

# 1. YOLO

## Dataset Description:

1. **Source: Combination of 2 datasets:**
- Dataset A: Provided by the project team; ~100 images were manually annotated using bounding boxes and gender labels (male / female).
  - Dataset B: Imported from Roboflow; pre-annotated with bounding boxes and gender classes.
2. **Annotation Format:**
- YOLOv11-compatible .txt files for each image.
  - Each line in a label file:
    - <class\_id> <x\_center> <y\_center> <width> <height>
3. **Object Classes:**
- 0 → Male
  - 1 → Female
4. **Image Resolution:**

Varied across datasets; resized to 640×640 during training with padding where necessary.

5. **Dataset Statistics:**

**Train:** 1872

**Validate:** 225

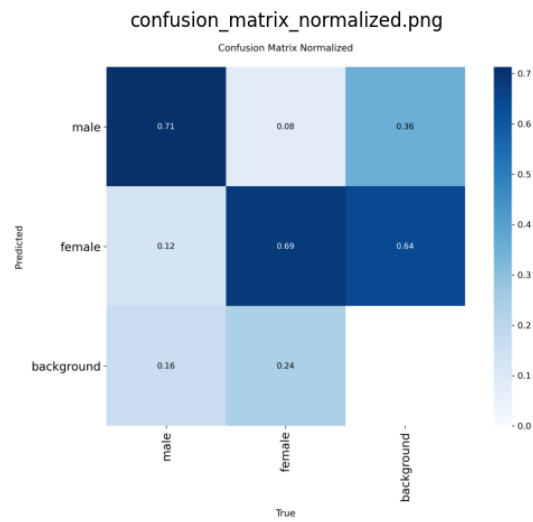
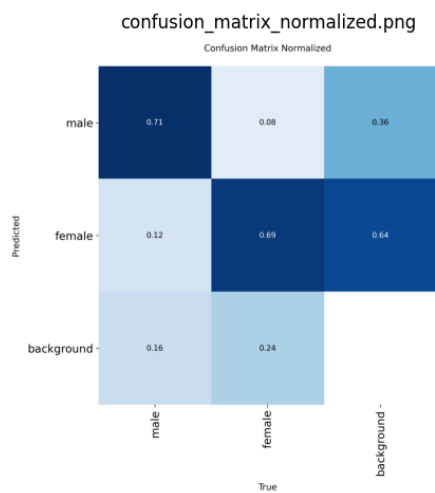
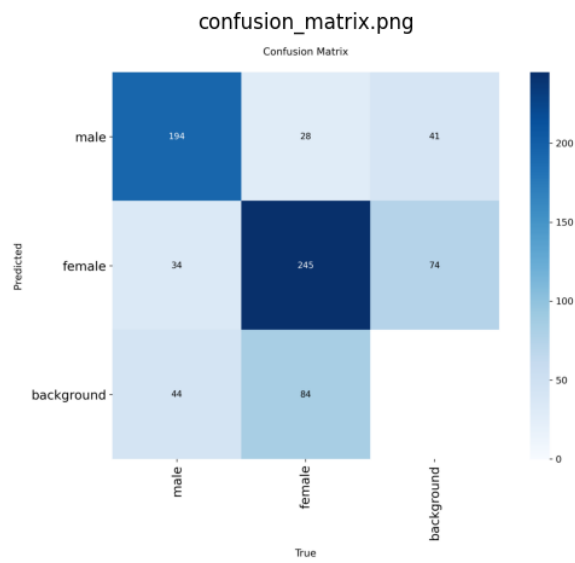
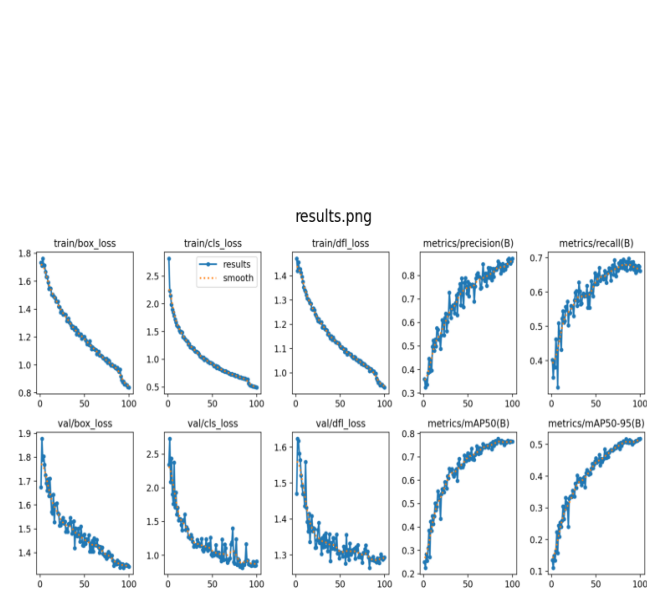
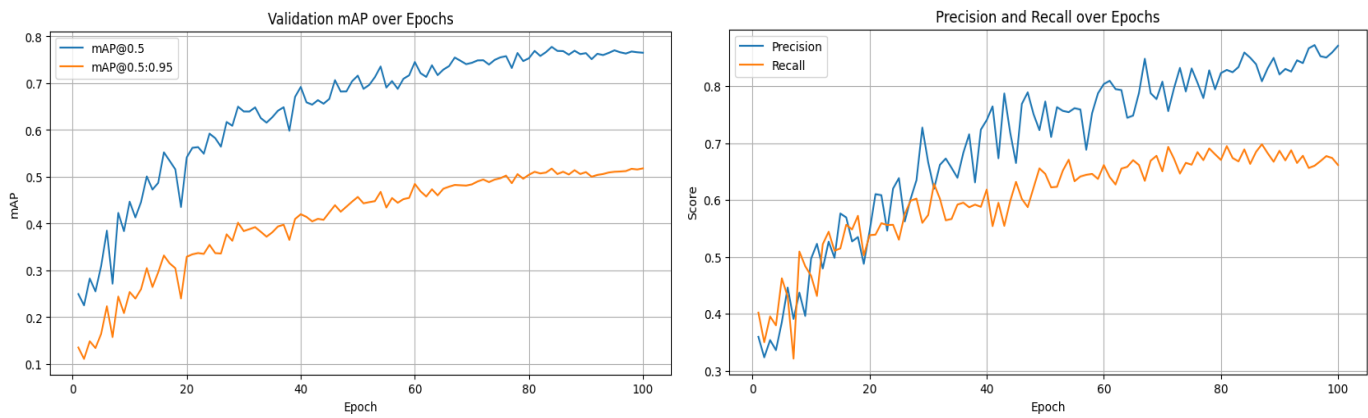
## Model Training Details:

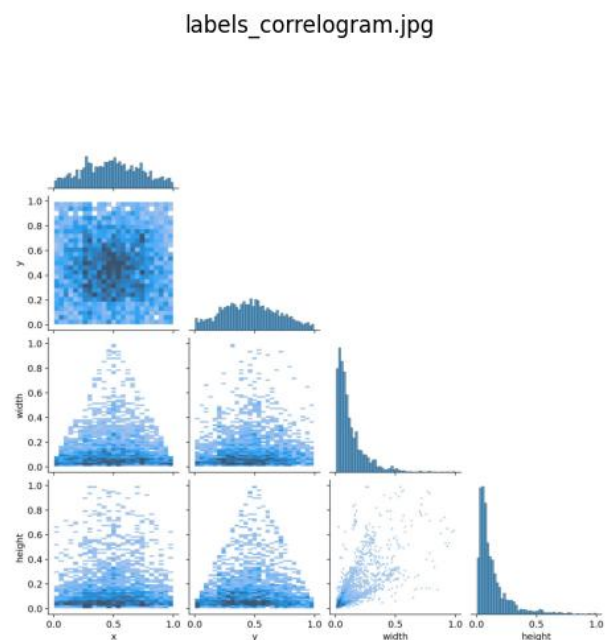
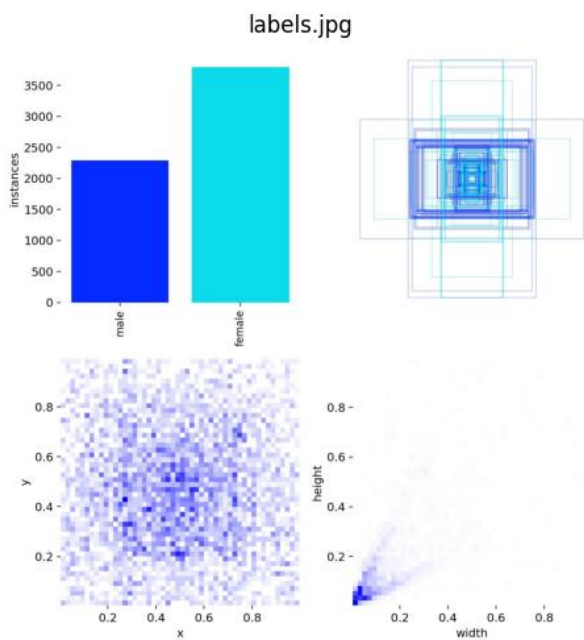
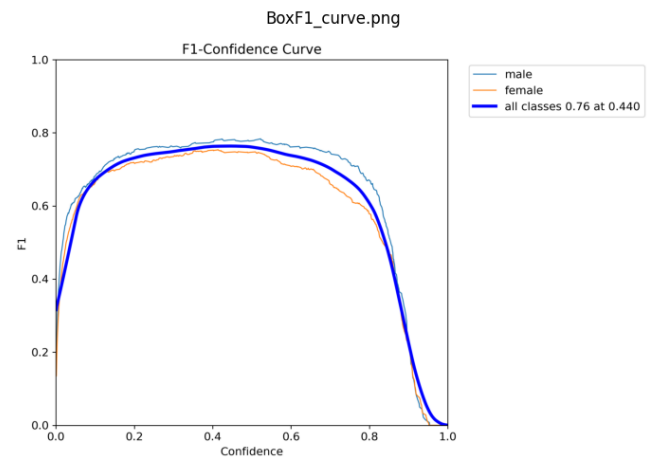
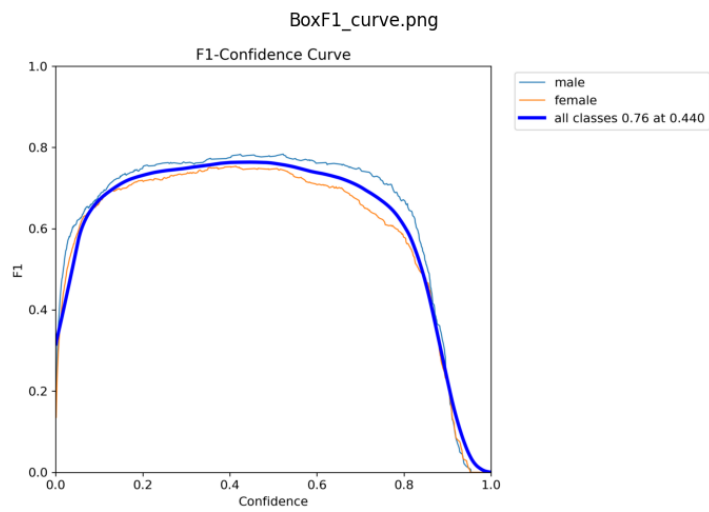
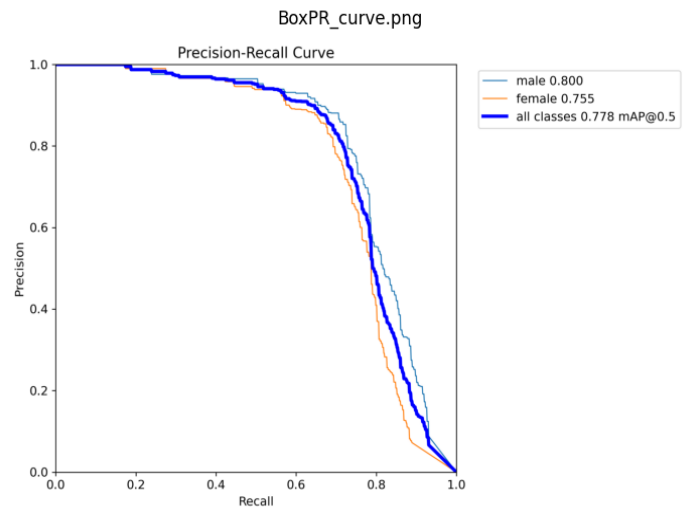
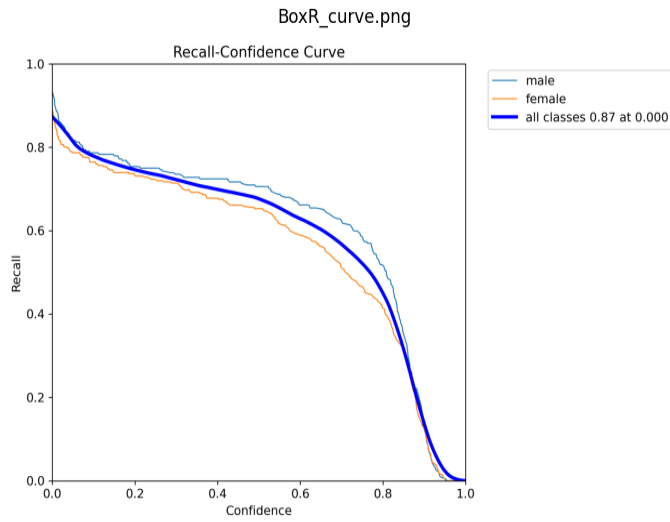
- Model Used:** YOLOv11n (YOLO version 11 nano)  
Chosen for its light weight and fast inference while maintaining good performance on small datasets.
  - Training Framework:** PyTorch using Ultralytics YOLOv11
  - Training Location:** Kaggle Notebook
  - Epochs:** 100
  - Image size:** 640×640
  - Batch size:** 16
  - Device:** GPU (Tesla T4, 16 GB VRAM)

## Training Performance Summary:

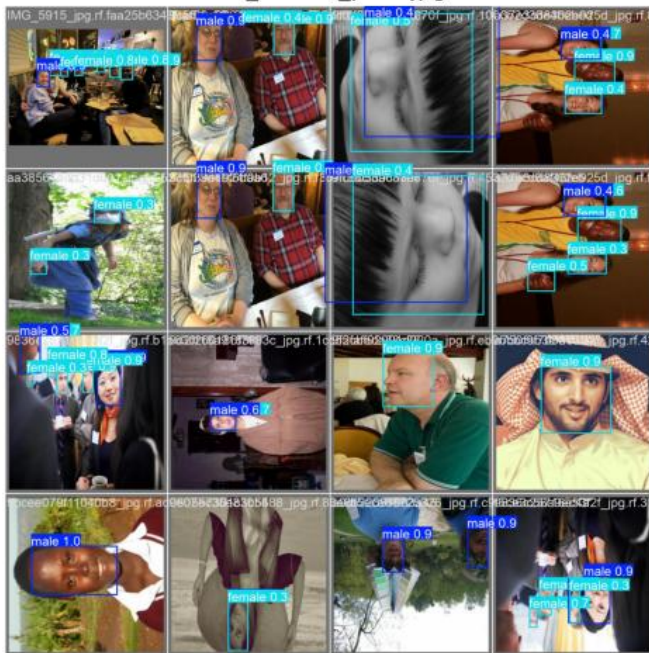
Class	Precision	Recall	mAP@0.5	mAP@0.5:0.95
All	0.857	0.689	0.778	0.518
Male	0.856	0.717	0.800	0.548
Female	0.858	0.661	0.755	0.487

# Training Results:

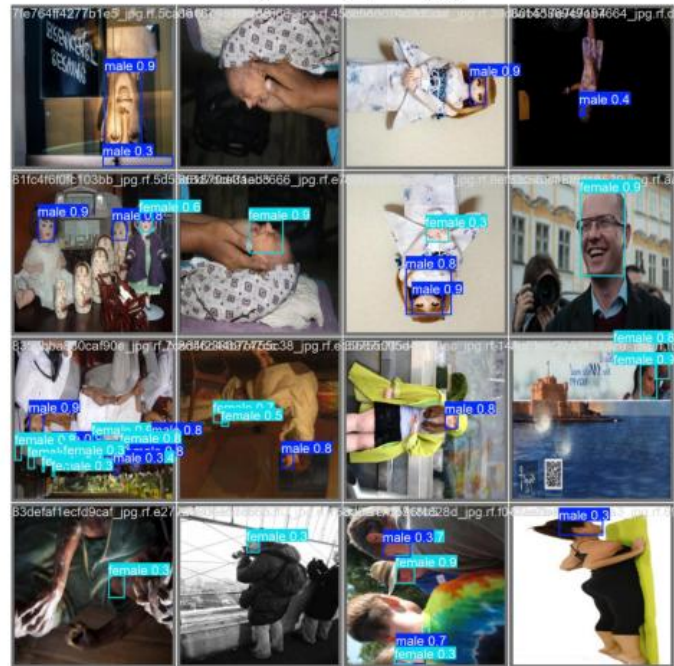




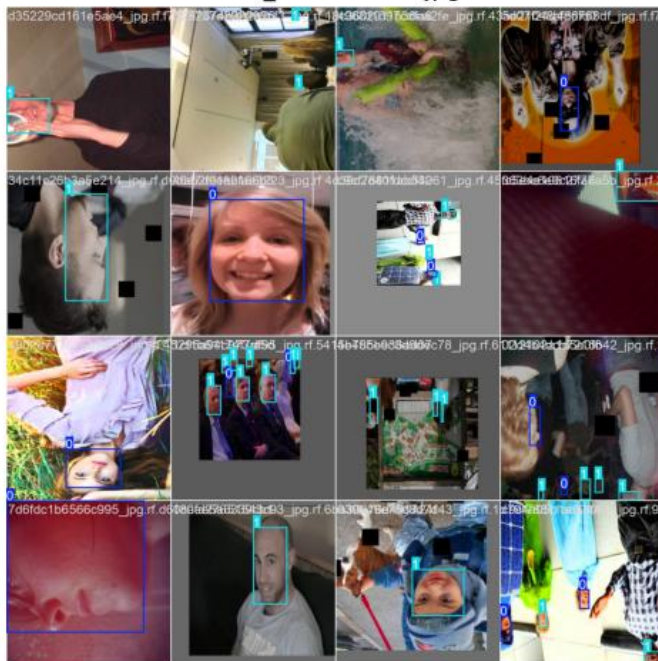
val\_batch0\_pred.jpg



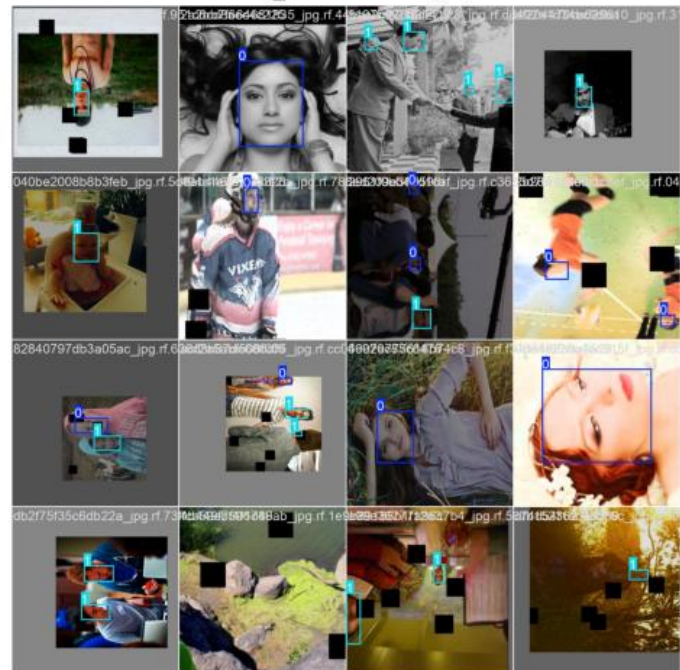
val\_batch1\_pred.jpg



train\_batch10532.jpg



train\_batch10531.jpg





## 2. Classifier

### Model Architecture:

Base Backbone: EfficientNet\_V2\_S, pre-trained on ImageNet1K (IMAGENET1K\_V1 weights).

- Custom Classifier Head:
- Dropout(0.5)
- Linear(1280  $\rightarrow$  512)
- ReLU
- Dropout(0.3)
- Linear (512  $\rightarrow$  16) (final output layer for 16 celebrity classes)

This structure balances transfer learning benefits with custom fine-tuning capacity for the specific celebrity face task.

### Dataset Details:

#### 1. Data Format:

- Data downloaded from Kaggle competition page.
- Images and labels loaded from .npy files (faces\_cropped.npy, labels\_cropped.npy)
- Images are already cropped face regions of shape (224, 224, 3), RGB.

#### 2. Label Remapping:

Labels with value -1 (used for unknowns) were reassigned to class index 15 to maintain label consistency within 0–15.

### Data Augmentations:

Data augmentation was applied using Albumentations to enhance generalization and robustness:

- HorizontalFlip (p=0.5)
- Rotate (limit=30°, p=0.5)
- HueSaturationValue (hue  $\pm 20$ , sat  $\pm 30$ , val  $\pm 20$ , p=0.7)
- RandomBrightnessContrast (brightness  $\pm 0.3$ , contrast  $\pm 0.3$ , p=0.7)
- CoarseDropout: 1–2 holes, each 10–20% of image area (black mask)
- Normalization: mean=[0.485, 0.456, 0.406], std=[0.229, 0.224, 0.225]
- Conversion to tensor with ToTensorV2()

### Training Strategy:

#### ➤ Phase 1: Train Classifier Head Only

- Freeze all layers except the classifier head.
- Optimizer: AdamW(lr=1e-3, weight\_decay=1e-4)
- Epochs: 100
- Loss Function: CrossEntropyLoss with MixUp
- Accuracy was monitored on a hold-out validation set (20%).

### ➤ Phase 2: Fine-Tune Last Feature Blocks + Head

- Only unfreeze the last 5 blocks of EfficientNet and the classifier.
- Optimizer: AdamW(lr=1e-4, weight\_decay=1e-5)
- Epochs: 50
- Same loss and sampling strategy used.

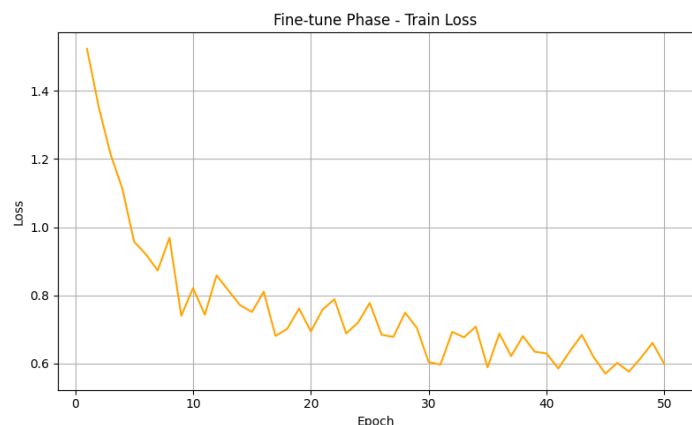
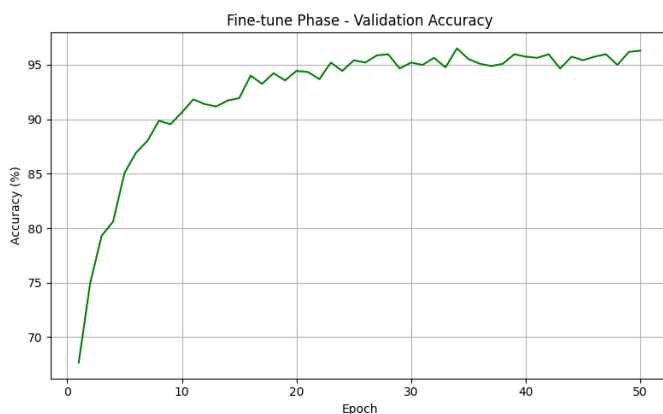
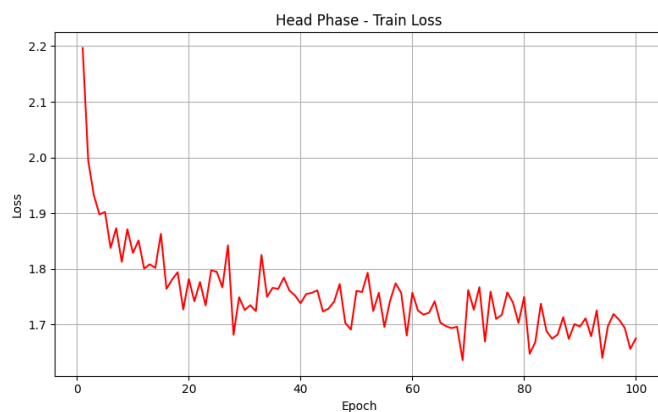
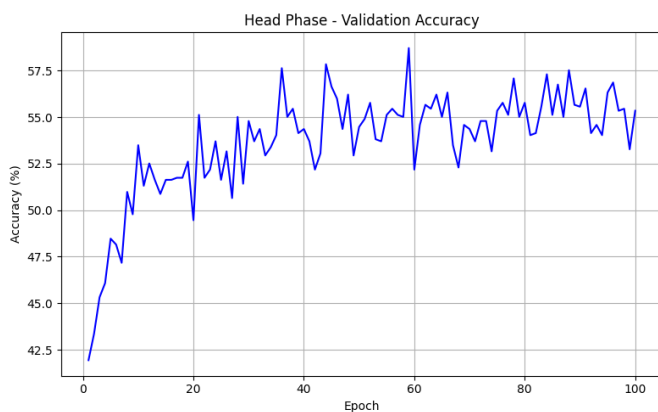
### ➤ MixUp Regularization

- MixUp was applied during training
- Hyperparameter: alpha=0.4 (Beta distribution)
- Loss computed as:

$$\text{lam} * \text{CE}(\text{preds}, y_a) + (1 - \text{lam}) * \text{CE}(\text{preds}, y_b)$$

## Training Results:

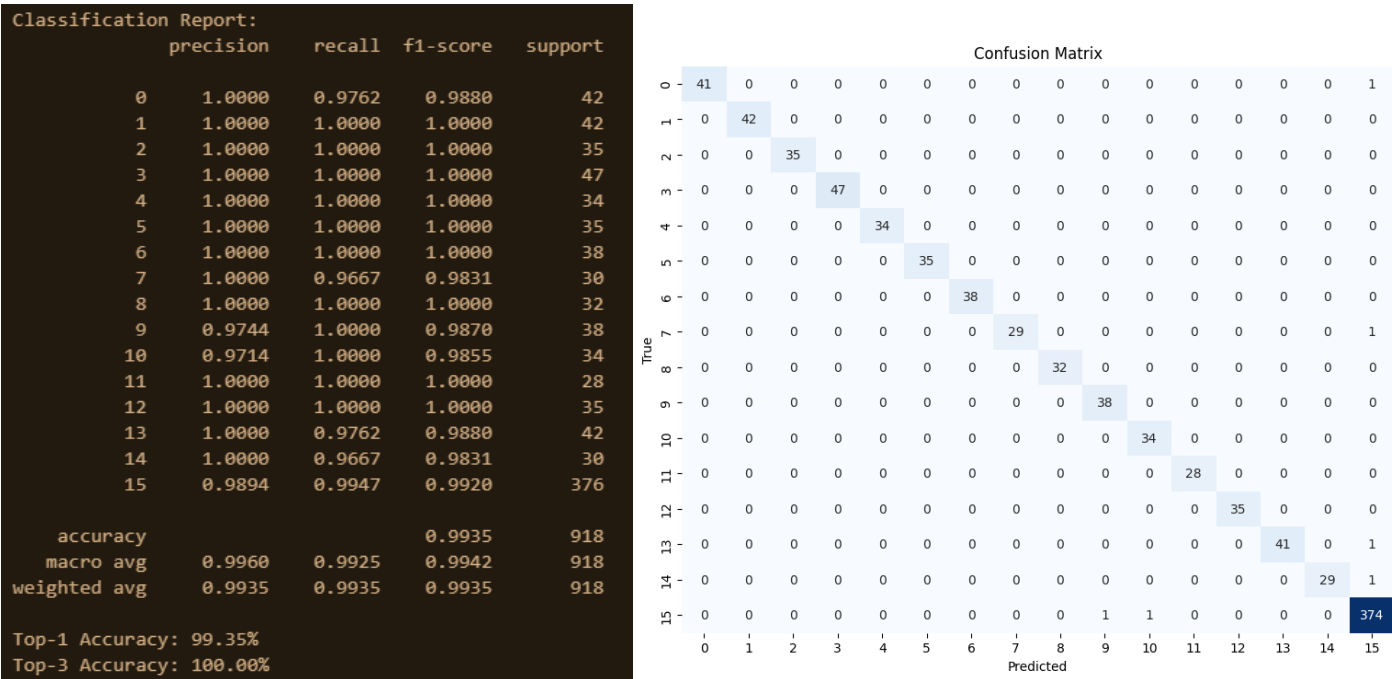
- Training ran for a total of 150 epochs across two phases.
- Model consistently achieved strong validation accuracy (>99%).
- This model was later used to classify cropped YOLO predictions during evaluation and visualization.



```
Top-1 Accuracy: 99.35%
Top-3 Accuracy: 100.00%
Macro F1 Score: 0.9935
```

→ Final results on validation set:

- Top-1 Acc: 99.35%
- Top-3 Acc: 100.00%
- Weighted F1: 0.9935



## Conclusion:

This project involved two integrated models: a YOLOv11n object detector for gender detection and an EfficientNet-based classifier for celebrity recognition.

1. The **YOLOv11n model**, trained on a custom gender-labeled dataset (~2,000 images), achieved solid detection performance with:
- Precision: 85.7%
  - Recall: 68.9%
  - mAP@0.5: 77.8%
  - mAP@0.5:0.95: 51.8%

These metrics indicate reliable detection, particularly for the male class, with lightweight architecture optimized for fast inference on small datasets.

2. The **EfficientNet\_V2\_S** classifier, fine-tuned in two stages using MixUp regularization and strong augmentations, reached >99% validation accuracy. It effectively learned to classify 16 celebrity identities from cropped facial images, even handling unknowns by remapping them to a dedicated class.

Together, these models form a robust pipeline for gender and identity recognition with efficient computation, high accuracy, and resilience to data variation—well-suited for real-world deployment or downstream analysis.

# Validation:

