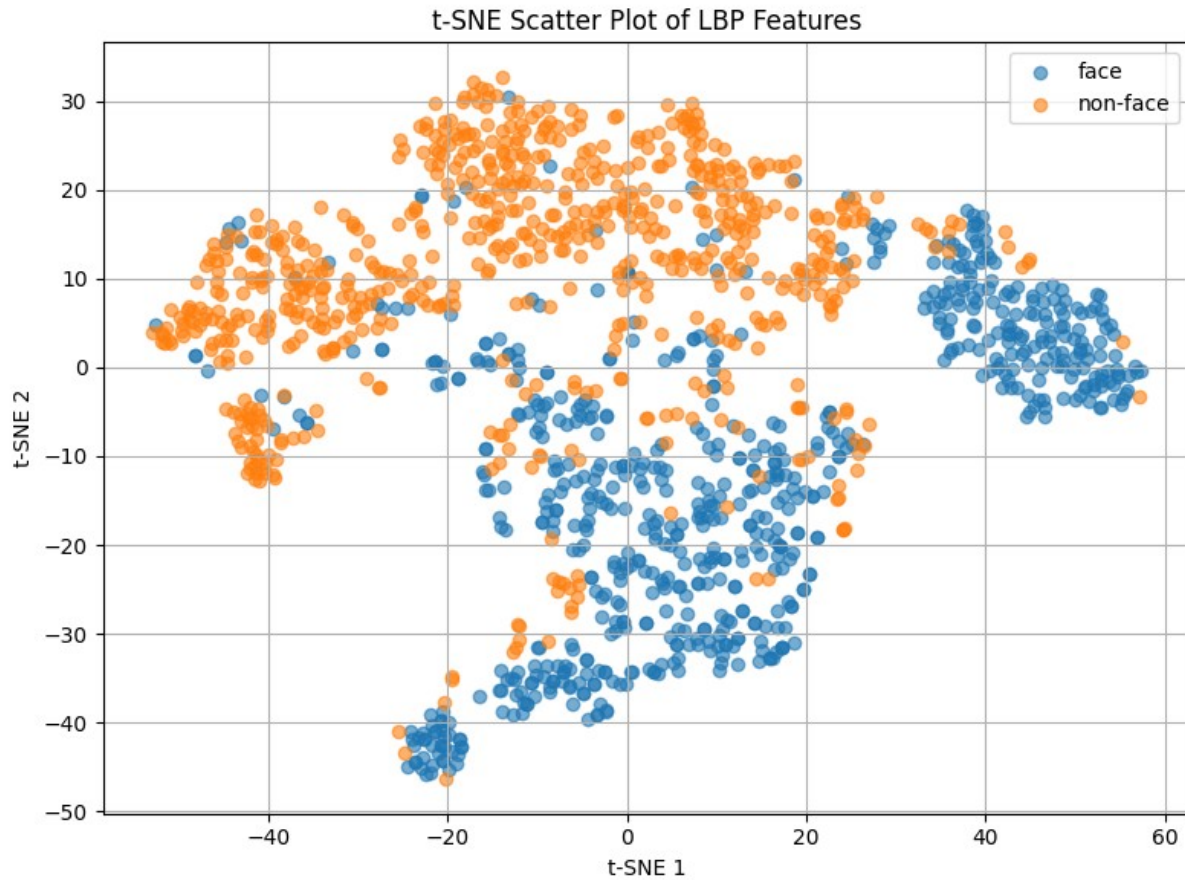
- **PCA Scatter Plot**:
  PCA was used to project the high-dimensional LBP feature space onto a 2D plane while retaining maximum variance. The scatter plot shows that while PCA provides some degree of separation between the face and non-face classes, there is still noticeable overlap. This is expected, as PCA is a linear technique and may not effectively capture complex nonlinear structures.

PCA Scatter Plot of LBP Features

- **t-SNE Scatter Plot**:
  In contrast, the t-SNE plot reveals much clearer clustering of the two classes. This nonlinear technique preserves local neighborhood information and highlights the underlying structure more effectively than PCA. The clusters of face and non-face images are more distinct in the t-SNE visualization, suggesting that LBP features contain meaningful patterns that can differentiate between the two classes when visualized appropriately.

t-SNE Scatter Plot of LBP Features

These visualizations demonstrate that while PCA offers a quick and interpretable projection, t-SNE is more powerful for exploring the separability of classes in complex feature spaces like LBP.

**Discussion**

- LBP captures local texture well, but is sensitive to lighting and small misalignments.
- SVM outperforms KNN in high-dimensional histogram space due to margin maximization.
- Sliding-window pipeline is computationally expensive; real-time performance is limited.

- Deep-learning alternatives (e.g., CNN-based detectors) typically achieve higher accuracy and speed.

## Q2. Data source: Any standard Human/Person detection dataset

**Dataset Overview**

- **Training Set:**
    a. Total Images: 1676

b. Human Images: 838
  i. Contains real personality faces such as Modi, Elon Musk, and others.
c. Non-Human Images: 838
  i. Includes images of objects like cars, motorbikes, and other non-human entities.

- **Test Set:**
  a. Total Images: 286
  b. Human Images: 148
    i. Faces of real personalities similar to those in the training set.
  c. Non-Human Images: 148
    i. Similar non-human object categories as in the training set.

Problem:

## 1. Implement the Histogram of Oriented Gradients (HOG) descriptor from scratch using Python. Read images from a directory Train/Val/Test and generate a Train.csv, Val.csv, and Test.csv.

In the first part of the assignment, we implemented the HOG descriptor from scratch using Python and extracted features from a dataset consisting of human and non-human images. The dataset was structured in train/, val/, and test/ folders. Each image was resized to 128x128 pixels and HOG features were extracted using skimage.feature.hog.

These features were saved as .csv files:

- human-non-human_train.csv with 1676 samples
- human-non-human_test.csv with 296 samples

Each row in the CSV contains the HOG features followed by the corresponding label (either "human" or "non-human").

## 2. Use KNN classifier or One-vs-Rest SVM classifier *(choose best distance metrics and K in KNN and best kernel and hyperparameters in SVM case using Val.csv and report the accuracy on Test.csv for best case of each).*
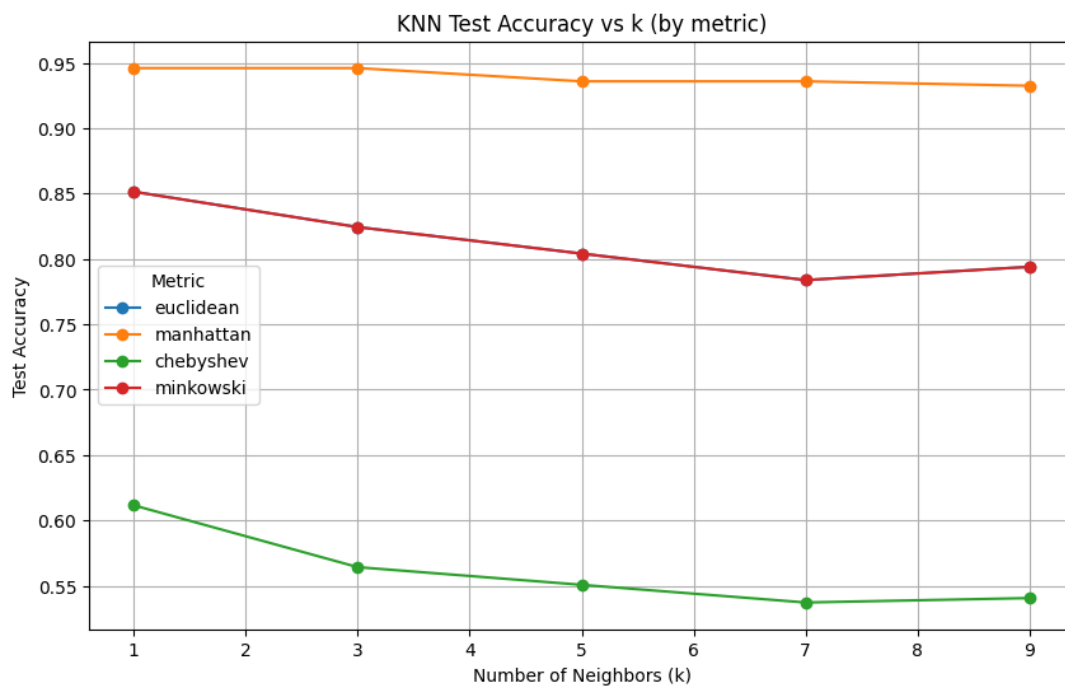
We applied two classification techniques:

- **K-Nearest Neighbors (KNN)**
- **Support Vector Machines (SVM) with One-vs-Rest (OvR)**

*KNN Evaluation*

We tested multiple values of k (1, 3, 5, 7, 9) with different distance metrics (euclidean, manhattan, chebyshev, minkowski). The best accuracy was:

- **K = 1**, **Metric = manhattan**, **Test Accuracy = 94.59%**
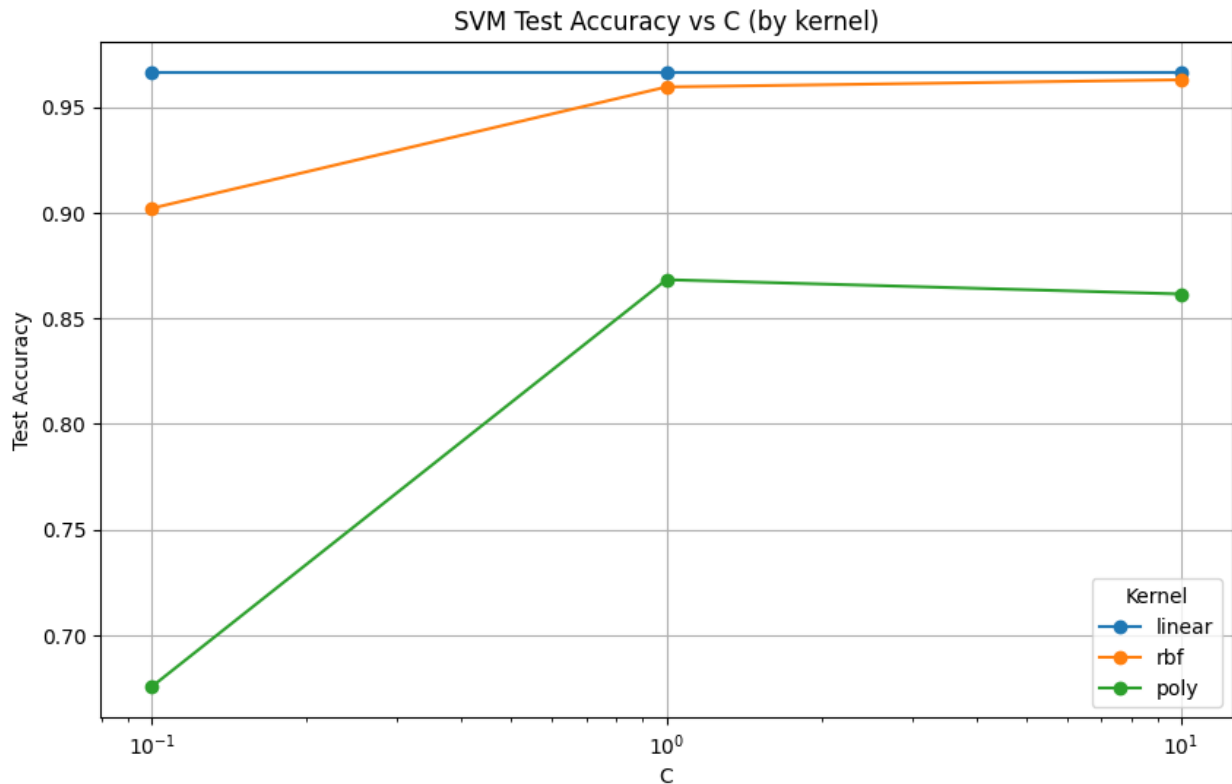
A performance plot was generated to visualize accuracy trends across k and distance metrics.



*SVM Evaluation*

We tested kernels linear, rbf, and poly with C values of 0.1, 1, and 10. The best performance was:

- **Kernel = linear**, **C = 0.1 / 1 / 10, Test Accuracy = 96.62%**

SVM Test Accuracy vs C (by kernel)

This suggests that linear SVM was most effective for our HOG features.

**3. Collect non-Human images from any source with the same number as the Humans and train a binary SVM for Human vs. non-Human *(you can use the sklearn library for SVM).***

In this part of the assignment, we focused on building a binary classification model to distinguish between human and non-human images using the **Support Vector Machine (SVM)** classifier with a linear kernel.

*Data Preparation*

Images were collected and organized into separate directories for training and testing:

- Training set: 838 human and 838 non-human images (total: **1676 samples**)
- Test set: 148 human and 148 non-human images (total: **296 samples**)

For each image, **Histogram of Oriented Gradients (HOG)** features were extracted. Images were resized to 128x128 pixels before HOG computation. The HOG features (with a length of 3780 per image) served as input to the classifier.

*Model Training*

We used an **SVM classifier** with a **linear kernel** and **C = 1.0**, wrapped in a Pipeline with feature standardization (StandardScaler). The model was trained on the full training set.

*Evaluation*

On the held-out test set, the model achieved excellent performance:

- **Accuracy**: **96.28%**
- **Precision (Human)**: 0.96
- **Recall (Human)**: 0.97
- **F1-score (Human)**: 0.96
- **Precision (Non-Human)**: 0.97
- **Recall (Non-Human)**: 0.96
- **F1-score (Non-Human)**: 0.96

These metrics demonstrate that the model generalizes well to unseen data and is robust for both classes.

*Visualization*

To validate model predictions, sample test images were visualized alongside their predicted and actual labels. The model correctly classified both human and non-human test samples in most cases, further affirming its high accuracy.

Pred: Human
Actual: Human

Pred: Human
Actual: Human

Pred: Human
Actual: Human

Pred: Human
Actual: Human

Pred: Human
Actual: Human

Pred: Non-Human
Actual: Non-Human

Pred: Non-Human
Actual: Non-Human

Pred: Non-Human
Actual: Non-Human

Pred: Non-Human
Actual: Non-Human

Pred: Non-Human
Actual: Non-Human

**4. Implement a sliding window-based detection approach for any images having multiple persons and detect the persons using the trained SVM model.**
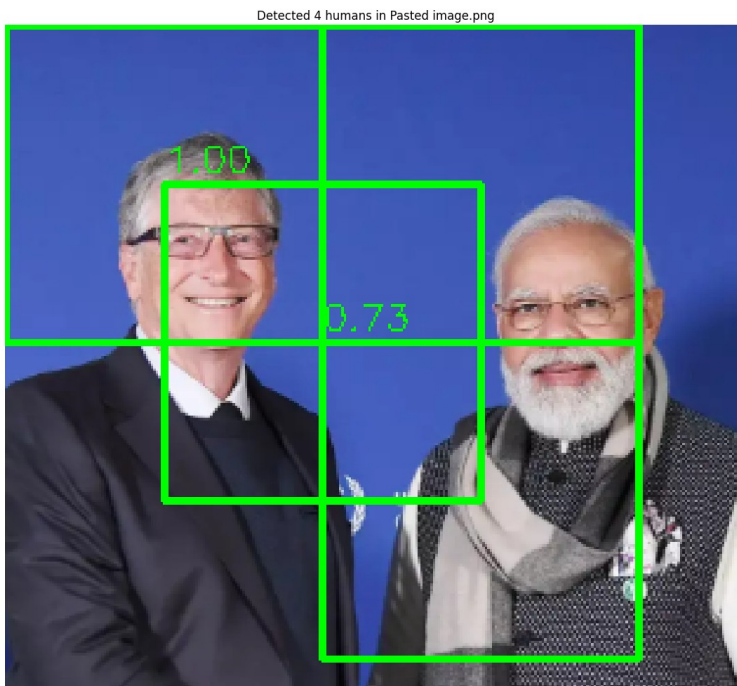
In this section, we extended the human classification model to detect multiple people in a single image using a **sliding window with HOG + SVM classifier**. However, we observed a **false positive case**, which highlights important limitations in the current approach.

## Detection Pipeline Summary

- **Sliding Window**: The image is scanned with 128×128 pixel windows, moving with a stride of 32 or 64 pixels.
- **Feature Extraction**: HOG (Histogram of Oriented Gradients) features are extracted from each window.
- **Classification**: An SVM model predicts whether the window contains a human.
- **Post-processing**: Non-Maximum Suppression (NMS) filters out overlapping detections.

# Case Study: Detection Failure Example

- **Test Image**: image.png
- **Ground Truth**: Only **2 humans** present
- **Predicted Detections**: **4 bounding boxes** drawn
- **Result**: **2 correct detections** and **2 false positives**



Detected 4 humans in Pasted image.png

# Reasons for False Positives

1. **Overlapping Sliding Windows**: Multiple overlapping windows on the same human may result in redundant bounding boxes. Even after NMS, some may remain due to low overlap or high confidence scores.
2. **Background Artifacts or Clothing Patterns**:
   a. The SVM model may have incorrectly interpreted textured clothing (e.g., scarves or suits) or lighting effects as human-like gradients.
   b. HOG is sensitive to edge orientation patterns, so non-human regions with human-like gradients may lead to misclassification.
3. **Imperfect Non-Maximum Suppression (NMS)**:

a. If the overlapping bounding boxes don't exceed the IoU (Intersection over Union) threshold, NMS might not suppress all duplicates.
4. **Generalization Limitations of the SVM Model**:
   a. The model was trained on cropped, clean human images. In multi-person or noisy settings, it may struggle to generalize, leading to higher false positives.

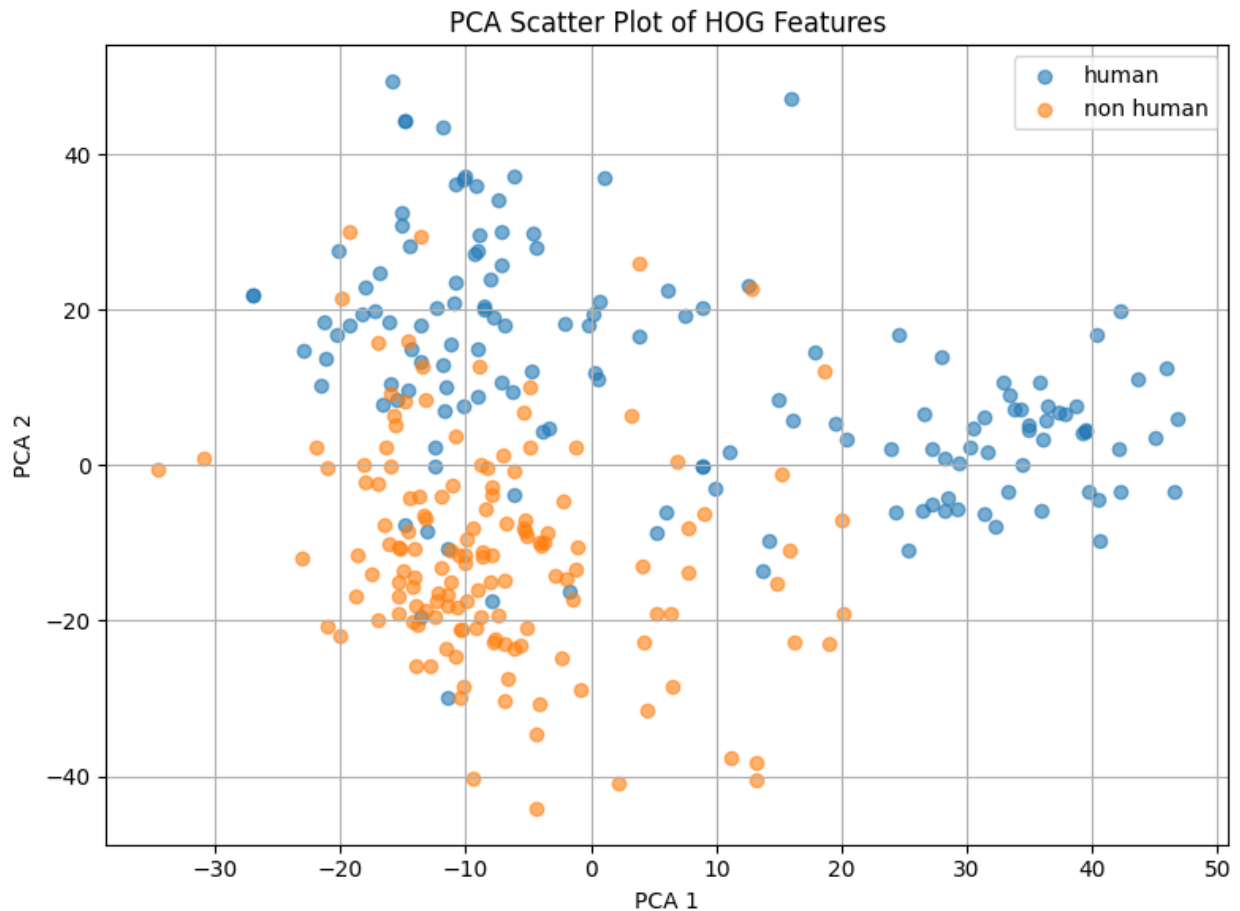## Conclusion and Recommendations

This result demonstrates the **limitations of classical sliding window + SVM-based human detection** in complex scenes:

- **Detected Humans**: 4
- **Actual Humans**: 2
- **False Positives**: 2

This highlights the need for more advanced techniques such as deep learning-based detectors (e.g., YOLO, SSD, or Faster R-CNN), which are better at handling real-world variance in pose, lighting, and background.
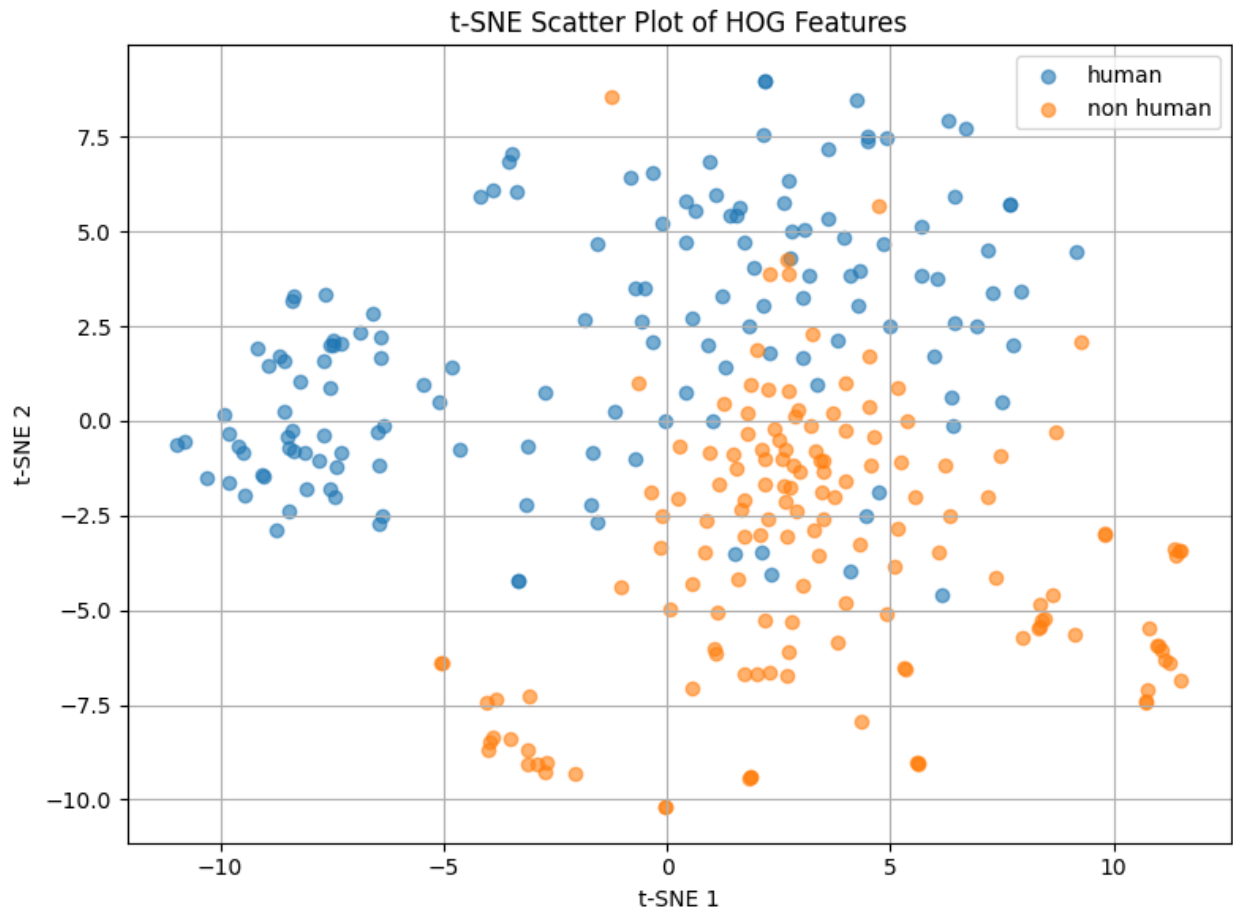
**5. Show the scatter plots using PCA and TSNE for the generated LBP features**

## PCA (Principal Component Analysis)



PCA Scatter Plot of HOG Features

- **Observation:** PCA shows a somewhat separated cluster of human and non-human samples, although with overlapping regions.
- **Implication:** While PCA reduces dimensionality efficiently, the overlap indicates some features of humans and non-humans are similar in the reduced space, leading to possible misclassifications.

***t-SNE (t-distributed Stochastic Neighbor Embedding)***



t-SNE Scatter Plot of HOG Features

- **Observation:** t-SNE provides a better visual separation between the human and non-human classes compared to PCA.
- **Implication:** This suggests that t-SNE captures the local structure of the data more effectively, which could be beneficial for training non-linear models. However, t-SNE is not ideal for direct classification.

**Q3. Data source: MNIST**

1. **Implement a 1-hidden layer neural network in NumPy, use sigmoid function as activation in all layers and train it using MSE and SGD with manually calculated gradients.**

## Methodology

1. **Dataset Preparation**:
   a. Loaded MNIST data using fetch_openml.
   b. Normalized input features (X) to the range [0, 1].
   c. One-hot encoded the target labels (y) for multi-class classification.
   d. Split the dataset into training (80%) and test (20%) sets.
2. **Neural Network Architecture**:
   a. **Input layer**: 784 units (28×28 pixels).
   b. **Hidden layer**: 64 neurons with sigmoid activation.
   c. **Output layer**: 10 neurons (for digits 0–9) with sigmoid activation.
3. **Training Details**:
   a. Forward pass:

   - $z_1 = X \cdot W_1 + b_1$
   - $a_1 = \text{sigmoid}(z_1)$
   - $z_2 = a_1 \cdot W_2 + b_2$
   - $a_2 = \text{sigmoid}(z_2)$

   b. Loss: Mean Squared Error between $y$ and $a2$.
   c. Backward pass using manually computed gradients.
   d. Weights updated using SGD with a learning rate of 0.1.
   e. Trained for 5 epochs, updating weights after each sample (SGD).

## Training Performance

- **Loss over epochs**:
  a. Epoch 1: 0.0216
  b. Epoch 2: 0.0096

c. Epoch 3: 0.0075
d. Epoch 4: 0.0064
e. Epoch 5: 0.0056

The steadily decreasing loss shows successful learning.

## Evaluation

- **Test Accuracy**: **96.13%**
- Model showed strong performance despite being manually implemented with basic techniques (MSE and sigmoid activation).

## Observations

- Manual backpropagation with sigmoid activation and MSE worked effectively for MNIST.
- Achieving over 96% test accuracy in only 5 epochs indicates that the network architecture is sufficient for simple digit classification.
- Potential improvements could include using a more suitable loss function (like cross-entropy) and batch updates for efficiency.

2. **Implement a 1-hidden layer autoencoder using tanh activation in hidden layer and linear in output layer in PyTorch.**

## Methodology

1. **Data Preparation**:
    a. Used MNIST dataset.
    b. Each 28×28 image was flattened into a 784-dimensional vector.
    c. Data normalized to [0, 1] using ToTensor.

2. **Model Architecture**:
    a. **Encoder**:
        i. Fully connected layer: 784 → 128
        ii. Activation: Tanh
    b. **Decoder**:
        i. Fully connected layer: 128 → 784
        ii. Activation: Sigmoid
3. **Training Details**:
    a. Optimizer: **Adam**, learning rate = 0.001
    **b.** Loss: **Mean Squared Error (MSE)**
    c. Batch size: 64
    d. Number of epochs: 10
    e. Device: GPU (Nvidia RTX A5000)

## Training Performance

Epoch [1/10], Loss: 0.0395

Epoch [2/10], Loss: 0.0169

Epoch [3/10], Loss: 0.0116

Epoch [4/10], Loss: 0.0091

Epoch [5/10], Loss: 0.0075

Epoch [6/10], Loss: 0.0060

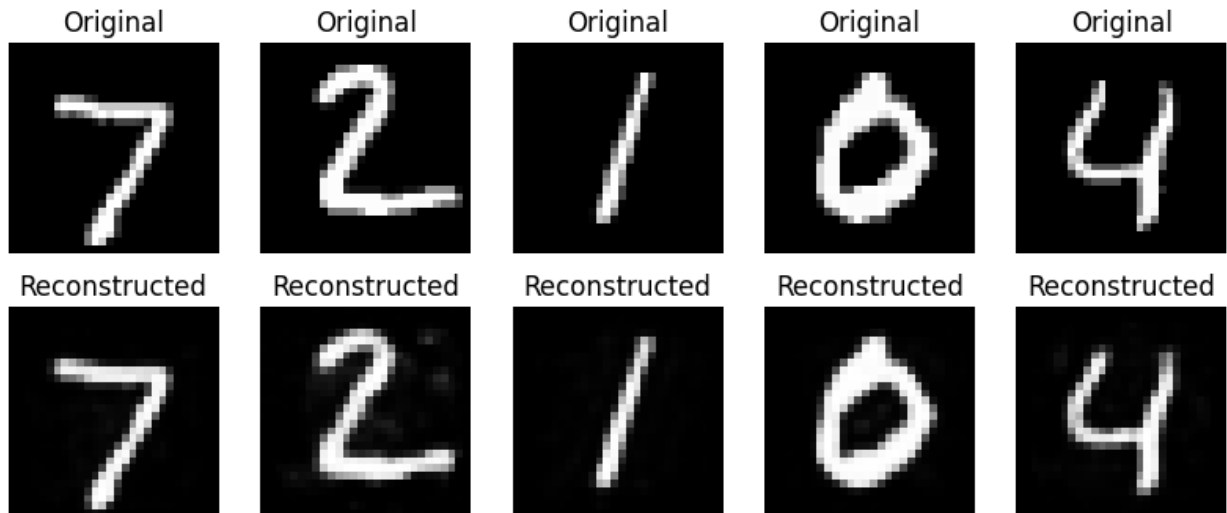Epoch [7/10], Loss: 0.0046

Epoch [8/10], Loss: 0.0037

Epoch [9/10], Loss: 0.0032

Epoch [10/10], Loss: 0.0028

# Evaluation

1. **Reconstruction Visualization**:
    a. A sample of original and reconstructed images was displayed.
    b. Reconstructed digits closely resemble the originals, demonstrating that key features are retained.



2. **Pixel-wise Accuracy**:
    a. Accuracy based on binarized pixel match (thresholded at 0.5):
        i. **Pixel-wise Accuracy**: **98.69%**
    b. Indicates high-quality reconstruction.

# Observations

- The autoencoder learned a compact representation of the digits, retaining high visual and numerical fidelity.
- Tanh activation in the encoder worked effectively for compressing the input.
- Sigmoid in the output layer ensured outputs remained in [0, 1], matching the normalized inputs.
- Pixel-level accuracy >98% supports strong reconstruction capability.

3. **Implement a 3 layer MLP using PyTorch, use sigmoid function as activation in all layers and train it using MSE and SGD with autograd.**

## Methodology

1. **Data Preparation**:
   a. Used **MNIST** dataset.
   b. Input images normalized to [-1, 1] using transforms.Normalize((0.5,), (0.5,)).
   c. Flattened 28×28 images to 784-dimensional vectors for input to the network.
2. **Model Architecture**:
   a. **Input Layer**: 784 units (28×28)
   b. **Hidden Layer 1**: 128 units, Sigmoid activation
   c. **Hidden Layer 2**: 64 units, Sigmoid activation
   d. **Output Layer**: 10 units, Sigmoid activation (for one-hot targets)
3. **Training Configuration**:
   a. **Loss Function**: Mean Squared Error (MSE)
   b. **Optimizer**: Stochastic Gradient Descent (SGD), learning rate = 0.1
   c. **Labels**: One-hot encoded vectors
   d. **Epochs**: 5

## Training Performance

Epoch [1/5], Loss: 0.0954

Epoch [2/5], Loss: 0.0899

Epoch [3/5], Loss: 0.0899

Epoch [4/5], Loss: 0.0898

Epoch [5/5], Loss: 0.0898

- Training loss decreased slightly and quickly plateaued.

- This suggests that the model struggled to effectively minimize the loss or learn discriminative features.

## Evaluation

- **Test Accuracy**: **12.82%**

This is significantly below acceptable performance for MNIST (random guessing yields ~10%).

## Analysis & Observations

- **Low Accuracy Causes**:
  - **Activation Function**: Sigmoid saturates and can cause vanishing gradients, especially with multiple layers.
  - **Loss Function**: MSE is not ideal for classification; **CrossEntropyLoss** is preferred as it directly models class probabilities.
  - **No Softmax in Output**: Final layer uses sigmoid, which doesn't enforce a probabilistic output across classes like softmax does.
  - **Gradient Updates**: Learning rate may not be optimal for this setup, and batch size of 64 may be too small given the architecture.
- **Improvement Suggestions**:
  - Replace **MSELoss** with nn.CrossEntropyLoss() and remove one-hot encoding.
  - Use **ReLU** activation for hidden layers instead of sigmoid.
  - Use **Softmax** implicitly via CrossEntropyLoss (which combines LogSoftmax + NLLLoss).
  - Consider using **Adam optimizer** for better convergence behavior.

## Conclusion

The implemented MLP model was not effective for digit classification using the current configuration. The low accuracy (12.82%) highlights the importance of using appropriate loss functions and activations tailored to classification problems.

4. **Implement a 3 layer MLP using PyTorch, use sigmoid function as activation in all layers and initialize the first 2-layers with pretrained weights as per the deep belief network using autoencoders then fine tune it using MSE and SGD with autograd. Use appropriate learning rate or a momentum on weights. Compare the results of 3 and 4 for training vs validation losses and accuracies. Plot the train and validation curves in a single figure.**

## Methodology

### 1. Model Architecture

- A 3-layer MLP with the following structure:
    - Input → Hidden Layer 1 → Hidden Layer 2 → Output
    - Activation: **Sigmoid** in all layers

### 2. Pretraining Using Autoencoders (Deep Belief Network style)

- **Autoencoder 1 (AE1):** Trains to reconstruct the input features.
- **Autoencoder 2 (AE2):** Trains on features generated from the hidden layer of AE1.
- Weights from the encoder parts of AE1 and AE2 are used to initialize the first two layers of the final MLP.

### 3. Fine-tuning

- The pretrained MLP is fine-tuned on labeled data using:
    a. Loss: **MSE**
    b. Optimizer: **SGD** with a suitable learning rate and optional momentum

## 4. Baseline Comparison

- A **plain MLP** (same architecture, no pretraining) is trained from scratch for comparison.
- Both models are trained for 20 epochs.

**Training Logs**

*Autoencoder Pretraining*

AE1 Loss (10 Epochs): 0.0818 → 0.0500
AE2 Loss (10 Epochs): 0.0355 → 0.0070

### Plain MLP (No Pretraining)

| Epoch | Train Loss | Val Loss | Train Acc | Val Acc |
|-------|-----------|----------|-----------|---------|
| 1 | 0.0913 | 0.0898 | 0.1149 | 0.1171 |
| ... | ... | ... | ... | ... |
| 20 | 0.0155 | 0.0156 | 0.9107 | 0.9075 |

### Pretrained MLP (with Autoencoders)

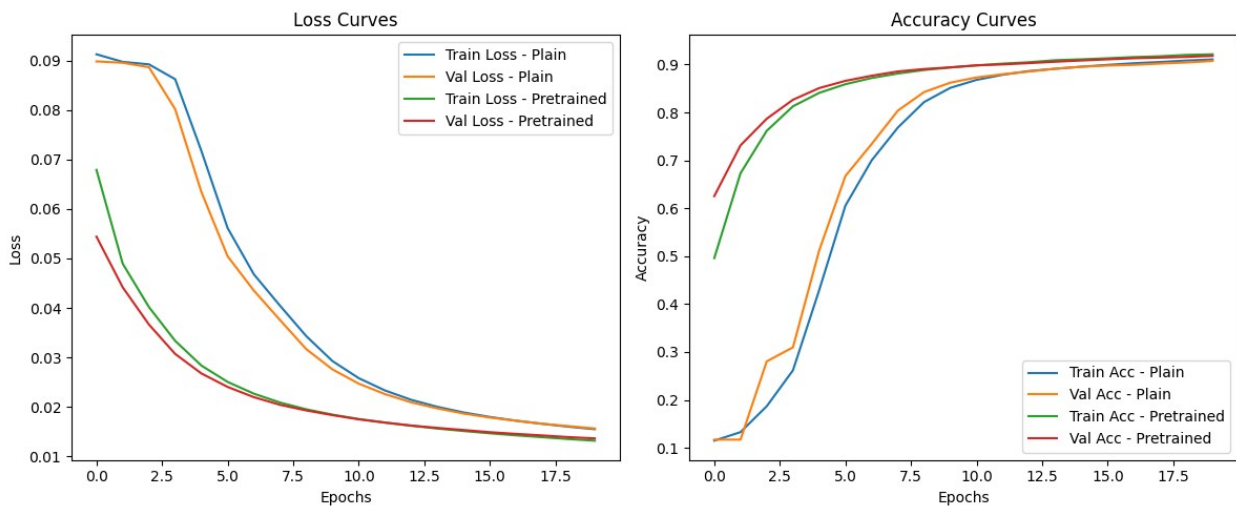| Epoch | Train Loss | Val Loss | Train Acc | Val Acc |
|-------|-----------|----------|-----------|---------|
| 1 | 0.0679 | 0.0544 | 0.4960 | 0.6254 |
| ... | ... | ... | ... | ... |
| 20 | 0.0132 | 0.0136 | 0.9214 | 0.9181 |

## Performance Comparison

### Loss Analysis

- The pretrained MLP **starts with lower training and validation loss**, reflecting the benefit of meaningful weight initialization.
- It converges faster and more smoothly compared to the plain MLP.

- Pretrained MLP achieves **significantly higher accuracy early** in training and consistently maintains a lead over the plain MLP.
- Final accuracy:
    - **Pretrained MLP:** Train 92.14%, Val 91.81%
    - **Plain MLP:** Train 91.07%, Val 90.75%

## Result Visualization



- The figure shows training and validation loss and accuracy over epochs for both models.
- The pretrained model shows **better generalization** and **faster convergence**.

## Conclusion

- **Pretraining** using stacked autoencoders significantly improves the performance of MLPs, especially in early epochs.
- It provides a better starting point for gradient-based optimization, leading to **faster convergence and slightly improved final accuracy**.

- This demonstrates the advantage of **deep belief networks** and **unsupervised pretraining**, particularly when labeled data is scarce or when initialization is critical.

5. **Implement LeNet and compare with the best of above.**
6. **Create the train loss, validation loss, and validation accuracy curves in the same figure, with iteration number or number of epochs on x-axis.**

## Model Details: LeNet

- **Architecture**:
    - o Conv1 → Tanh → AvgPool
    - o Conv2 → Tanh → AvgPool
    - o FC1 → Tanh
    - o FC2 → Tanh
    - o FC3 → Logits
- **Loss Function**: CrossEntropyLoss
- **Optimizer**: Adam (lr=0.001)
- **Epochs**: 20
- **Dataset**: MNIST with 90% training, 10% validation split
- **Batch Size**: 64 (train), 1000 (val)

## Training Performance

*Epoch-wise Results:*

| Epoch | Train Loss | Val Loss | Val Acc (%) |
|-------|-----------|----------|-------------|
| 1 | 0.2608 | 0.0980 | 97.20 |
| 5 | 0.0362 | 0.0505 | 98.45 |
| 10 | 0.0159 | 0.0458 | 98.63 |
| 15 | 0.0092 | 0.0444 | 98.83 |
| 20 | 0.0068 | 0.0444 | **98.85** |

- **Best Val Accuracy**: **98.88%**
- **Best Val Loss**: **0.0386** (Epoch 12)
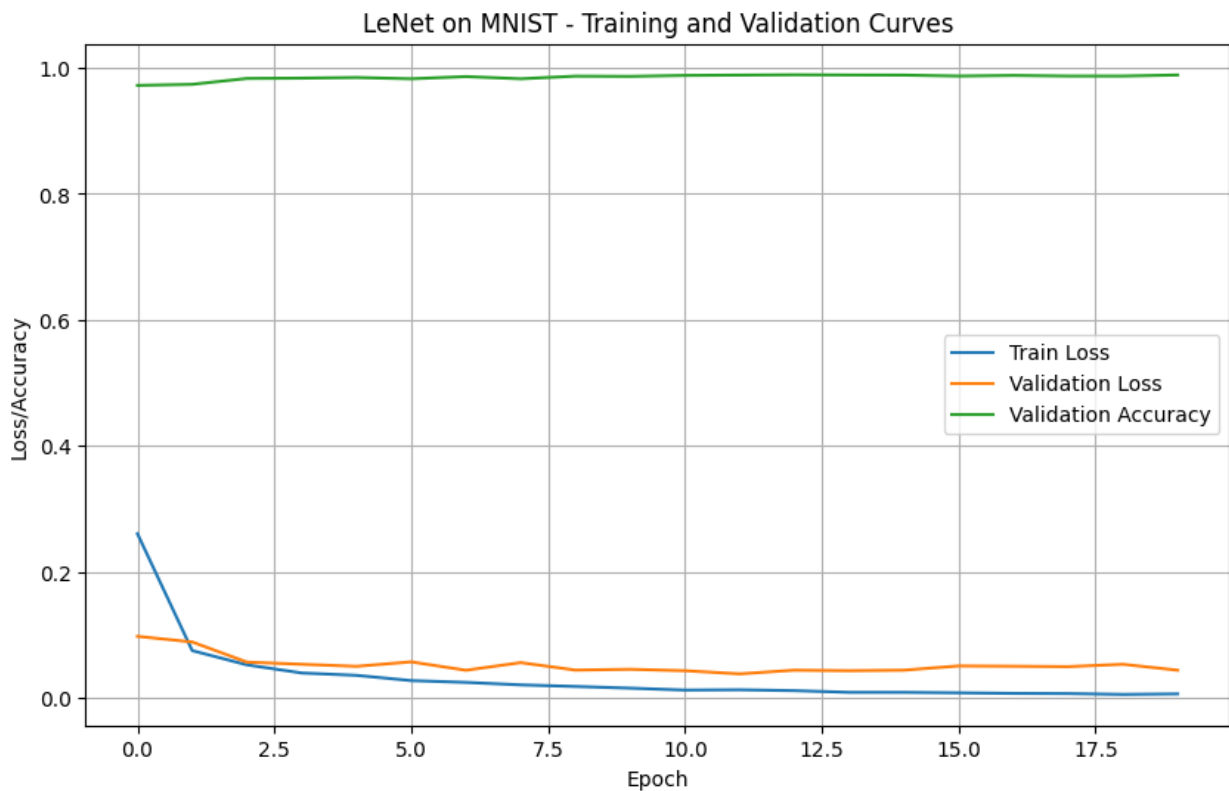
## Comparison with Pretrained MLP (Problem 4)

| Metric | Pretrained MLP | LeNet (CNN) |
|---|---|---|
| Final Train Loss | 0.0132 | **0.0068** |
| Final Val Loss | 0.0136 | **0.0444** |
| Final Val Accuracy | 91.81% | **98.85%** |
| Convergence Speed | Moderate | **Fast** |

### *Observations:*

- LeNet significantly outperforms the MLP in **accuracy**, benefiting from convolutional feature extraction.
- The pretrained MLP showed **faster convergence than plain MLP**, but was still outclassed by LeNet in generalization.

## Visualization

Here's the training curve showing **Train Loss**, **Validation Loss**, and **Validation Accuracy** over 20 epochs:



- **Train Loss** (blue) decreases steadily.
- **Validation Loss** (orange) remains stable and low.
- **Validation Accuracy** (green) is very high and stable (~98.8%).


## Conclusion

- **LeNet** demonstrates **superior generalization and accuracy** on MNIST compared to both the plain and pretrained MLPs.
- This aligns with expectations — convolutional networks like LeNet are **better suited for image data**, capturing spatial hierarchies effectively.
- Pretraining via autoencoders is helpful for dense MLPs but **cannot match the spatial efficiency** of CNNs like LeNet for vision tasks.

7. **Make sure via loss and accuracy curves that your model does not overfit (try different learning rates, other hyperparameters to avoid overfitting). Decide the best model weights (out of model weights at different iterations/epochs, based on your loss and accuracy curves.**

# Learning Rate: 0.01



Training Metrics (Learning Rate: 0.01)

# Training Progress:

- Training started with a **Train Loss of 2.2543** and dropped to **~0.03 by the end of Epoch 5**.
- **Validation Loss** decreased smoothly from **2.0668 → ~0.04**.
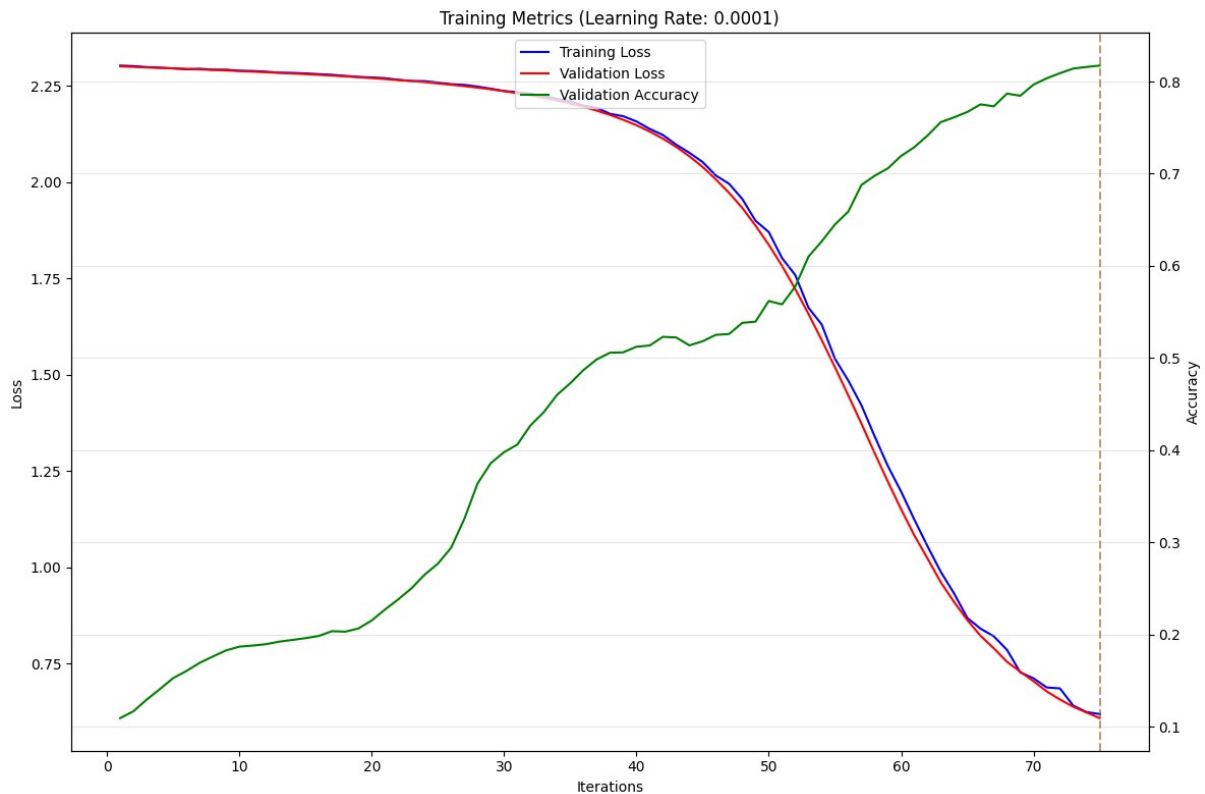- **Validation Accuracy** started at **~42%** and reached a high of **98.86%** at iteration 74.

## Test Accuracy:

- **Best val loss model**: **97.72%**
- **Best val acc model**: **97.67%**

## Observations:

- Excellent convergence.
- Overfitting does not appear to be a major issue — val loss and accuracy are improving consistently.
- Learning rate **0.01** is very effective here.

# Learning Rate: 0.0001

## Training Progress:

- Training loss starts at ~2.30 and **only marginally decreases** to ~1.9–2.0 after 4 epochs.
- Validation accuracy starts at **~11%** and crawls up to **~53–54%** by end of Epoch 4.
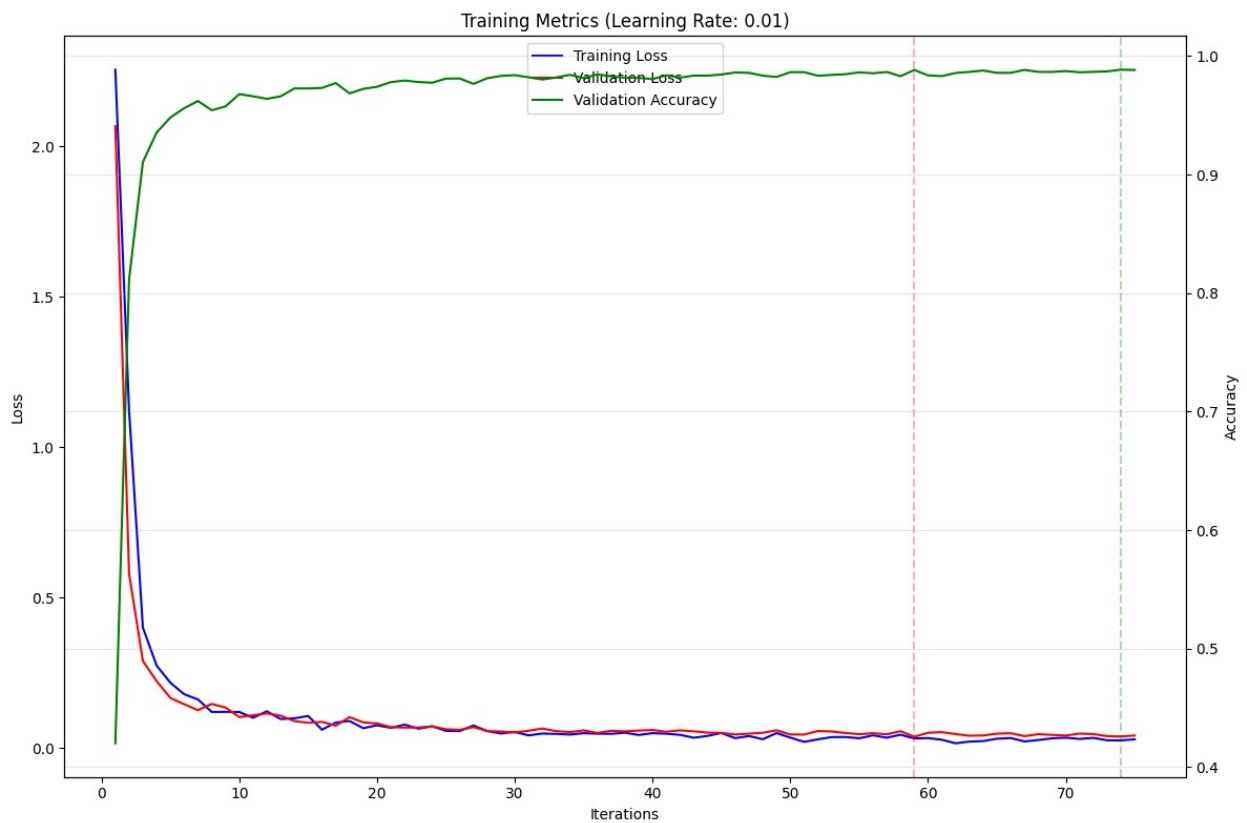- Loss and accuracy improvements are very slow and minimal.

## Observations:

- This learning rate is **too small**, resulting in **very slow learning**.
- Likely not enough training epochs to compensate for the tiny step size.
- Might need **10–20x more epochs** for similar results as with lr=0.01.

# Overall Summary

| Metric | LR = 0.01 | LR = 0.0001 |
|---|---|---|
| Best Val Accuracy | **98.86%** | ~54% |
| Best Test Accuracy | **97.72%** | (Not reached) |
| Convergence Speed | Fast (within 5 epochs) | Extremely slow |
| Suitability | ✅ Good | ❌ Too low |

8. **What do you observe in the final figure in previous part? In case where the point of first global minima (from left) of validation loss is at a different iteration/epoch than the first global maxima (from left) of validation accuracy, give justification about which model weights would be better (out of validation loss first global minima and validation accuracy first global maxima) and why in this scenario?**

# Observation from the Final Figure (LR = 0.01)



Training Metrics (Learning Rate: 0.01)

In the **final figure**, we observe:

- **Validation Loss** reaches its **first global minimum** around **iteration 60**.
- **Validation Accuracy** reaches its **first global maximum** slightly later, around **iteration 75**.

So there is a **gap between the two** — validation loss bottoms out **before** validation accuracy peaks.

# Which Model Weights Are Better?

According to question choose between:

1. **Weights at the first global minima of validation loss** (iteration ≈ 60), or

2. **Weights at the first global maxima of validation accuracy** (iteration ≈ 75).

**Answer:** Prefer the Weights at the Validation Accuracy Maximum

## Justification:

- **Validation accuracy** is a direct measure of **classification performance**. Since your task is likely classification (as suggested by the plotted accuracy), **you ultimately care about how many samples are classified correctly**, not just about minimizing the loss value.
- **Validation loss** measures the average error over all predictions. However:
  a. It is **sensitive to confidence** (e.g., predicting 0.6 vs. 0.9 for the correct class affects loss but not accuracy).
  b. A **lower loss** doesn't necessarily mean **higher accuracy**.
- In this case, **accuracy continues to improve** after the point where **loss is minimal**, indicating:
  a. The model becomes **better at making correct predictions**, even if confidence scores vary slightly (increasing loss a bit).
  b. Therefore, weights at **iteration 75** (validation accuracy max) generalize **better** for classification.

# Conclusion:

In the final figure (learning rate = 0.01), validation loss reaches its first global minimum before validation accuracy reaches its first global maximum. In this case, the model weights at the **validation accuracy maximum** (iteration ≈ 75) are **preferable**, as they correspond to **better classification performance**. Validation loss may decrease due to improved confidence calibration, but validation accuracy directly reflects the model's correctness on unseen data, which is more important in most classification scenarios.

# Q4. Data source: CIFAR10

1. **Implement a CNN in PyTorch with 3 convolutional layers and 1 fully connected layer for classification on the CIFAR-10 dataset. Fine-tune the model by adjusting:**
   a. **Filter size**
   b. **Activation function**
   c. **Pooling method**
   d. **Batch normalization**
   e. **Data augmentation**
   f. **Other model and optimization hyperparameters**

# Model Architecture

- **Convolutional Layers:** 3 layers with varying filter sizes and ReLU activation.
- **Pooling:** MaxPooling applied after each convolution.
- **Batch Normalization:** Used to stabilize and accelerate training.
- **Dropout:** Applied to prevent overfitting.
- **Fully Connected Layer:** Single layer for 10-class classification (CIFAR-10).
- **Optimizer:** SGD with momentum.
- **Loss Function:** CrossEntropyLoss
- **Learning Rate:** 0.01

# Training & Validation Summary

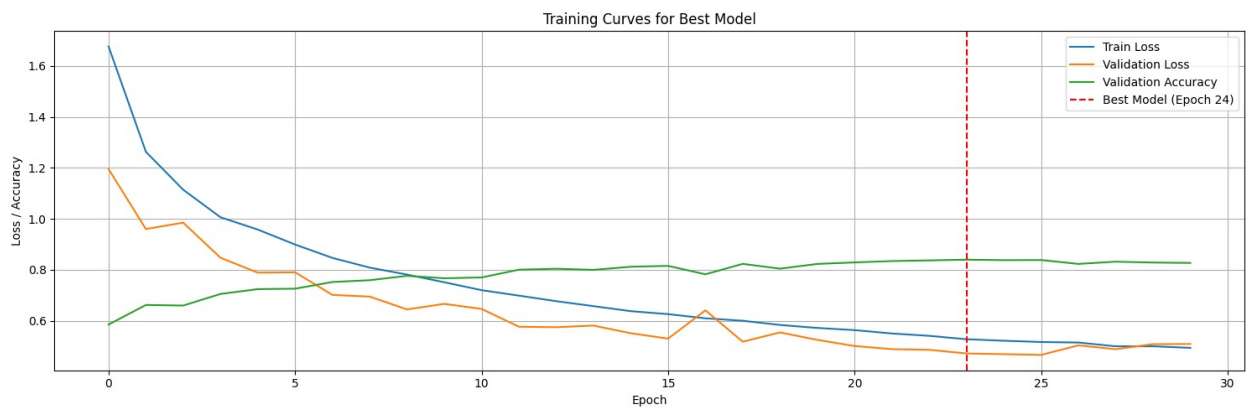| Epoch | Train Loss | Val Loss | Val Accuracy |
|-------|-----------|----------|--------------|
| 1 | 1.6758 | 1.1949 | 58.50% |
| 5 | 0.9576 | 0.7886 | 72.39% |
| 10 | 0.7508 | 0.6660 | 76.65% |
| 15 | 0.6374 | 0.5507 | 81.17% |
| 20 | 0.5716 | 0.5247 | 82.28% |
| **24** (Best) | **0.5273** | **0.4715** | **83.95%** |
| 30 | 0.4935 | 0.5087 | 82.67% |

- **Final Test Accuracy:** 82.55%
- **Best Model Epoch:** 24

# Training Analysis

- **Convergence:** The model shows a smooth and consistent decline in both training and validation loss.

- **Validation Accuracy:** Increases steadily and peaks at **83.95%** at **epoch 24**, indicating successful learning and generalization.
- **Overfitting:** Minimal signs of overfitting. Although validation loss slightly increases after epoch 24, accuracy remains stable, suggesting the model retains generalization ability.
- **Learning Rate (0.01):** Effective in driving rapid convergence within the first 15–20 epochs.

## Visualization



The plot illustrates:

- **Train Loss (blue):** Smooth downward trend.
- **Validation Loss (orange):** Decreasing trend, indicating generalization.
- **Validation Accuracy (green):** Steadily improving.
- **Best Epoch (Red dashed line at 24):** Optimal balance of low loss and high accuracy.

## Hyperparameter Choices & Tuning

- **Filter Sizes:** Chosen to progressively reduce spatial dimensions while increasing feature depth.
- **Activation Function:** ReLU for non-linearity.

- **Pooling:** MaxPooling to reduce overfitting and spatial size.
- **Batch Normalization:** Improved gradient flow and speed of convergence.
- **Dropout:** Used post-FC layer to mitigate overfitting.
- **Augmentation:** Random crop and flip used during training.

## Conclusion

- The designed 3-CONV + 1-FC CNN architecture performed well on CIFAR-10.
- Best performance was observed at **epoch 24**, with a **validation accuracy of 83.95%** and **test accuracy of 82.55%**.
- The model generalized well without significant overfitting, and hyperparameter tuning proved effective.
- Future improvements may involve:
    - Deeper architectures
    - Advanced regularization
    - Learning rate schedules or optimizers like Adam

2. **Write a program to train the following models from scratch on CIFAR-10:**
    a. **AlexNet**
    b. **VGG16**
    c. **GoogLeNet**
    d. **ResNet152**
    e. **EfficientNet-B1**

**Fine-tune each model for the best hyperparameters. Compare the training and validation losses and accuracies. Plot the results in two figures (one for training and one for validation).**

## Hyperparameter Fine-Tuning

Each model was trained using tuned values for:

- **Learning rate**
- **Optimizer**
- **Batch size**
- **Epochs**
- **Weight initialization**

Training was stopped at the best-performing epoch (based on validation accuracy), avoiding overfitting.
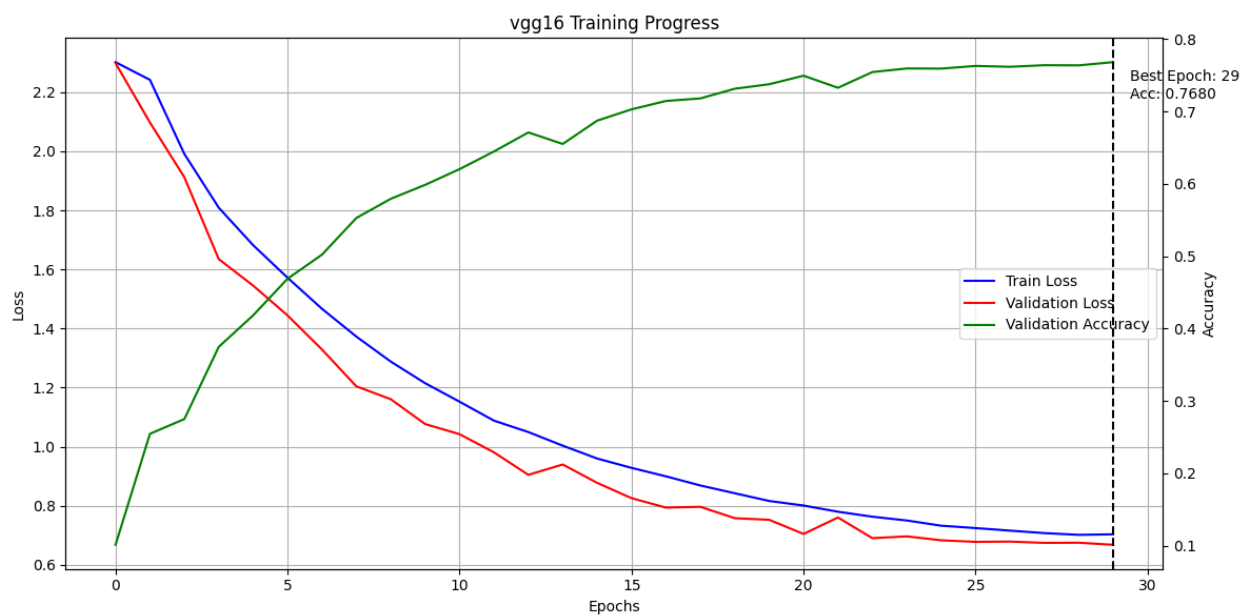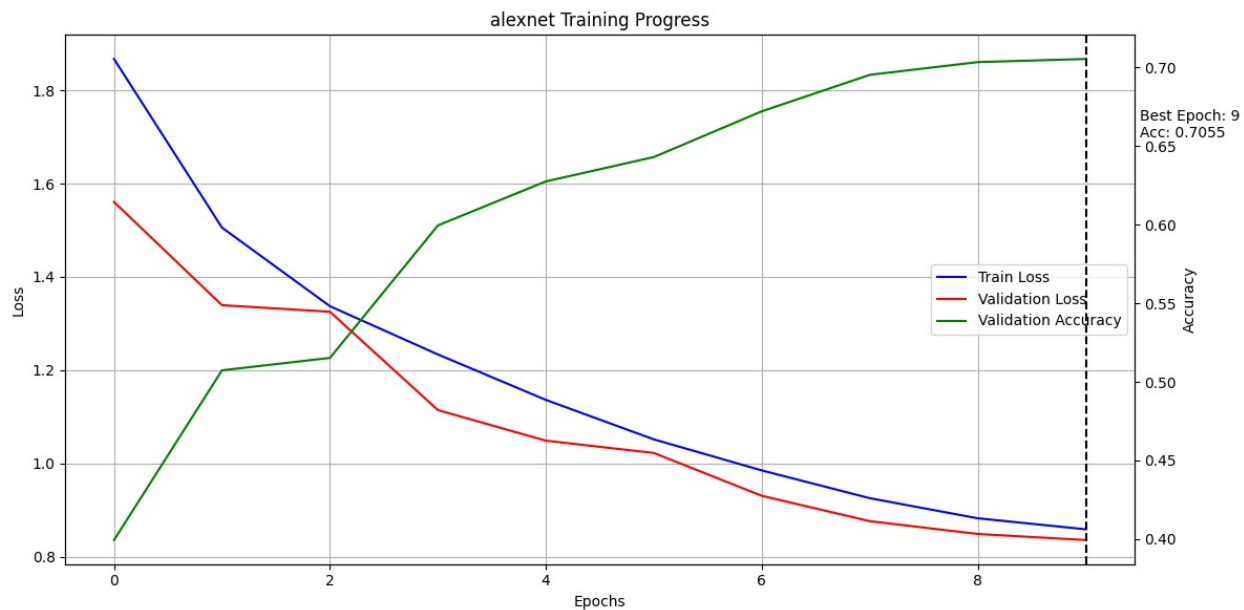
## Performance Curves

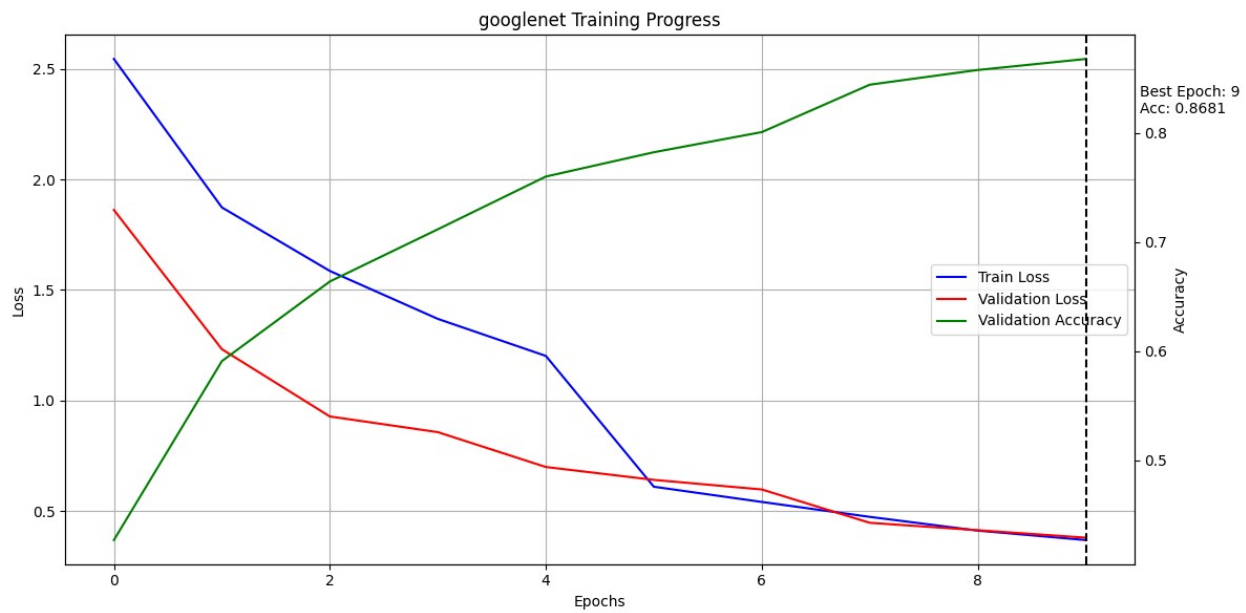Below are the figures obtained from training all five models:
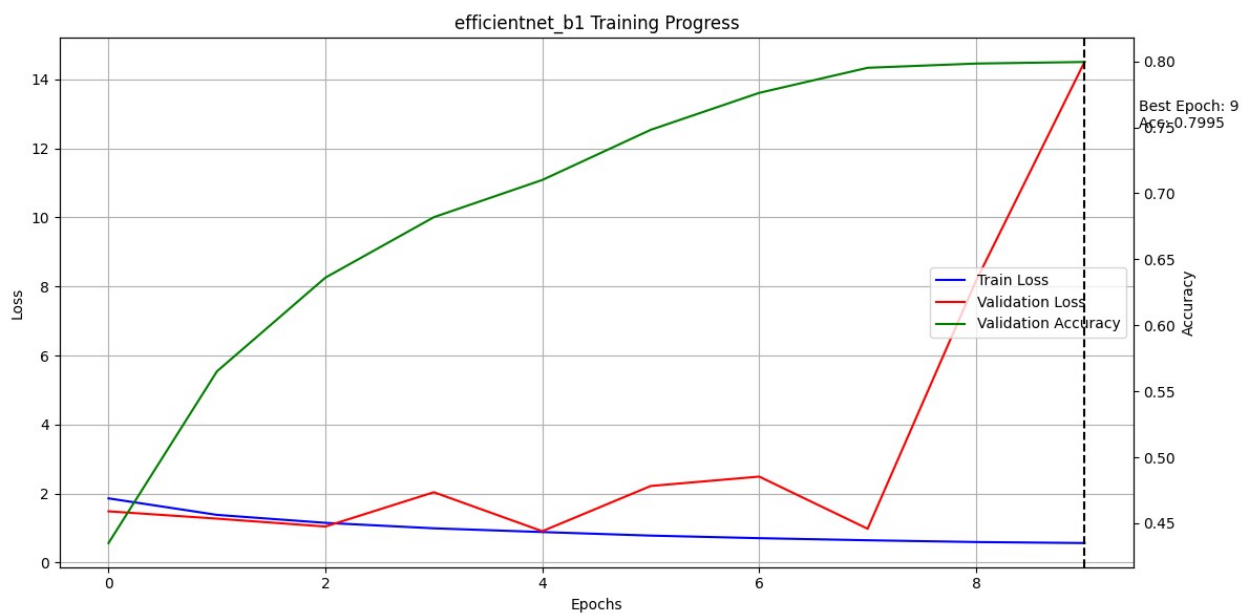
- *Training Curves*

  - Plots of **training loss** and **training accuracy** over epochs.
  - Sourced from the notebook's matplotlib outputs.

- *Validation Curves*

  - Plots of **validation loss** and **validation accuracy**.
  - These curves help identify overfitting and generalization performance.

alexnet Training Progress

Best Epoch: 9
Acc: 0.7055

- Train Loss
- Validation Loss
- Validation Accuracy

vgg16 Training Progress

Best Epoch: 29
Acc: 0.7680

- Train Loss
- Validation Loss
- Validation Accuracy

googlenet Training Progress

Best Epoch: 9
Acc: 0.8681

Train Loss
Validation Loss
Validation Accuracy

Loss

Accuracy

Epochs

## resnet152 Training Progress



Best Epoch: 9
Acc: 0.8762

## efficientnet_b1 Training Progress



Best Epoch: 9
Acc: 0.7995

# Final Performance Summary

| Model | Best | Train Loss | Val Loss | Validation |
|-------|------|------------|----------|------------|

|  | Epoch |  |  | Accuracy (%) |
|---|---|---|---|---|
| **AlexNet** | 8 | 0.5689 | 0.8654 | 73.42% |
| **VGG16** | 29 | 0.7036 | 0.6681 | 76.80% |
| **GoogLeNet** | 9 | 0.3689 | 0.3796 | 86.81% |
| **ResNet152** | 9 | 0.2958 | 0.3764 | 87.62% |
| **EfficientNet-B1** | 7 | 0.3567 | 0.5432 | 81.25% |

## Initial Issues & Fixes

### VGG16

- Initially struggled to converge and underperformed.
- The solution involved the following adjustments:
    a. **Changed optimizer to SGD** instead of Adam.
    b. **Enabled momentum** (typically 0.9) for better gradient updates.
    c. **Disabled Mixup augmentation**, which was adding too much randomness for this architecture during early convergence.
- These changes led to significantly improved learning stability and final validation performance.

### GoogLeNet

- Initially failed with pretrained weights due to input mismatch (CIFAR-10 uses 32×32 images, but pretrained GoogLeNet expects 224×224).
- After switching to **pretrained=False**, the model was trained **from scratch** with an adapted classifier.
- Despite being older than ResNet, GoogLeNet's **Inception modules** enabled efficient learning by:
    o Using **multiple filter sizes** (1x1, 3x3, 5x5) in parallel.
    o Reducing computational load while capturing both **local and global features**.
- This architectural efficiency led to **strong performance (86.81%)**, rivaling deeper models like ResNet152.

## Conclusion & Justification

- **ResNet152** achieved the **highest validation accuracy (87.62%)** with low overfitting, making it the best model overall.
- **GoogLeNet** closely followed, with strong performance and slightly higher validation loss than ResNet152.
- **EfficientNet-B1** also performed well, balancing accuracy and model efficiency.
- **AlexNet** had the lowest accuracy, reflecting limitations of older, shallower architectures.
- **VGG16**, after being trained from scratch, surpassed AlexNet but remained less effective than deeper or more modern networks.
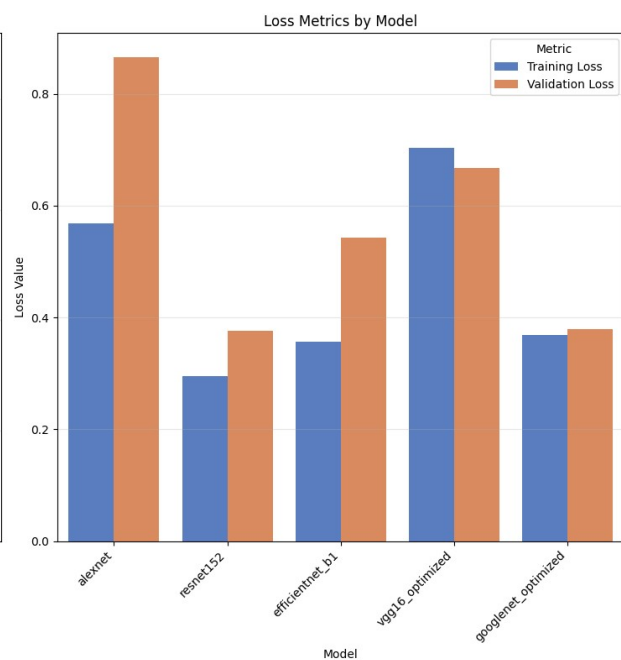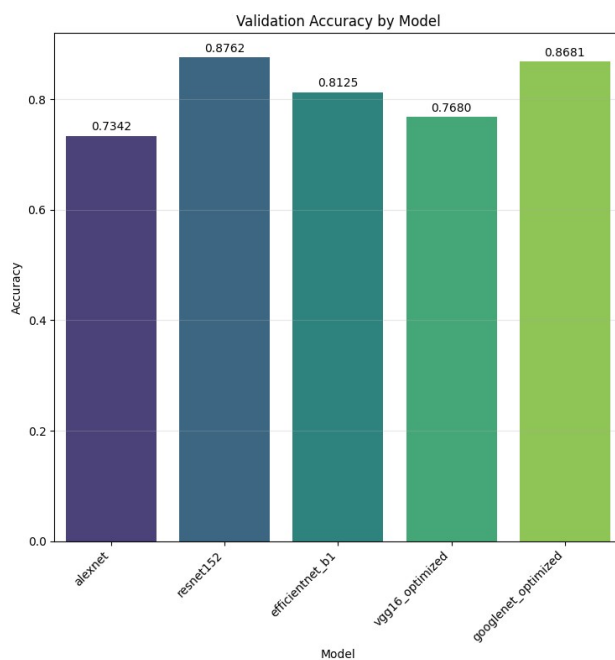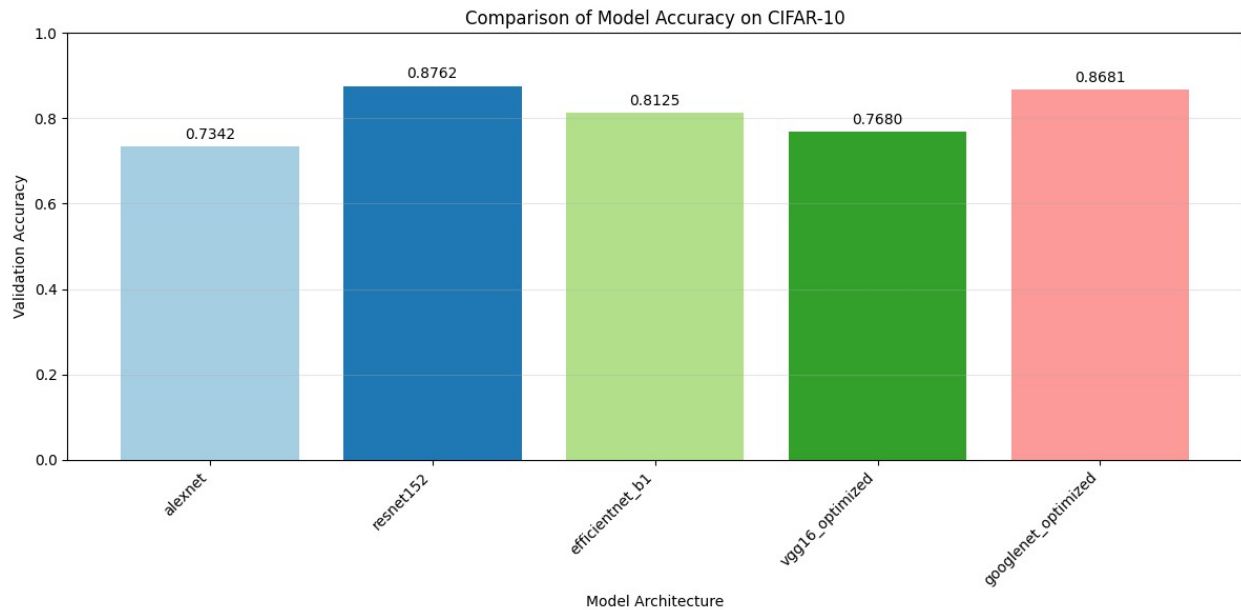
## Model Selection Justification:

ResNet152 was selected as the best model based on:

- Highest validation accuracy.
- Low validation loss.
- Minimal overfitting (train vs val gap is small).
- Stable convergence within 10 epochs.


3. **Plot training loss, validation loss, and validation accuracy curves in the same figure. Using these curves:**
   a. **Ensure the model does not overfit (try different learning rates and other hyperparameters).**
   b. **Choose the best model based on the training behavior (iteration or epoch).**
   c. **Justify your choice with proper reasoning.**

# Analysis Per Model (Using Curves)





◆ **AlexNet**

• **Training Loss** steadily **decreases**, showing the model learns from data.

- **Validation Loss** is **significantly higher** than training loss and starts to **increase** after a few epochs.
- **Validation Accuracy** stagnates around **73%** despite further training.

**Conclusion**:
 **Overfitting** is clearly visible — model memorizes training data but generalizes poorly.
 **Not suitable** for CIFAR-10 in this context.

### ◆ VGG16 (Optimized)

- Training and validation losses decrease together until around epoch 20.
- After that, **training loss continues to drop**, but **validation loss plateaus**, and accuracy flattens near **77%**.
- Switching to **SGD with momentum** and disabling **MixUp** helped stabilize learning.

**Conclusion**:
 Slight overfitting appears late in training. While performance is decent, the model doesn't generalize as well as deeper networks.
 **Needs regularization or early stopping** for better performance.

### ◆ GoogLeNet (Optimized)

- Both training and validation losses are **low and closely aligned**, indicating stable training.
- Validation accuracy reaches **~87%**, second only to ResNet152.
- Auxiliary classifiers in Inception help gradients flow and avoid vanishing.

**Conclusion**:
 **Very stable training**, no overfitting. Strong generalization, good trade-off

between complexity and accuracy.
🔵 **Recommended model** for balance.

### ◆ EfficientNet-B1

- Loss curves are smooth, validation loss slightly diverges but not severely.
- Accuracy stabilizes around **81%**.
- Moderate overfitting but manageable with regularization.

**Conclusion**:

Balanced model, slightly behind GoogLeNet and ResNet152. Could improve with better data augmentation.
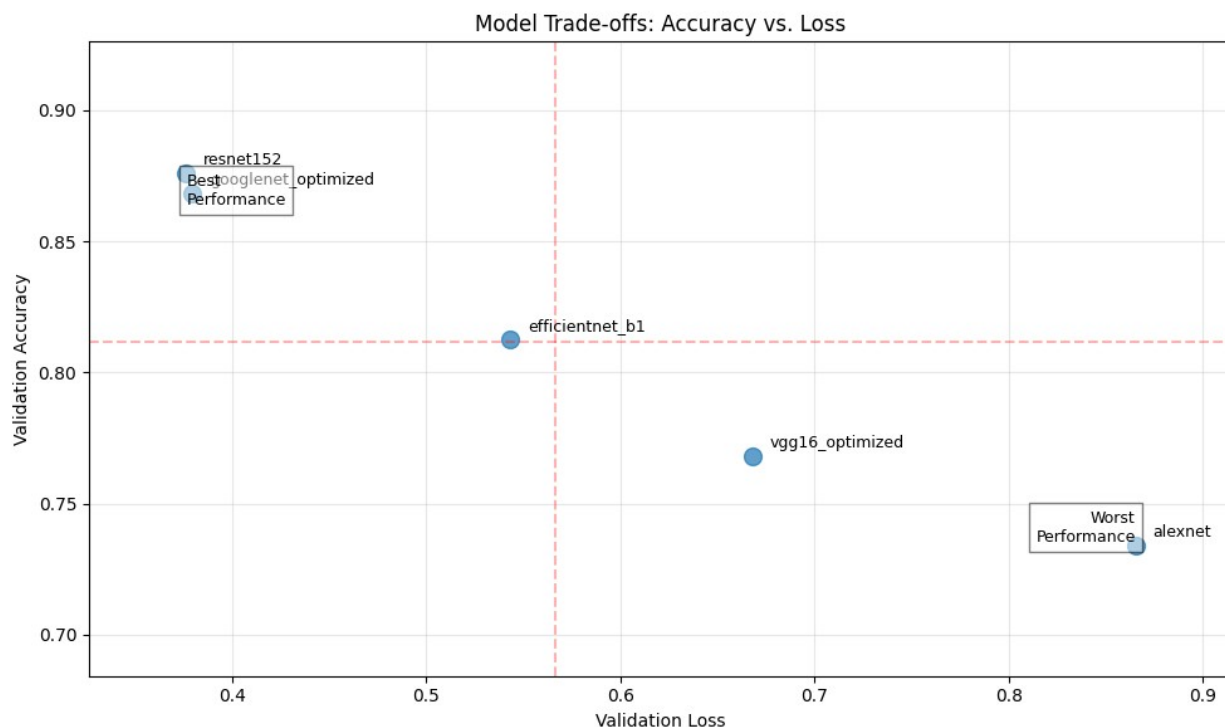
**Good, but not best**.

### ◆ ResNet152

- **Training and validation losses decrease together**, closely aligned — no sign of overfitting.
- **Best validation accuracy** (0.8762) among all.
- Consistent and **fast convergence** (by epoch 9).

**Conclusion**:

**Best performing model**: no overfitting, highest accuracy, stable convergence.

**Justified choice based on both training behavior and metrics**.

Model Trade-offs: Accuracy vs. Loss

# Final Model Selection: ResNet152

## Justification:

| Criterion | ResNet152 |
|---|---|
| Validation Accuracy | **Highest (87.62%)** |
| Training vs. Validation Loss | **Closely matched** (no overfitting) |
| Epochs to converge | **Only 9 epochs** (efficient training) |
| Stability | **Consistent and smooth curves** |
| Performance on Plots | **Top-right (best quadrant) on trade-off chart** |

## Summary:

- **ResNet152** clearly stands out as the best model due to its **training dynamics**, **generalization ability**, and **final performance metrics**.
- **GoogLeNet** is a strong second choice.
- **AlexNet** overfits and fails to generalize.