# Bangladesh University of Business and Technology (BUBT)

## PROJECT REPORT

**Project Name:** ASL Gesture Recognition

### Supervised By

Name: Sondip Poul Singh
Department of: CSE
Assistant Professor
Bangladesh University of Business & Technology

### Submitted By

Name: Md. Famidul Islam Pranto
ID No: 21225103054
Intake: 49
Section: 04
Course Title:  Neural Network & Fuzzy Systems Lab
Course Code: CSE 478

## Dataset
I used Kaggle **American Sign Language Dataset,**

https://www.kaggle.com/datasets/ayuraj/asl-dataset

The dataset contains 36 classes, which include the 26 alphabets (a-z) and 10 digits (0-9).
The data was split into 80% training, 10% validation, and 10% testing sets with the following image counts:
• Training: 2012 images
• Validation: 251 images
• Testing: 252 images
All images were resized to 224x224 pixels for uniformity.
**My Notebook:** https://www.kaggle.com/code/fipro054/asl-gesture-recognition

## Approach
I experimented with three different models to classify the hand gestures:
• A basic CNN model built from scratch, using Conv2D, MaxPooling, BatchNormalization, and Dense layers
• Transfer learning using **DenseNet121** pretrained on ImageNet
• Transfer learning using **EfficientNetB0** pretrained on ImageNet

## Data Preprocessing and Augmentation
For the training data, I applied data augmentation techniques including rotation, width and height shifts, zoom, and horizontal flips to increase variability but it caused issues like two different class became same sometime ,which made the model confused so I removed them.
For all the model, I rescaled the pixel values to be between 0 and 1.

## Training
I used two callbacks to optimize training:
• EarlyStopping to stop training if validation loss didn't improve for 5 consecutive epochs
• ReduceLROnPlateau to reduce the learning rate when validation accuracy plateaued
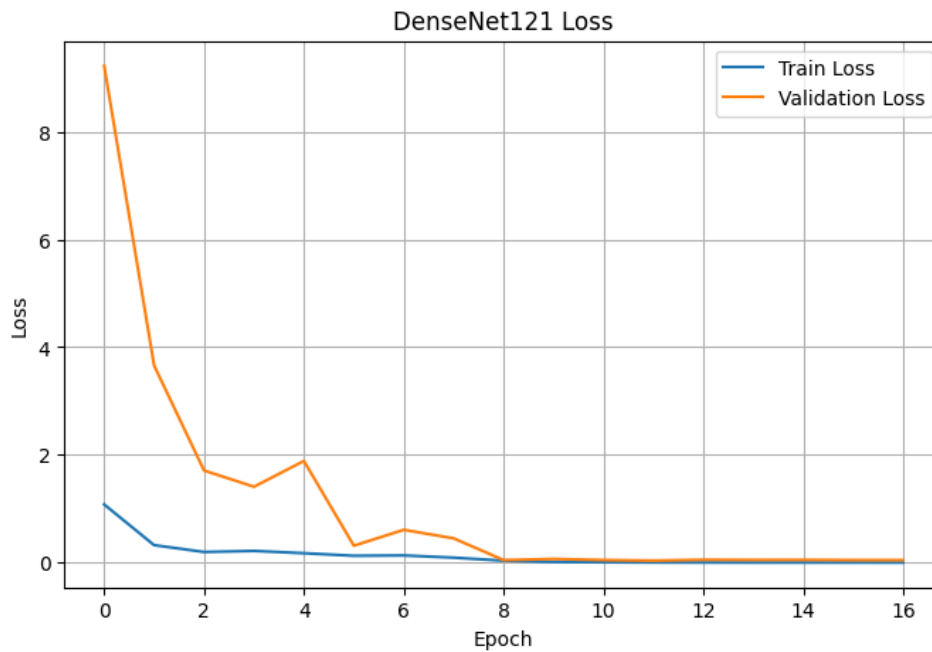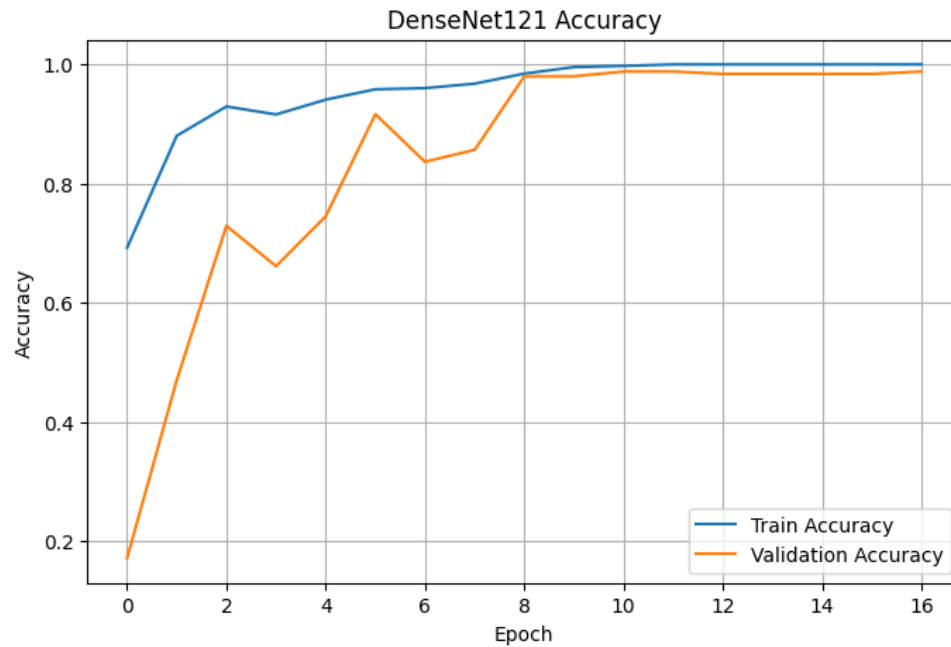All models were trained with a batch size of 32 for up to 20 epochs.

## Results
Among the three models evaluated, **DenseNet121 outperformed both the basic CNN and EfficientNetB0** on the ASL dataset. This superior performance is likely due to DenseNet's ability to efficiently reuse features through dense connections, which helps it generalize well even on relatively simple datasets with plain backgrounds.
The **DenseNet model achieved the highest test accuracy of 96.83%**, making it the best-performing model in this study. The **basic CNN model also performed well**, achieving a test accuracy of **94.05%**, which demonstrates that a simpler architecture can still be highly effective on clean, well-structured datasets. In contrast, **EfficientNetB0 significantly underperformed**, achieving only **2.78% accuracy**, likely due to suboptimal preprocessing or architectural mismatches.
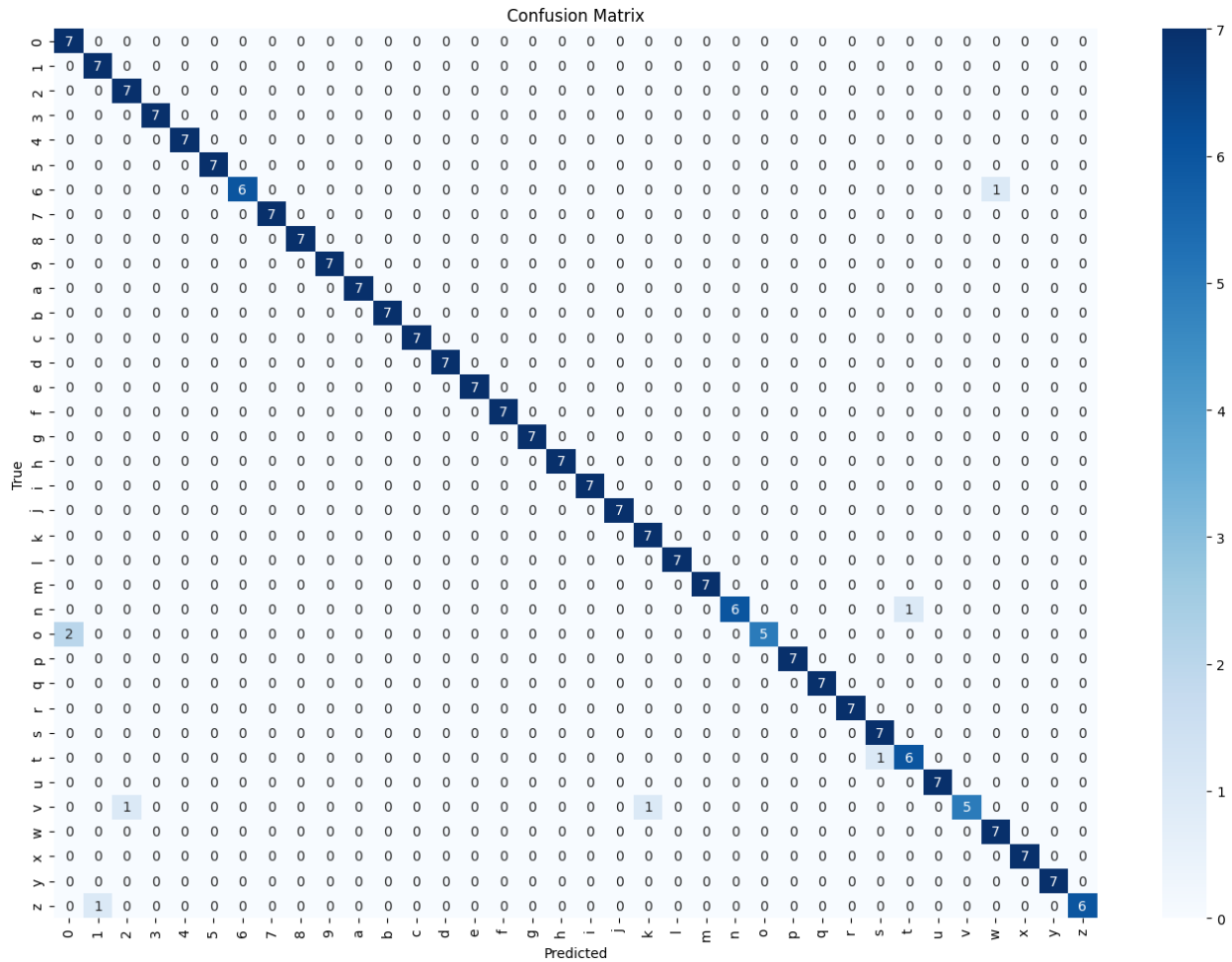
## Visualizations

I visualized the DenseNet model's training using accuracy and loss curves, which showed consistent performance across epochs. To evaluate class-wise predictions, I also created a confusion matrix. It helped reveal which classes the model confused the most.

**Training and Validation Accuracy/Loss Curves**

**Confusion Matrix**



Confusion Matrix

## Conclusion

In this project, I developed and evaluated several models for ASL hand gesture classification. Among them, DenseNet121 delivered the highest performance, outperforming both the basic CNN and EfficientNetB0 models. This suggests that deeper transfer learning models can effectively capture complex patterns, especially when paired with clean and well-structured datasets. With the right architecture and training approach, high accuracy is achievable even in relatively straightforward classification tasks.