

# Project Name

## Alphabet Classification

**Teaching-Assistant:** Eng/ Ahmed Ashraf

### Inception V1 — Advanced Inception Architecture

**Reference:** Szegedy et al., 2015/2016 (Rethinking

**Overview / Core Concept:** Inception V1 builds on the original Inception (v1) by factorizing convolutions (e.g., replacing  $5 \times 5$  convs with two  $3 \times 3$ , or  $n \times n$  with  $1 \times n$  followed by  $n \times 1$ ), incorporating Batch Normalization, label smoothing, and RMSProp optimization. These improvements achieve higher accuracy without significantly increasing computation.

#### **Forward Pass Flow (High-Level)**

**Stem:** Optimized initial convolution layers reduce spatial dimensions efficiently.

**Inception Modules:** Multiple branches with factorized convolutions (e.g.,  $1 \times n$  followed by  $n \times 1$ ) to approximate larger kernels efficiently.

**Auxiliary Classifiers:** Often removed or simplified; aggressive regularization and batch normalization are applied throughout.

**Classification Head:** Global Average Pooling  $\rightarrow$  Dropout  $\rightarrow$  Fully Connected Layer  $\rightarrow$  Softmax.

## Key Details / Hyperparameters

Factorized and asymmetric convolutions reduce computation cost

(e.g.,  $1 \times 7 \rightarrow 7 \times 1$ ).

Training strategies such as batch normalization and label smoothing improve convergence.

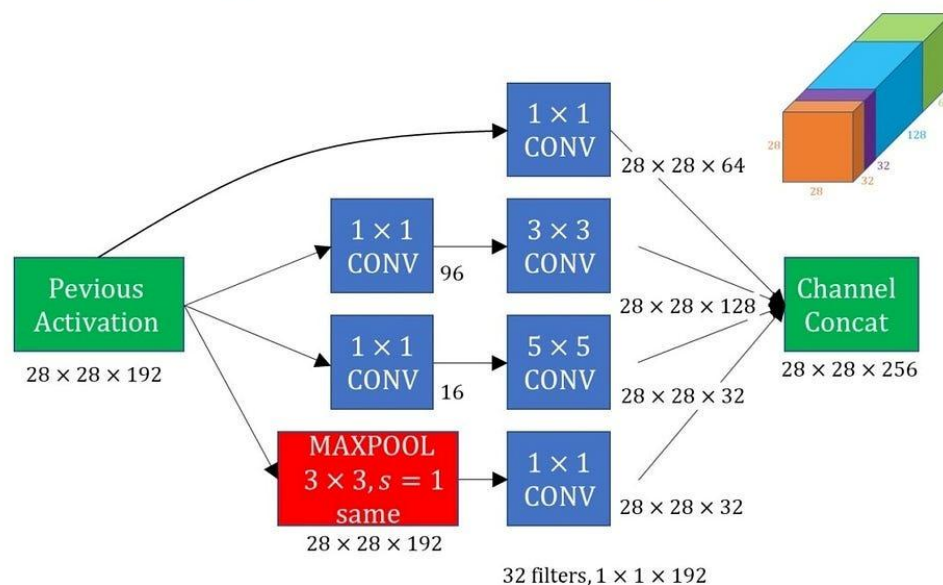
Inception V1 is widely used as a high-performance backbone for image classification.

## Advantages and Limitations

**Pros:** Excellent accuracy-to-computation ratio; widely adopted in practice.

**Cons:** More complex architecture to implement and tune than standard residual networks; numerous hyperparameters and module variants make it harder to modify.

## An example of an Inception module



## 2] VGG19 — Very Deep Convolutional Network

---

**Reference:** Simonyan & Zisserman, 2014

### **Overview / Core Concept :**

**VGG employs a straightforward and consistent**

**architecture:** sequences of  $3 \times 3$  convolutional layers (activated by ReLU) interspersed with occasional  $2 \times 2$  max-pooling layers. The network's depth is increased to enhance its feature representation capabilities . The VGG19, model specifically contains 19 layers with learnable weights composed of 16 convolutional layers followed by 3 fully connected layers.

### **Forward Pass Flow (VGG19)**

**Input:** Image of size  $224 \times 224 \times 3$  (RGB).

**Block 1:** Two conv layers ( $3 \times 3$ , 64 filters)  $\rightarrow$  max-pool ( $2 \times 2$ , stride 2)  
 $\rightarrow$  output spatial dimensions halved.

**Block 2:** Two conv layers ( $3 \times 3$ , 128 filters)  $\rightarrow$  max-pool  $\rightarrow$  output halved.

**Block 3:** Four conv layers ( $3 \times 3$ , 256 filters)  $\rightarrow$  max-pool.

**Block 4:** Four conv layers ( $3 \times 3$ , 512 filters)  $\rightarrow$  max-pool.

**Block 5:** Four conv layers ( $3 \times 3$ , 512 filters)  $\rightarrow$  max-pool.

**Classification Head:** Flatten  $\rightarrow$  FC (4096)  $\rightarrow$  ReLU  $\rightarrow$  Dropout  $\rightarrow$  FC (4096)  $\rightarrow$  ReLU  $\rightarrow$  Dropout  $\rightarrow$  FC (1000)  $\rightarrow$  Softmax.

## Key Details / Hyperparameters:

All convolution layers use  $3 \times 3$  kernels, stride 1, and padding 1.

Max-pooling uses  $2 \times 2$  windows with stride 2.

Fully connected layers are large and account for the majority of the network's parameters.

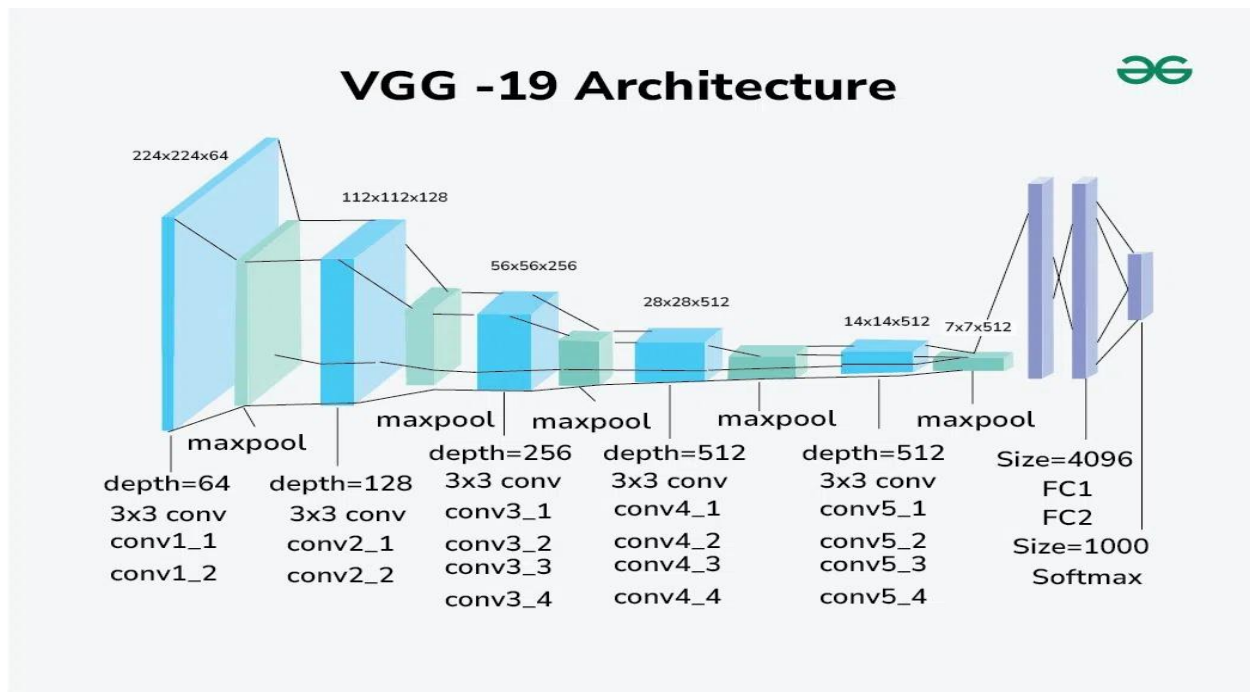
Typical input resolution:  $224 \times 224$ .

Total parameters: ~144 million (depends on exact FC layer configuration).

## Advantages and Limitations:

**Pros:** Simple and uniform design, strong performance baseline, easy to implement and fine-tune for transfer learning.

**Cons:** Very heavy in terms of parameters and memory; computationally slow; not ideal for mobile deployment or extremely deep scaling



### 3] ResNet50 — Deep Residual Network

---

**Reference:** He et al., 2015 (Deep Residual Learning)

#### **Overview / Core Concept:**

ResNet introduces residual (skip) connections, which learn a residual mapping  $F(x) = H(x) - x$

$F(x)=H(x)-x$  instead of directly learning  $H(x)$

$H(x)$ . This design allows very deep networks (e.g., 50, 101, 152 layers) to be trained effectively by mitigating the vanishing gradient problem.

#### **Forward Pass Flow (ResNet50)**

**Input:** Image of size  $224 \times 224 \times 3$ .

**Initial Layer:** Conv ( $7 \times 7$ , stride 2)  $\rightarrow$  BatchNorm  $\rightarrow$  ReLU  $\rightarrow$  MaxPool ( $3 \times 3$ , stride 2).

**Stage 1 (conv2\_x):** 3 residual bottleneck blocks, each composed of:

$1 \times 1$  conv (channel reduction)  $\rightarrow$  BatchNorm  $\rightarrow$  ReLU

$3 \times 3$  conv (spatial)  $\rightarrow$  BatchNorm  $\rightarrow$  ReLU

$1 \times 1$  conv (channel expansion)  $\rightarrow$  BatchNorm

Add input via identity or projection shortcut  $\rightarrow$  ReLU

**Stage 2 (conv3\_x):** 4 bottleneck blocks.

**Stage 3 (conv4\_x):** 6 bottleneck blocks.

**Stage 4 (conv5\_x):** 3 bottleneck blocks.

**Classification Head:** Global Average Pooling → Fully Connected Layer → Softmax for class probabilities.

### Key Details / Hyperparameters:

Bottleneck structure reduces computation while maintaining network depth.

Typical ResNet50 has ~25 million parameters.

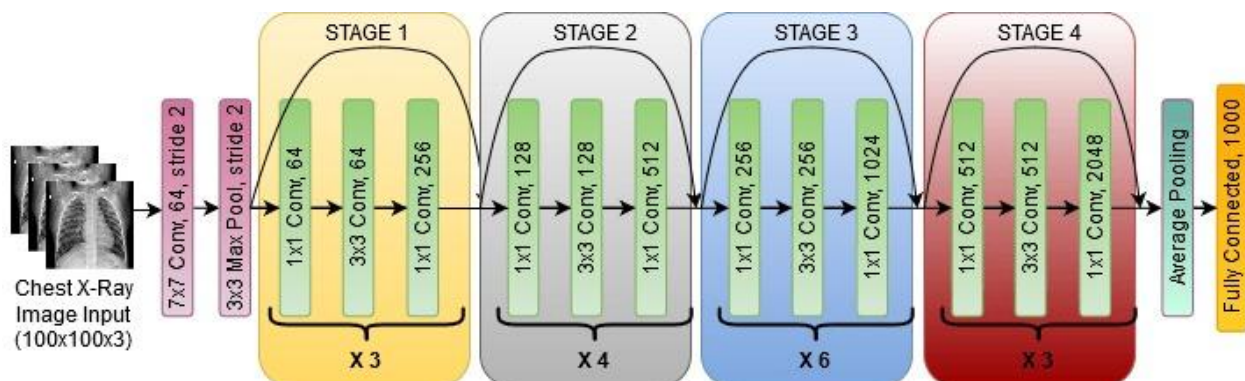
Uses Batch Normalization and He initialization.

**Shortcuts:** Identity connections are used when input/output dimensions match; projection shortcuts handle dimension changes.

### Advantages and Limitations

**Pros:** Easily trainable for very deep networks, strong accuracy-to-computation ratio, widely adopted as a backbone for many tasks.

**Cons:** Still relatively heavy compared to lightweight mobile networks; residual connections may require task-specific adjustments (e.g., in detection or segmentation).



## 4] MobileNet — Efficient Networks with Depthwise Separable Convolutions

---

**Reference:** Howard et al., 2017 (MobileNets)

### **Overview / Core Concept:**

MobileNet replaces standard convolution layers with depthwise separable convolutions, which split computation into two steps:

**Depthwise convolution:** performs spatial convolution independently on each input channel.

**Pointwise convolution ( $1 \times 1$ ):** combines channels across depth.

This decomposition significantly reduces both FLOPs and the number of parameters. Two main hyperparameters control the model size and accuracy: width multiplier ( $\alpha$ ) and resolution multiplier ( $\rho$ ).

### **Forward Pass Flow (MobileNet v1)**

**Input:** e.g.,  $224 \times 224 \times 3 \rightarrow$  Initial  $3 \times 3$  convolution (stride 2).

**Repeated Blocks:**  $N \times (\text{Depthwise } 3 \times 3 \text{ conv} \rightarrow \text{BatchNorm} \rightarrow \text{ReLU} \rightarrow \text{Pointwise } 1 \times 1 \text{ conv} \rightarrow \text{BatchNorm} \rightarrow \text{ReLU})$ . Stride can be 1 or 2 depending on the block.

**Output:** Global Average Pooling  $\rightarrow$  Fully Connected Layer  $\rightarrow$  Softmax for classification.

## Hyperparameters:

Width multiplier  $\alpha \in (0,1]$

$\alpha \in (0,1]$  scales channel numbers.

Resolution multiplier  $\rho$

$\rho$  adjusts input spatial dimensions.

## Key Details / Hyperparameters

MobileNet v1 typically uses ReLU (or ReLU6) and Batch Normalization.

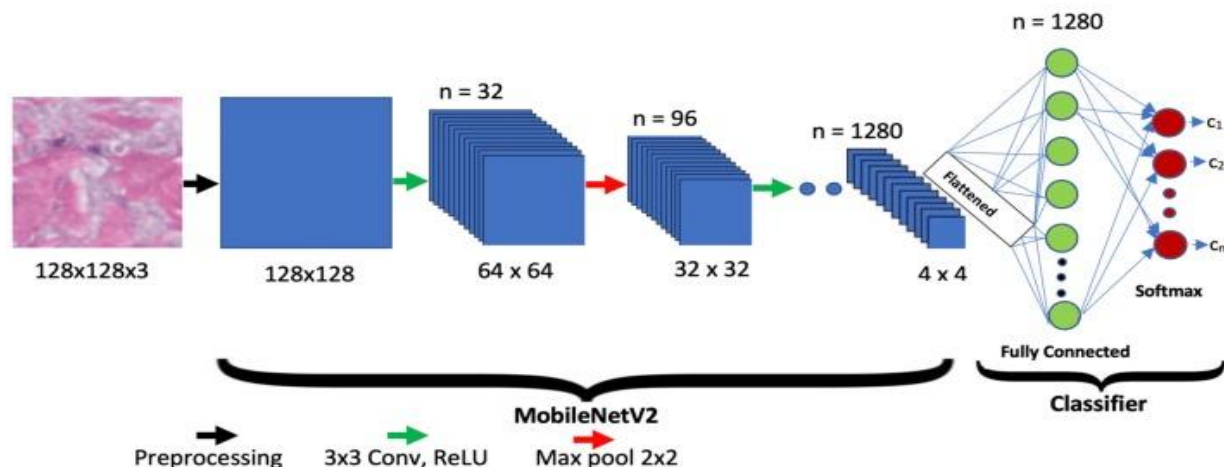
$\alpha = 1.0$  corresponds to the full model;  $\alpha < 1.0$  produces a smaller, lighter version for mobile applications.

Designed to trade off accuracy for latency and memory efficiency.

## Advantages and Limitations

**Pros:** Extremely lightweight and fast, ideal for mobile and embedded devices.

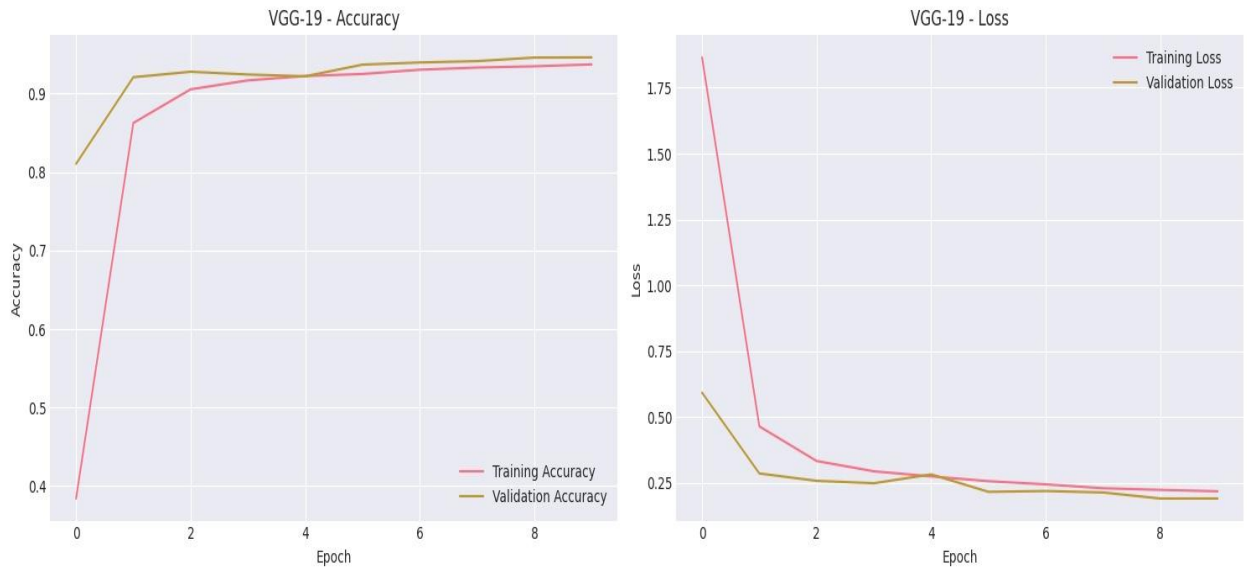
**Cons:** Lower accuracy compared to large standard CNNs; later versions (MobileNetV2/V3) improved efficiency using inverted residual blocks and other enhancements



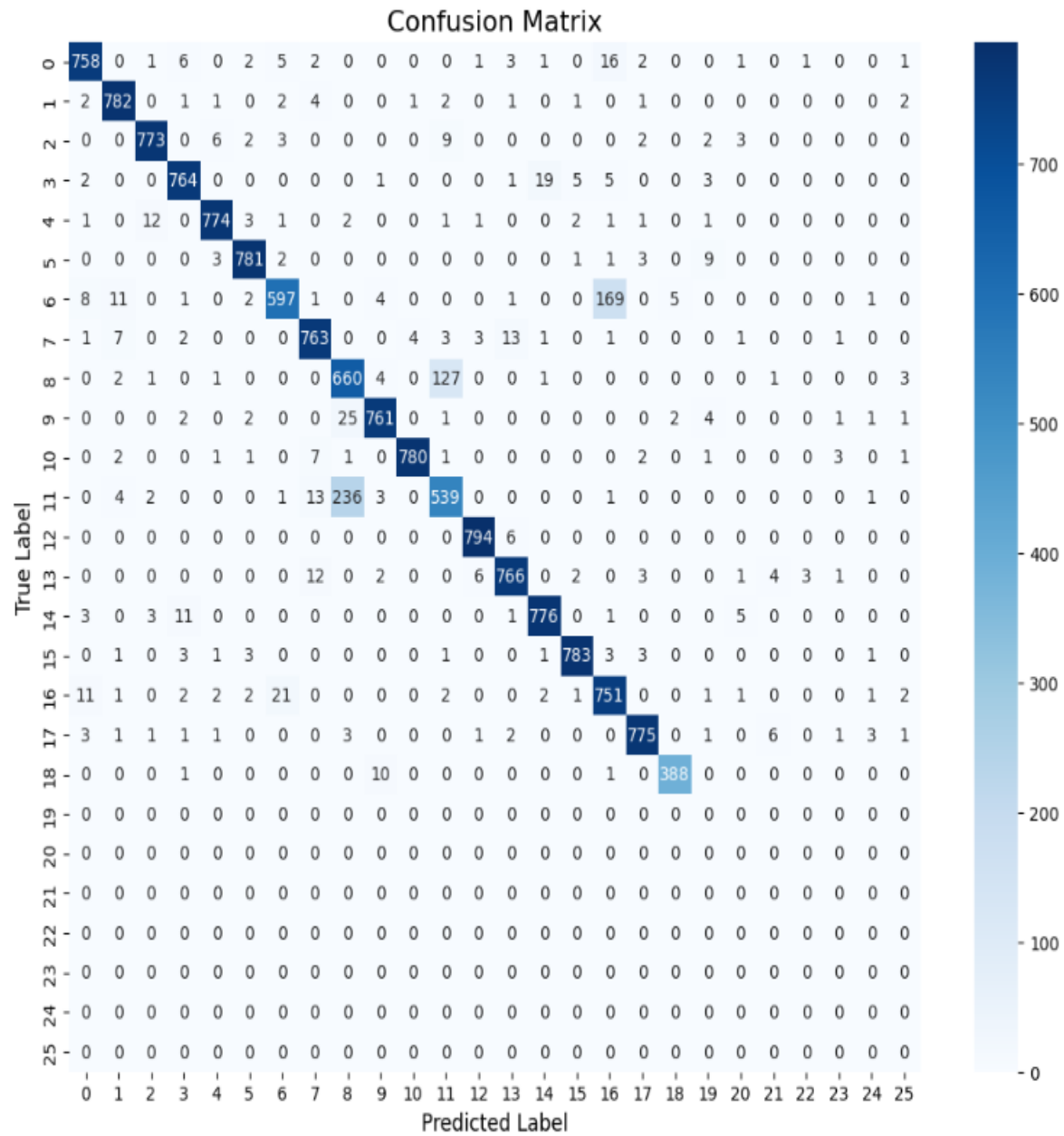


## Results for each model:

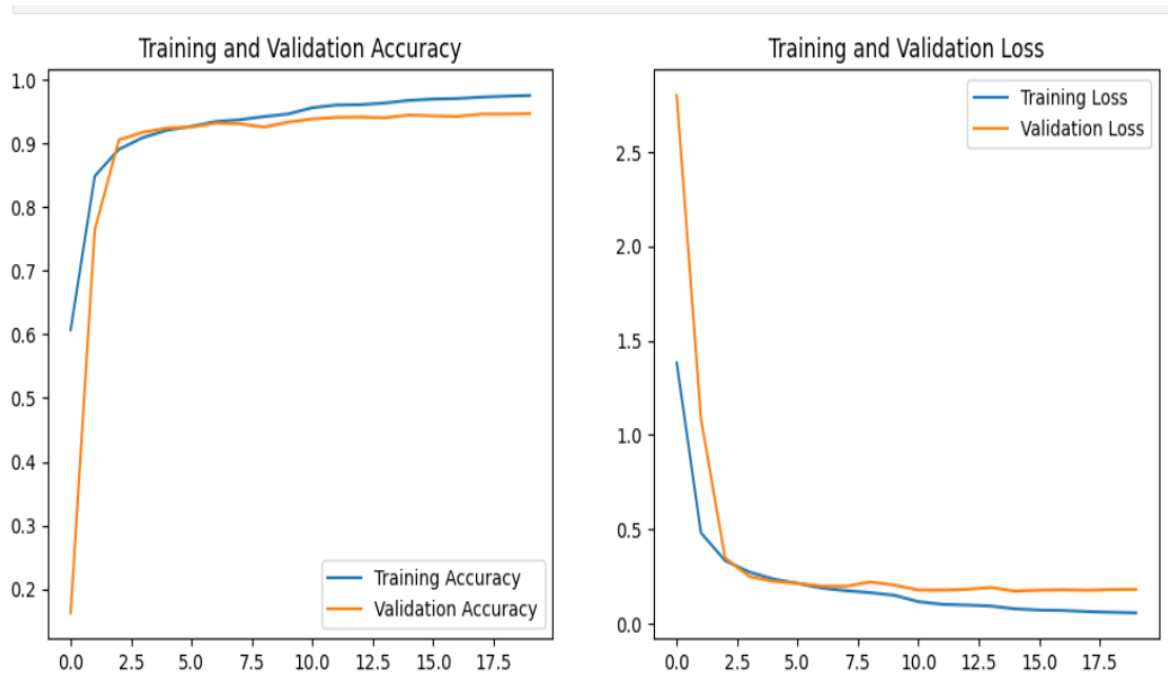
### 1) VGG19 :



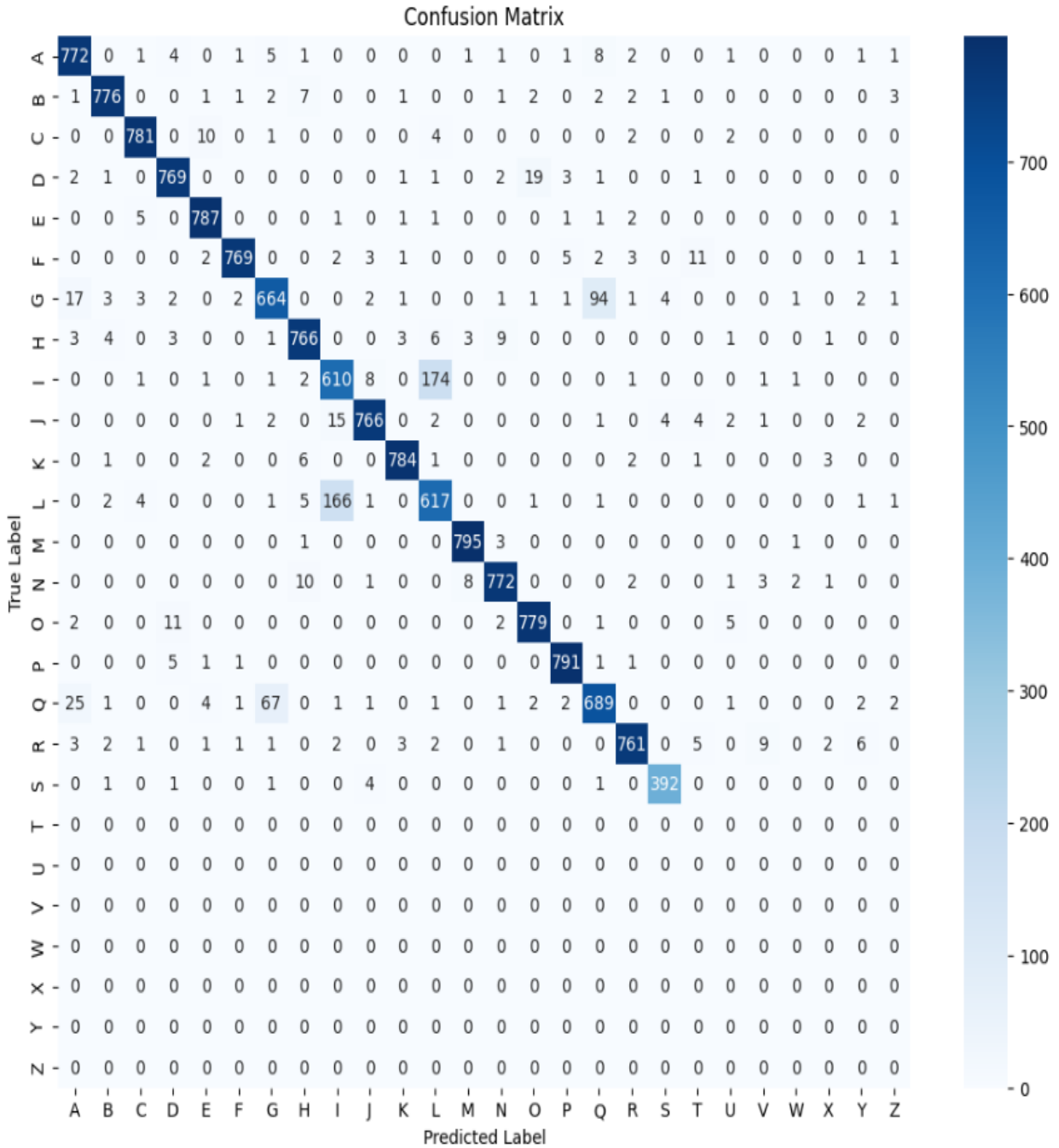
**VGG -> accuracy: accuracy: 0.9371 - loss: 0.2178 - val\_accuracy: 0.9461 - val\_loss: 0.1904 - learning\_rate: 1.0000e-04**



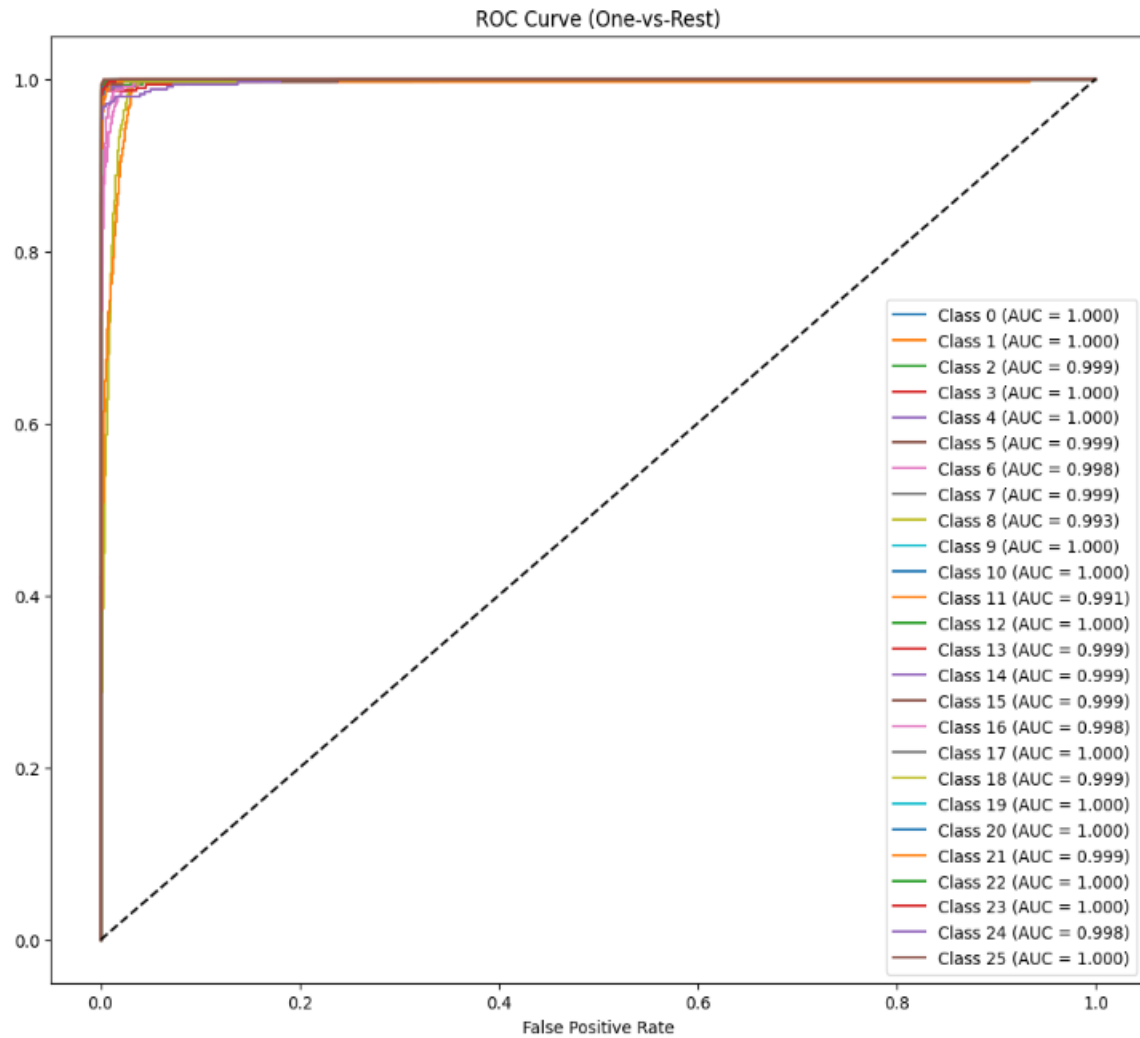
## 2) MobileNet:



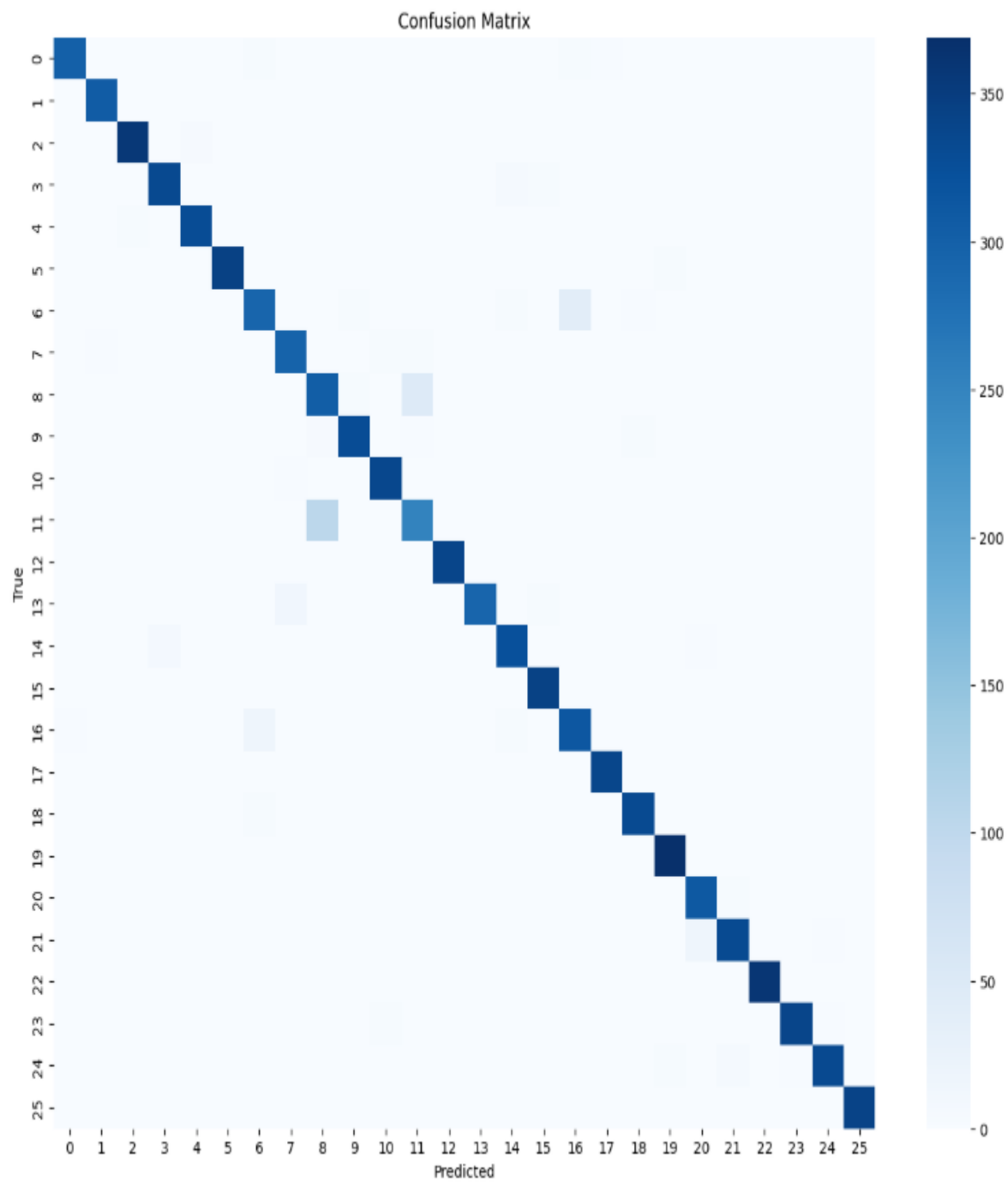
**MobileNet : accuracy: 0.9748 - loss: 0.0573 - val\_accuracy: 0.9467 –  
val\_loss: 0.1814 - learning\_rate: 6.2500e-06**



### 3) Inception V1:



**Inception V1: Epoch 10: Train Acc=0.977, Val Acc=0.947**



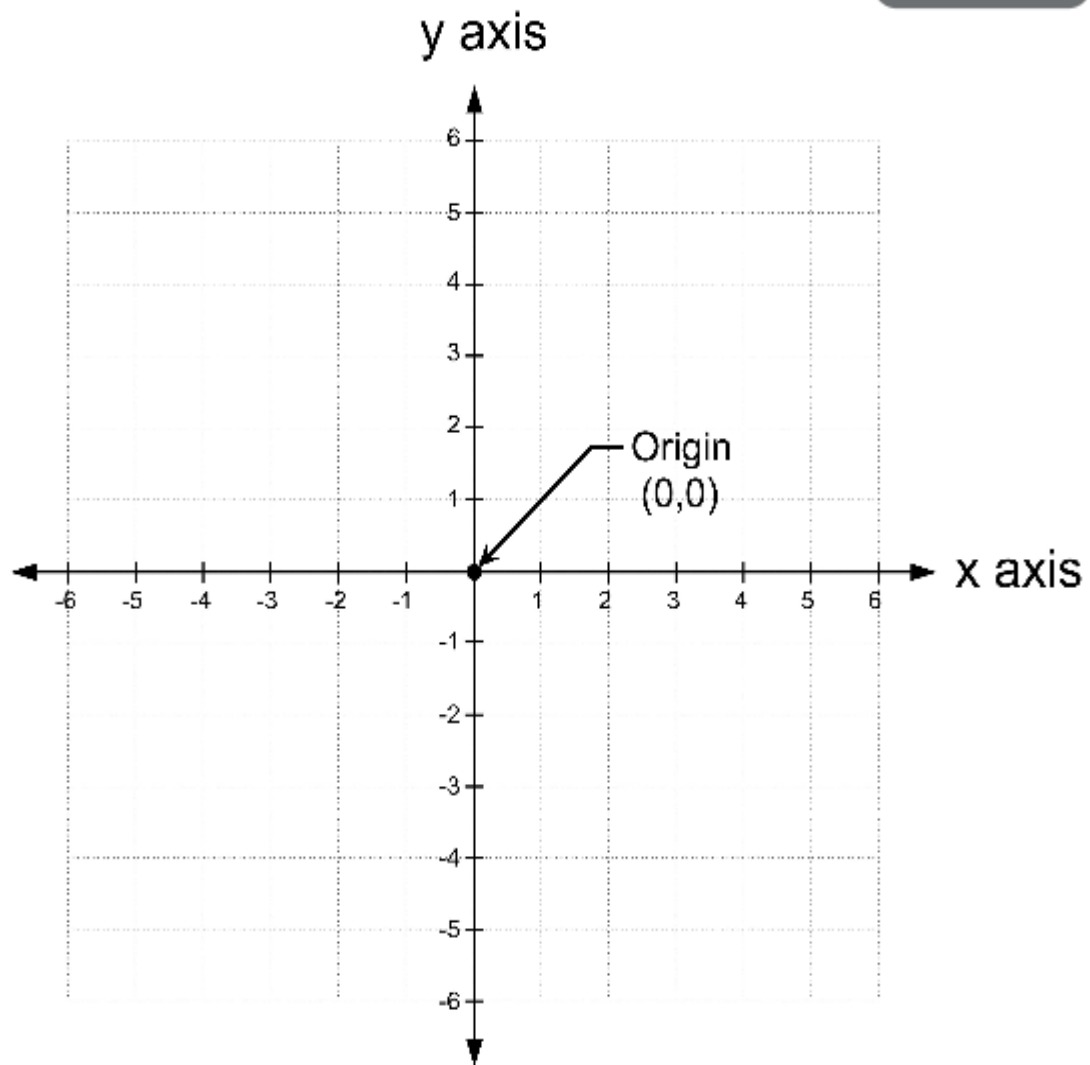
## comparison :

---

| Model Architecture                                      | Training Approach                            | Test Accuracy                 | Test Loss                      | Key Observations  |
|---|--|-------------------------------|--------------------------------|---|
| <b>VGG-style CNN</b><br>(Custom)                        | Trained from scratch                         | 0.9420<br>(Val Acc)           | 0.1652<br>(Val Loss)           | High performance, fast training due to custom, smaller architecture. Excellent generalization ( Overfitting gap $\approx 0.015$ ).  |
| <b>ResNet50</b><br>(Transfer Learning)                  | Feature Extraction<br>(Frozen Base)          | 0.8889<br>(Val Acc)           | 0.3442<br>(Val Loss)           | Significantly lower performance than VGG/MobileNet. Scaling up to 224x224 and conversion to 3 channels likely introduced artifacts/complexity that the frozen base could not handle optimally for this domain.                        |
| <b>MobileNetV2</b><br>(Transfer Learning + Fine-Tuning) | Fine-Tuning after initial feature extraction | 0.9351                        | 0.2081                         | Strong performance, very fast inference due to efficient architecture. Fine-tuning provided a slight improvement over initial results.  |
| <b>Vision Transformer (ViT)</b> (Trained from scratch)  | Trained from scratch with Data Augmentation  | 0.8900<br>(Val Acc, Epoch 12) | 0.3124<br>(Val Loss, Epoch 12) | Started slowly but showed strong learning ability. Performance is comparable to basic transfer learning models, but lower than the custom CNN/MobileNet, possibly due to limited training time and dataset size for a full ViT model. |
| <b>Inception V1/GoogLeNet</b><br>(Transfer Learning)    | Fine-Tuning<br>(Cut short due to OOM)        | (N/A)                         | (N/A)                          | Performance could not be fully evaluated due to <b>CUDA Out of Memory</b> errors, highlighting its high memory requirement on this hardware configuration.  |

## Pros and Cons of Each Architecture:

Shutterstock





## Pros and Cons of Model Architectures:

| Architecture                              | Pros (+)  | Cons (-)  |
|---|---|---|
| <b>VGG-style CNN</b><br>(Custom)          | <ul style="list-style-type: none"> <li>+ Very high accuracy on this dataset.</li> <li>+ Robust against overfitting (small generalization gap).</li> </ul>                   | <ul style="list-style-type: none"> <li>- More parameters than MobileNet (less efficient).</li> <li>- Design is relatively manual and less modular than ResNet.</li> </ul>                               |
| <b>ResNet50</b> (Transfer Learning)       | <ul style="list-style-type: none"> <li>+ Utilizes powerful pre-trained features (Transfer Learning).</li> <li>+ Skip connections mitigate vanishing gradients.</li> </ul>   | <ul style="list-style-type: none"> <li>- Low performance without full fine-tuning.</li> <li>- Requires input image resizing to 224x224, which alters the small 28x28 EMNIST images.</li> </ul>          |
| <b>MobileNetV2</b><br>(Transfer Learning) | <ul style="list-style-type: none"> <li>+ Extremely efficient (low number of parameters).</li> <li>+ Good balance between speed and accuracy.</li> </ul>                     | <ul style="list-style-type: none"> <li>- Requires input image conversion to 3 channels (RGB) for pre-trained weights.</li> <li>- Slightly lower peak accuracy than the custom VGG-style CNN.</li> </ul> |
| <b>Vision Transformer (ViT)</b> (Scratch) | <ul style="list-style-type: none"> <li>+ Excellent for capturing global relationships through self-attention.</li> <li>+ State-of-the-art for many vision tasks.</li> </ul> | <ul style="list-style-type: none"> <li>- Very high data requirement for training from scratch.</li> <li>- Slower convergence than CNNs on smaller datasets.</li> </ul>                                  |
| <b>Inception V1</b><br>(GoogLeNet)        | <ul style="list-style-type: none"> <li>+ High parameter efficiency using Inception Modules.</li> <li>+ Uses auxiliary classifiers for better gradient flow.</li> </ul>      | <ul style="list-style-type: none"> <li>- Very high memory footprint; failed due to <b>OOM error</b>.</li> <li>- Complex architecture, sensitive to hyperparameter changes.</li> </ul>                   |