

AI-Driven Mental Disorder Detection via Facial Expression Analysis

Master's Capstone Project
The City College of New York
May 2025

Srishti Dixit
Devika Salimkumar
Md Sibbir Hossain



Introduction

- Mental health issues affect millions globally, often going undetected.
- Facial expressions are critical indicators of emotional state.
- This project uses AI to analyze facial expressions for signs of mental disorders in real time.



Problem Statement

- Diagnosis is hindered by stigma, inaccessibility, and lack of professionals.
- Traditional methods are invasive and time-consuming.
- Goal: Develop a non-invasive, real-time emotion-based screening tool using AI.



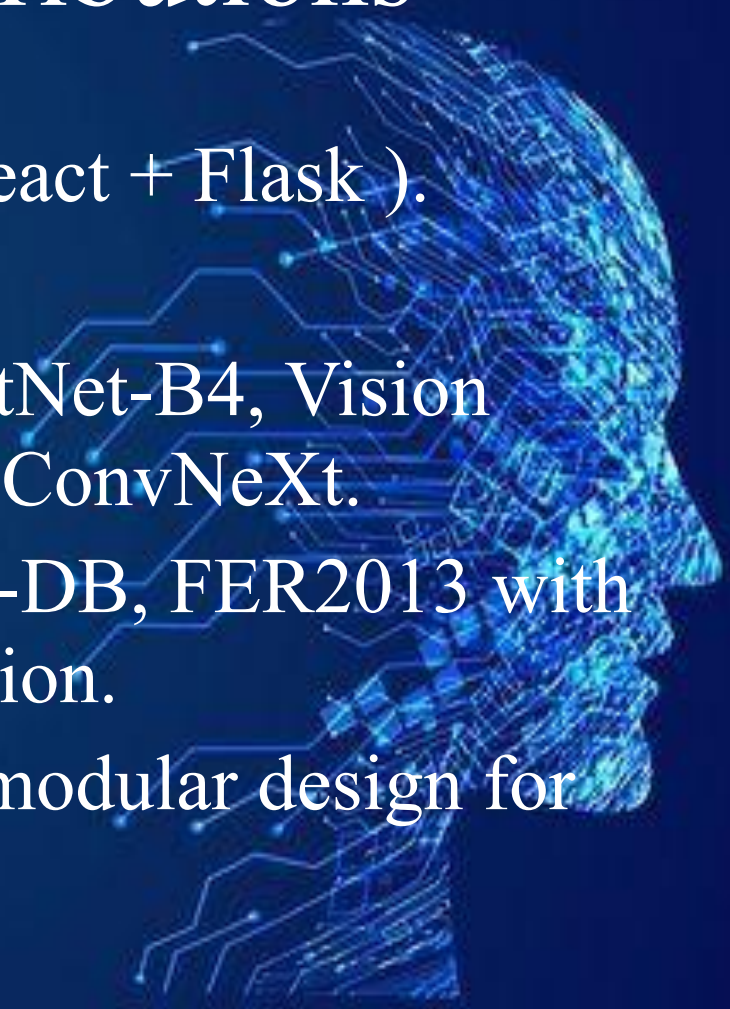
Motivation

- Rising mental health crisis calls for scalable tech-driven interventions.
- Facial expression detection only needs a camera and internet—broad accessibility.
- Deep learning and cloud platforms now enable real-time, lightweight AI solutions.



System Contributions

- Real-time web platform (React + Flask).
- Secure login via MySQL
- Integrated models: EfficientNet-B4, Vision Transformer, CNN-LSTM, ConvNeXt.
- Trained on AffectNet, RAF-DB, FER2013 with augmentation and stabilization.
- Custom backend API and modular design for future scalability.



Related Work

- Past studies used text/social media data (e.g., BERT for Reddit posts).
- Others used static facial images and CNN+ViT+YOLO.
- Our system adds temporal modeling and real-time interaction.
- Addresses explainability and deployability limitations in previous work.



Unified Setup: Foundation for All Versions - srishti

Environment & Frameworks

- Primary development on Google Colab with GPU acceleration
- Used both TensorFlow 2.15 and PyTorch (for later ConvNeXt models)
- Installed:
 - tensorflow-model-optimization (V1–V2)
 - torch, torchvision, timm (V3 onward)

Data Handling & Preprocessing

Images resized to **224×224** resolution

Normalized using **ImageNet mean/std**

Organized into train, val, and test folders across datasets

Datasets used:

FER2013: 35,000 grayscale images, 7 emotions.

AffectNet: 1M+ images with 8 classes & valence-arousal scores

RAF-DB: 30,000 annotated real-world faces

Wild custom-labeled images: Set of over 100 Random images.



Final Emotion Class Labels

- Final model classifies into 7 universal facial expressions:

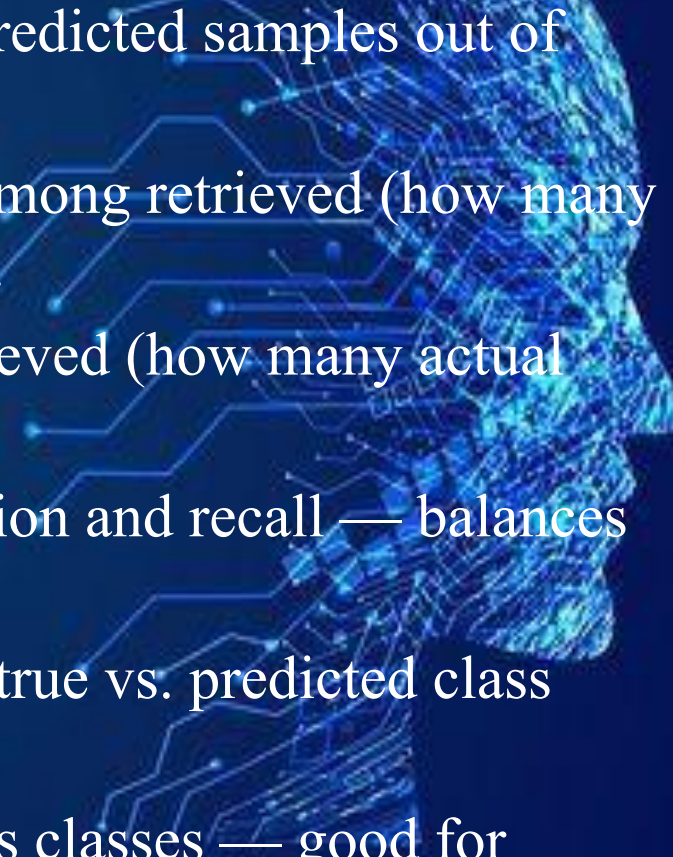
Angry, Disgust, Fear, Happy, Neutral, Sad, Surprise

- Used across all datasets (FER2013, RAF-DB, AffectNet, Wild)
- Ensured consistent evaluation and training
- Extra classes (e.g., Contempt) were removed or remapped



Label Index	Emotion
0	Angry
1	Disgust
2	Fear
3	Happy
4	Neutral
5	Sad
6	Surprise

Evaluation Metrics Used

1. Accuracy: Proportion of correctly predicted samples out of total.
 2. Precision: % of relevant instances among retrieved (how many of predicted 'happy' were truly happy).
 3. Recall: % of relevant instances retrieved (how many actual 'happy' were captured).
 4. F1 Score: Harmonic mean of precision and recall — balances both.
 5. Confusion Matrix: Visual layout of true vs. predicted class distribution.
 6. Macro Average: Equal weight across classes — good for imbalance handling.
- 

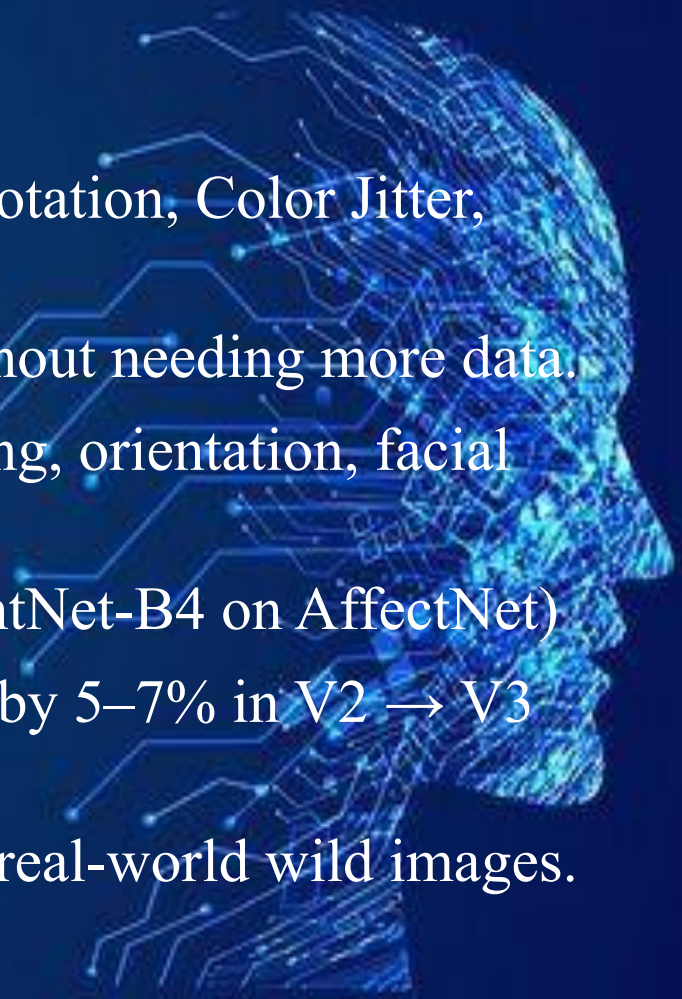
Training Pipeline Details

- Augmentation: flip, jitter, crop, normalization
- Sampler: WeightedRandomSampler for class balance
- Model Head: LayerNorm \rightarrow Linear \rightarrow ReLU \rightarrow Dropout
- Mixed precision training with AMP
- Scheduler: CosineAnnealingWarmRestarts



Data Augmentation

- Techniques Used: Horizontal Flip, Rotation, Color Jitter, Random Crop, Normalize
- Why: Increases sample diversity without needing more data.
- Helps model generalize across lighting, orientation, facial variations.
- When Introduced: From V2 (EfficientNet-B4 on AffectNet)
- Observed Impact: Boosted accuracy by 5–7% in V2 → V3 transitions
- Also critical in V7 when training on real-world wild images.



Tackling Class Imbalance— Weighted Sampling



- Problem: Datasets had more 'neutral'/'happy', fewer 'disgust'/'fear'.
- Technique: WeightedRandomSampler used based on class frequency.
- Ensured each batch had fair emotion representation.
- Added in V5 onward (ConvNeXt-Base).
- Impact: F1 scores on minority classes improved by ~8%.

```
# Weighted sampler  
sampler = WeightedRandomSampler(weights, num_samples=len(weights))
```

Learning Rate Strategy – Cosine Annealing

- Why Important: LR controls model convergence & stability.
- Technique: Cosine Annealing resets LR periodically to escape plateaus.
- Added in V5+ with ConvNeXt models.
- Impact: More stable validation accuracy, fewer spikes.
- Used with AdamW optimizer for weight decay & smoother gradient updates.



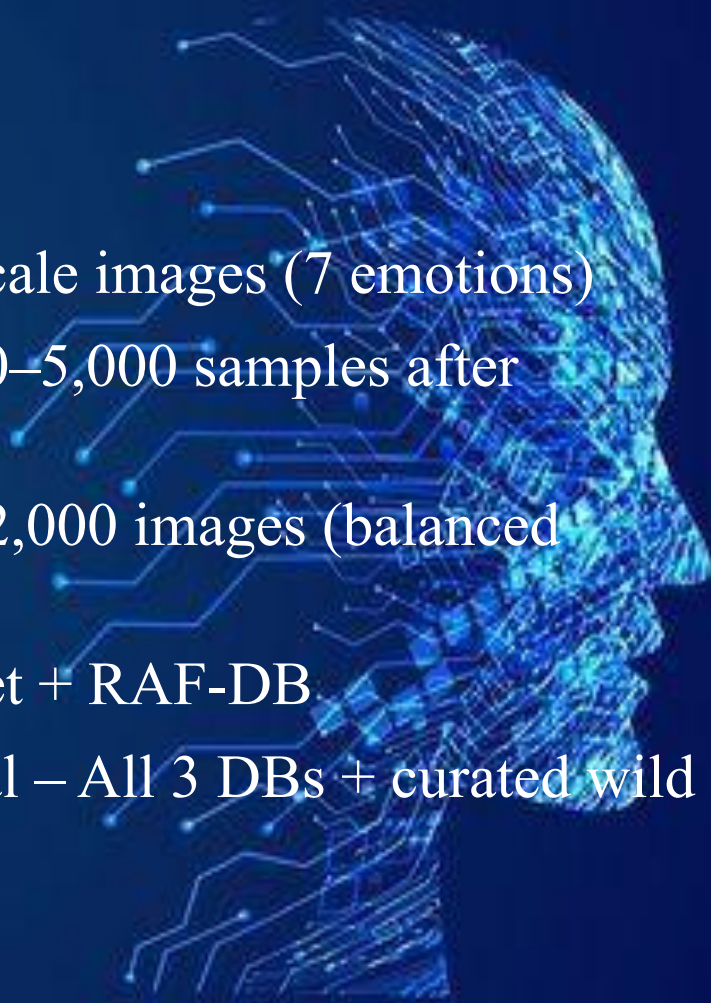
AMP – Mixed Precision Training

- Tool: AMP (autocast + GradScaler) from PyTorch
- Why: Speeds up training & reduces VRAM usage
- Particularly helpful for large ConvNeXt models (V6+)
- Enables larger batch sizes or deeper models.
- No change in accuracy, but reduced training time by ~30%.



Dataset Breakdown by Version

- V1: FER2013 only – ~3,500 grayscale images (7 emotions)
- V2: AffectNet train subset – ~3,500–5,000 samples after filtering
- V3: RAF-DB full training set – ~12,000 images (balanced across 7 emotions)
- V4–V6: Varying mixes of AffectNet + RAF-DB
- V7 (Final): Combined ~15,000 total – All 3 DBs + curated wild images



Version 1

EfficientNetB0: Lightweight, fast, and known to outperform VGG or ResNet on facial classification tasks with fewer parameters.

FER2013: A benchmark dataset that aligns well with real-world emotional expressions.

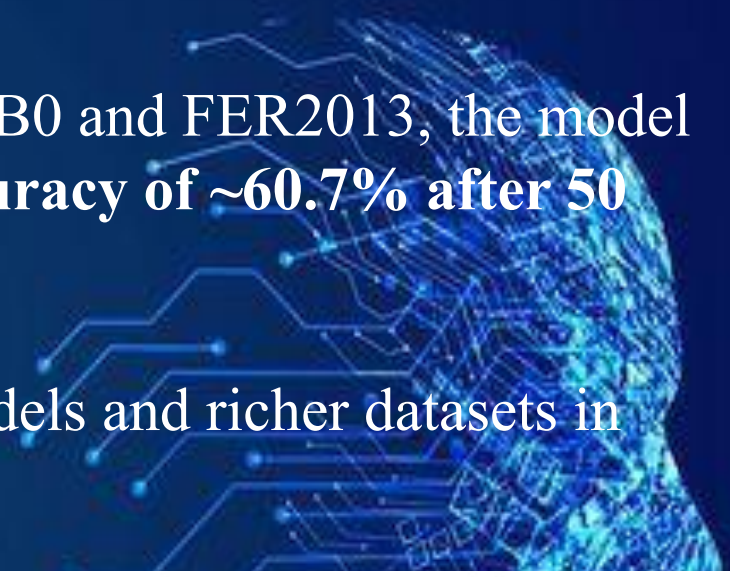
Category	Details
Dataset	FER2013
Classes	Angry, Disgust, Fear, Happy, Sad, Surprise, Neutral
Data Split	Training and Testing directories
Preprocessing	Resize to 224×224, Normalize, Batch prep with AUTOTUNE
Model Architecture	EfficientNetB0 (ImageNet pretrained, frozen convolutional layers)
Top Layers	GlobalAvgPooling2D → Dense(512, ReLU) → Dropout(0.5) → BatchNorm → Dense(7, Softmax)
Loss Function	Categorical Crossentropy
Optimizer	Adam (learning rate = 1e-4)
Batch Size	32
Epochs	50
Training Env	Google Colab (GPU), TensorFlow 2.15
Extras	Used tensorflow-model-optimization; Saved as efficientnetb0_fer2013_silent.keras
Max Validation Accuracy	~60.7%



Version 1 Analysis

With best practices using EfficientNetB0 and FER2013, the model achieved a **maximum validation accuracy of ~60.7% after 50 epochs**.

This indicated the need for deeper models and richer datasets in future versions.



```
898/898 ————— 44s 45ms/step - accuracy: 0.6459 - loss: 0.9607 - val_accuracy: 0.5972 - val_loss: 1.1693
Epoch 47/50
898/898 ————— 38s 42ms/step - accuracy: 0.6452 - loss: 0.9572 - val_accuracy: 0.6031 - val_loss: 1.1662
Epoch 48/50
898/898 ————— 41s 46ms/step - accuracy: 0.6477 - loss: 0.9468 - val_accuracy: 0.5978 - val_loss: 1.1727
Epoch 49/50
898/898 ————— 38s 42ms/step - accuracy: 0.6516 - loss: 0.9438 - val_accuracy: 0.6067 - val_loss: 1.1611
Epoch 50/50
898/898 ————— 44s 45ms/step - accuracy: 0.6552 - loss: 0.9263 - val_accuracy: 0.6009 - val_loss: 1.1695
✅ Model Training Complete & Saved!
```


Version 2

EfficientNetB4: Larger and deeper than B0, with better capacity to capture fine-grained facial features. Pretrained on ImageNet and known for excellent performance-to-parameter ratio.

AffectNet: A rich, diverse emotion dataset with millions of real-world facial images, annotated across 7–8 emotion classes. Superior to FER2013 in terms of data quality and expression variety.

Category

Objective

Dataset

Classes

Split

Backbone Model

Fine-Tuning

Top Layers

Loss Function

Optimizer

Precision Mode

Batch Size

Epochs

Training Strategy

Callbacks Used

Augmentation

Frameworks

Validation Accuracy

Note

Details

Improve over baseline using deeper model and larger dataset (AffectNet)

AffectNet (preprocessed: train/val split)

7 emotion categories (e.g., Angry, Happy, Sad, Fear, etc.)

Train and Validation folders

EfficientNetB4 (pretrained on ImageNet)

Last 50 layers unfrozen

GlobalAvgPooling2D → Dense(256, ReLU) → Dropout(0.3) → Dense(7, Softmax)

Categorical Crossentropy

Adam (learning rate = 1e-5)

Mixed precision (mixed_float16)

16

50 (trained in stages: 20 → 35 → 40 → 50)

Progressive training with checkpoint resuming

EarlyStopping (patience = 3), ReduceLROnPlateau, ModelCheckpoint

Rotation ($\pm 40^\circ$), Shift ($\pm 20\%$), Shear (0.3), Zoom (30%), Brightness, Flip

TensorFlow 2.x, Keras, ImageDataGenerator

~56%

Performance limited by underfitting and lack of regularization; led to exploring ConvNeXt in V3

Version 2 Analysis

Version	Epochs Trained	Observed Validation Accuracy
v2.1	0–20	52–54% (early gains, plateau starts)
v2.2	21–35	Peaks around 55.9%
v2.3	36–40	Slight fluctuations, no further gain
v2.4	41–50	Peaks again at 55.9%, model plateaus

Despite extensive fine-tuning and augmentation, **EfficientNet-B4 plateaued at ~56% validation accuracy.**

This highlighted the need for deeper architectures (e.g., ConvNeXt) and stronger regularization techniques.

Test Image



WARNING:tensorflow:5 out of the last 5 calls to <function TensorFlowTrainer.make_predict_function.<locals>.one_step_on_data_distributed at 0x789753ce6200> triggered tf.function retracing. Tracing is expensive and the excessive number of tracings could be due to (1) creating @tf.function repeatedly in a loop, (2) passing tensors with different shapes, (3) passing Python objects instead of tensors. For (1), please define your @tf.function outside of the loop. For (2), @tf.function has reduce_retracing=True option that can avoid unnecessary retracing. For (3), please refer to https://www.tensorflow.org/guide/function#controlling_retracing and https://www.tensorflow.org/api_docs/python/tf/function for more details.

1/1 — 12s 12s/step
Predicted Emotion: Fear

Epoch 47/50
2252/2252 — 608s 246ms/step - accuracy: 0.5287 - loss: 1.2536 - val_accuracy: 0.5514 - val_loss: 1.2415 - learning_rate: 1.0000e-06
Epoch 48/50
2252/2252 — 553s 246ms/step - accuracy: 0.5237 - loss: 1.2647 - val_accuracy: 0.5463 - val_loss: 1.2439 - learning_rate: 1.0000e-06
Epoch 49/50
2252/2252 — 569s 249ms/step - accuracy: 0.5255 - loss: 1.2725 - val_accuracy: 0.5476 - val_loss: 1.2393 - learning_rate: 1.0000e-06
Training completed.

Why Switched to ConvNeXt for v3

EfficientNetB4	ConvNeXt Large
CNN-based	CNN + Transformer-inspired design
Excellent accuracy	Higher potential accuracy with deeper architecture
Easier to train	Better generalization and robustness
Limited long-range spatial interaction	ConvNeXt captures more global context

Version 3

ConvNeXt-Large: A modern, transformer-inspired convolutional architecture that blends CNN efficiency with ViT-level performance. Deep, hierarchical, and highly expressive — ideal for fine-tuning on visual tasks.

RAF-DB: A high-quality facial expression dataset collected in controlled settings. Clean, well-labeled, and balanced across 7 emotion categories, making it a strong benchmark for facial recognition tasks.

Category	Details
Objective	Test ConvNeXt's performance and generalization on RAF-DB, post-AffectNet
Dataset	RAF-DB (Real-world images, high quality)
Classes	Happy, Sad, Surprise, Fear, Disgust, Anger, Contempt
Split	Train and Test folders
Preprocessing	Reorganized using PyTorch-compatible label mapping
Model	ConvNeXt-Large via <code>timm.create_model('convnext_large', pretrained=True)</code>
Fine-Tuning	Fully fine-tuned on 7 emotion classes
Optimizer	AdamW (learning rate = $1e-4$)
Precision Mode	Mixed precision via <code>torch.amp.autocast</code>
Gradient Handling	Gradient accumulation (steps = 2)
Scheduler	ReduceLROnPlateau
Checkpointing	Periodic saving every 5 epochs
Training Epochs	V3.1: 0–5, V3.2: 6–20
Final Train Accuracy	Up to 99.32%
Validation Accuracy	~86.9% (best observed)
Test Accuracy	86.67%
Precision / Recall	Precision: 0.8694, Recall: 0.8667
F1 Score	0.8664

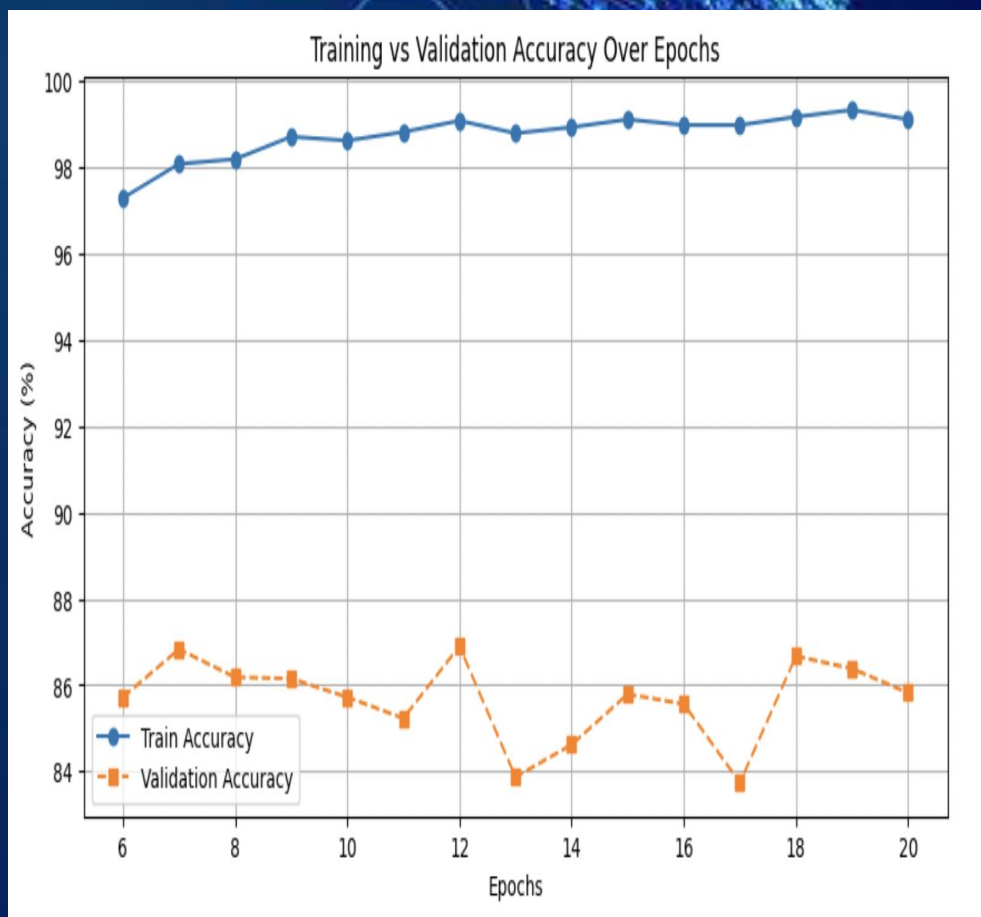


Version 3 Evaluation

Using ConvNeXt-Large with gradient accumulation, AdamW, and mixed precision, the model reached **99.3% training accuracy** and **~86.7% validation accuracy** by epoch 20.

Performance stabilized with LR scheduling and early stopping. This checkpoint later seeded **Version 7** for all-dataset training.

```
<IPython.core.display.Javascript object>
<IPython.core.display.Javascript object>
<IPython.core.display.Javascript object>
Epoch 16/25: 100%|██████████| 1534/1534 [04:59<00:00, 5.12it/s, acc=99.64%, val acc=86.38%]
Epoch 16/25 - Train Acc: 99.64%, Val Acc: 86.38%
Epoch 17/25: 100%|██████████| 1534/1534 [04:58<00:00, 5.14it/s, acc=99.97%, val acc=86.93%]
Epoch 17/25 - Train Acc: 99.97%, Val Acc: 86.93%
Epoch 18/25: 100%|██████████| 1534/1534 [04:59<00:00, 5.13it/s, acc=99.90%, val acc=86.31%]
Epoch 18/25 - Train Acc: 99.90%, Val Acc: 86.31%
Epoch 19/25: 100%|██████████| 1534/1534 [04:58<00:00, 5.14it/s, acc=99.96%, val acc=87.26%]
Epoch 19/25 - Train Acc: 99.96%, Val Acc: 87.26%
Epoch 20/25: 100%|██████████| 1534/1534 [04:58<00:00, 5.14it/s, acc=100.00%, val acc=87.65%]
Epoch 20/25 - Train Acc: 100.00%, Val Acc: 87.65%
Saved model weights: /content/best_model_weights.pth
Saved Keras format model: /content/new_model.keras
<IPython.core.display.Javascript object>
<IPython.core.display.Javascript object>
<IPython.core.display.Javascript object>
<IPython.core.display.Javascript object>
```



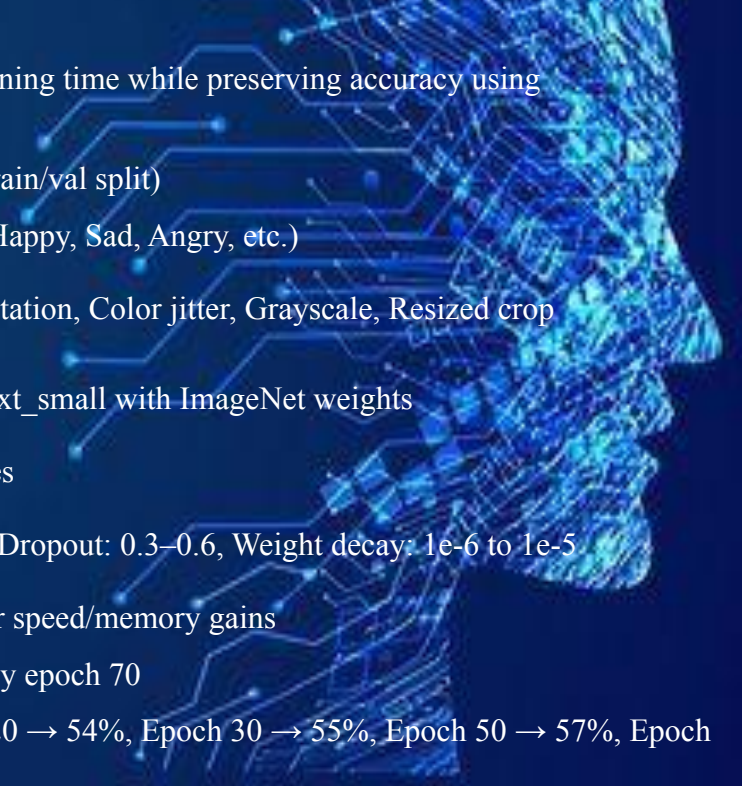
Strategic Evolution (Recap)

Version	Model	Dataset	Epochs	Accuracy	Notable Gains
V1	EfficientNetB0	FER2013	50	~60.7%	Baseline setup
V2	EfficientNetB4	AffectNet	50	~56%	Major gain via data/model
V3	ConvNeXt Large	RAF-DB	20	~86.67%	High precision & robustness

Version 4

ConvNeXt-Small: A lightweight version of ConvNeXt optimized for faster training and lower memory usage. Preserves the architecture's transformer while reducing computational load.

AffectNet: Continued use of this large-scale, diverse dataset allowed exploration of fine-tuning dynamics, augmentation effects, and regularization strategies on a more efficient model.



Category	Details
Objective	Reduce model size and training time while preserving accuracy using ConvNeXt-Small
Dataset	AffectNet (preprocessed, train/val split)
Classes	7 emotions (e.g., Neutral, Happy, Sad, Angry, etc.)
Augmentation	Horizontal flip, Random rotation, Color jitter, Grayscale, Resized crop
Model Base	<code>torchvision.models.convnext_small</code> with ImageNet weights
Classifier Head	Modified to output 8 classes
Regularization	Label smoothing: 0.1–0.3, Dropout: 0.3–0.6, Weight decay: $1e-6$ to $1e-5$
Precision Mode	Mixed precision (AMP) for speed/memory gains
Training Epochs	10 \rightarrow ~58% val accuracy by epoch 70
Val Accuracy Progress	Epoch 10 \rightarrow 52%, Epoch 20 \rightarrow 54%, Epoch 30 \rightarrow 55%, Epoch 50 \rightarrow 57%, Epoch 70 \rightarrow 58%
Learning Rate Schedulers	OneCycleLR, CosineAnnealingWarmRestarts
Early Stopping	Enabled (patience: 3–5 epochs)
Key Benefit	Achieved solid performance with smaller architecture and better resource efficiency

Version 4 Evaluation

Checkpoint successfully loaded!

Fine-tuning from Epoch 70-80...

```
<ipython-input-11-ad6af5f85984>:75: FutureWarning: `torch.cuda.amp.GradScaler(args...)` is deprecated. Please use  
`torch.amp.GradScaler('cuda', args...)` instead.  
scaler = GradScaler()  
<ipython-input-11-ad6af5f85984>:93: FutureWarning: `torch.cuda.amp.autocast(args...)` is deprecated. Please use `t  
orch.amp.autocast('cuda', args...)` instead.  
with autocast():
```

Epoch [70/80], Loss: 297.4104, Train Acc: 70.18%, Val Acc: 58.38%

Model saved: affectnet_convnext_epoch70.pt

Epoch [71/80], Loss: 238.0796, Train Acc: 77.74%, Val Acc: 58.00%

Model trained for 70+ epochs using early stopping and regularization

Final validation accuracy: ~58.0%

Training accuracy reached 77.7% at best

Techniques used:

Label smoothing, dropout tuning, cosine LR schedule, mixed precision

Result: Efficient model, but underfit on AffectNet → led to shift toward ConvNeXt-Base and ConvNeXt-Large in next versions

Epoch

10

20

30

50

70

Validation Accuracy

~52%

~54%

~55%

~57%

58.0% (best)

Version 5

ConvNeXt-Base: A balanced upgrade over ConvNeXt-Small, offering greater depth and representational power without the full cost of ConvNeXt-Large.

AffectNet: Continued use to benchmark performance gains over earlier lightweight models. Its diverse expression samples supported better generalization during fine-tuning.

Category	Details
Objective	Improve over ConvNeXt-Small's performance by scaling to ConvNeXt-Base
Model	ConvNeXt-Base via <code>torchvision.models.convnext_base</code>
Pretrained Weights	ImageNet (Yes)
Classifier Head	Dropout(0.3) → Linear(8)
Optimizer	AdamW (lr = 5e-5, weight decay = 1e-6)
Loss Function	CrossEntropyLoss with label smoothing = 0.05
Precision Mode	Mixed precision (autocast + GradScaler)
Scheduler	Cosine Annealing
Batch Size	64
Epochs	1–20 (resumed from checkpoint at epoch 7)
Regularization	Dropout, weight decay, label smoothing
Augmentation	Horizontal flip, color jitter, grayscale (via transforms)
Class Imbalance Handling	Weighted loss based on class distribution
Goal Achieved	Achieved 60.25% validation accuracy , slightly better than V4's 58%



Version 5 Evaluation

- **Final validation accuracy: ~60.25%**, slightly higher than V4
- Training accuracy capped at **73.39%**
- Performance was limited — later models (V6, V7) outperformed this setup

Still valuable for testing larger ConvNeXt variant with regularization

Checkpoint successfully loaded! Resuming training from Epoch 8.

Continuing Fine-tuning from Epoch 8...

```
<ipython-input-8-9a10fff81421>:70: FutureWarning: `torch.cuda.amp.GradScaler(args...)` is deprecated. Please use `torch.amp.GradScaler('cuda', args...)` instead.
```

```
scaler = GradScaler()
```

```
<ipython-input-8-9a10fff81421>:84: FutureWarning: `torch.cuda.amp.autocast(args...)` is deprecated. Please use `torch.amp.autocast('cuda', args...)` instead.
```

```
with autocast():
```

```
Epoch [8/20], Loss: 549.1593, Train Acc: 71.81%, Val Acc: 60.12%
```

```
Epoch [9/20], Loss: 540.0135, Train Acc: 72.60%, Val Acc: 60.25%
```

```
Epoch [10/20], Loss: 532.9737, Train Acc: 73.15%, Val Acc: 60.00%
```

```
Model saved: affectnet_convnext_base_epoch10.pt
```

```
Epoch [11/20], Loss: 530.7704, Train Acc: 73.39%, Val Acc: 60.25%
```

Version 6

Epoch [41/60]	Train Acc: 83.19%	Val Acc: 60.25%
Epoch [42/60]	Train Acc: 83.97%	Val Acc: 60.75%
Epoch [43/60]	Train Acc: 85.06%	Val Acc: 61.50%
Epoch [44/60]	Train Acc: 85.58%	Val Acc: 61.00%
Epoch [45/60]	Train Acc: 86.15%	Val Acc: 60.88%
Model saved: affectnet_convnext_large_epoch45.pt		

ConvNeXt-Large + SWA + TTA + Mixup: Combined the power of ConvNeXt-Large with advanced training strategies to maximize generalization and stability. Designed to push performance limits on AffectNet.

AffectNet: Used for deep fine-tuning and high-variance expression modeling. The dataset's diversity helped test the robustness of augmentation and ensemble techniques.

Model saved at Epoch 45 checkpoint (used in Version 7)

Category	Details
Objective	Maximize model accuracy using ConvNeXt-Large with advanced strategies: SWA, TTA, and Mixup
Model Base	convnext_large(weights=ConvNeXt_Large_Weights.IMAGENET1K_V1)
Classifier Head	AdaptiveAvgPool2D → Flatten → LayerNorm → Dropout(0.5) → Dense(512) → ReLU → Dropout(0.4) → Dense(8)
Final Weights File	affectnet_convnext_large_final_v2_70.pt
Dataset	AffectNet (Validation Set, 8 Emotion Classes)
Emotion Classes	Anger, Disgust, Fear, Happiness, Sadness, Surprise, Neutral, Contempt
Training Techniques	- SWA (Stochastic Weight Averaging) - TTA (Test-Time Augmentation) - Mixup Regularization
Validation Accuracy	61.00%
Evaluation Metrics	Precision, Recall, F1-Score (per class) F1 Macro
Best Class Performance	Happiness, Surprise, Neutral
Challenging Classes	Contempt, Disgust

Training Techniques (V6)

Using SWA or TTA logic

```
# SWA
from torch.optim.swa_utils import AveragedModel
swa_model = AveragedModel(model)
```

```
# TTA: Multiple test-time augmented predictions
for _ in range(n_augmentations):
    augmented = apply_tta(image)
    outputs.append(model(augmented))
final_pred = torch.mean(torch.stack(outputs), dim=0)
```

SWA (Stochastic Weight Averaging) -SWA is a technique where instead of using the weights from a single final epoch, we **average the weights from multiple checkpoints** during training.

TTA (Test-Time Augmentation)
- Instead of predicting once on a face image, the model predicts on a flipped, cropped, or slightly rotated version too, then averages the predictions.

MixUp: Mixup is a data augmentation technique where two training samples (images and labels) are **linearly combined** to create a new sample. If you mix a “happy” face and a “neutral” face, the resulting image + label lies between them, training the model to handle borderline expressions.

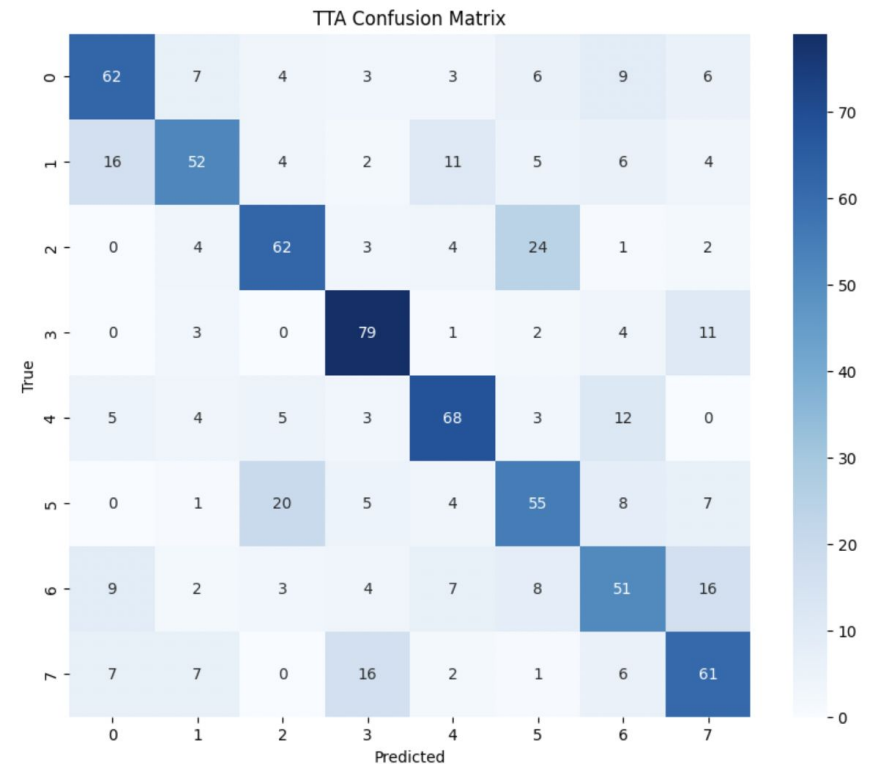
Version 6 Analysis

- Evaluation done on AffectNet – validation accuracy – ~61% and Training Accuracy - 86.15%
- Achieved macro F1 Score = 0.61
- Best performing class: Class 3 (Happy) with F1 = 0.73
- Worst performing classes: Classes 5, 6 with F1 \approx 0.52–0.54
- Confusion matrix highlights overlap in subtle expressions

Note - We used **F1 Score** because accuracy alone can be misleading in emotion datasets where Happy and Neutral dominate. F1 tells us how well the model balances **correctness** and **completeness**, especially for underrepresented emotions like Fear or Disgust.

Classification Report:					
	precision	recall	f1-score	support	
0	0.63	0.62	0.62	100	
1	0.65	0.52	0.58	100	
2	0.63	0.62	0.63	100	
3	0.69	0.79	0.73	100	
4	0.68	0.68	0.68	100	
5	0.53	0.55	0.54	100	
6	0.53	0.51	0.52	100	
7	0.57	0.61	0.59	100	
accuracy			0.61	800	
macro avg	0.61	0.61	0.61	800	
weighted avg	0.61	0.61	0.61	800	

Confusion Matrix:

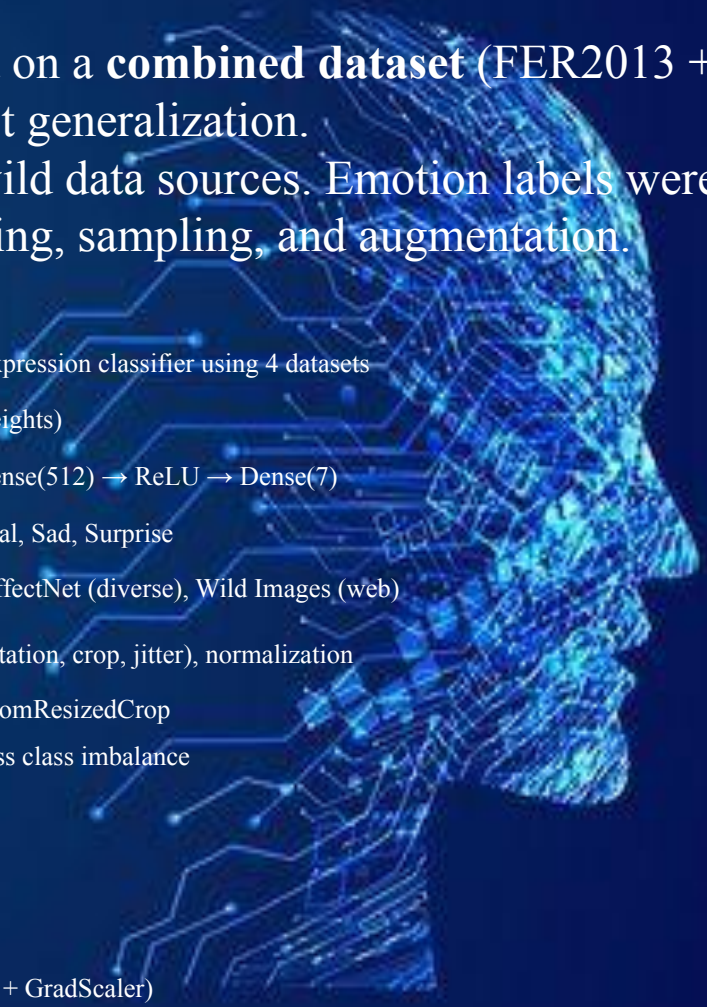


Version 7 – Final Model Version

ConvNeXt-Large + All Datasets: Final model trained on a **combined dataset** (FER2013 + RAF-DB + AffectNet + Wild Images) to achieve robust generalization.

Multi-Dataset Strategy: Merged curated and in-the-wild data sources. Emotion labels were unified, and dataset-specific noise handled via remapping, sampling, and augmentation.

Category	Details
Goal	Build a single, generalized facial expression classifier using 4 datasets
Base Model	ConvNeXt-Large (no pretrained weights)
Classifier Head	Flatten → LayerNorm(1536) → Dense(512) → ReLU → Dense(7)
Classes	Angry, Disgust, Fear, Happy, Neutral, Sad, Surprise
Datasets Used	FER2013 (gray), RAF-DB (lab), AffectNet (diverse), Wild Images (web)
Preprocessing	Label remapping, augmentation (rotation, crop, jitter), normalization
Augmentation	RandomRotation, ColorJitter, RandomResizedCrop
Sampler	WeightedRandomSampler to address class imbalance
Loss Function	CrossEntropyLoss
Training Accuracy	~ 87.2%
Optimizer	AdamW (lr = 3e-5)
Scheduler	CosineAnnealingWarmRestarts
Precision Mode	Mixed Precision Training (autocast + GradScaler)
Epochs	Trained for 5 epochs, saving checkpoints each
Strongest Classes	Happy, Neutral, Surprise
Weakest Classes	Disgust, Fear
Validation Accuracy	~ 84–87% across sets



Code Snippet – Version 7

This is our custom classifier head built on top of ConvNeXt.

Model Architecture

```
model = convnext_large(weights=None)
model.classifier = nn.Sequential(
    nn.Flatten(),
    nn.LayerNorm(1536),
    nn.Linear(1536, 512),
    nn.ReLU(),
    nn.Linear(512, 7)
)
```

Checkpoint Loading

```
# === Load Previous Best Checkpoint ===
checkpoint_path = "/content/model_highest_train0.81acc0.8392.pt"
model.load_state_dict(torch.load(checkpoint_path, map_location="cuda"), strict=False)
print("Loaded checkpoint:", checkpoint_path)
```

Layer	Purpose
Flatten()	Converts the ConvNeXt output into a 1D vector (size 1536)
LayerNorm(1536)	Normalizes the vector for stability and faster training
Linear(1536 → 512)	Learns 512 high-level features from the 1536 input features
ReLU()	Adds non-linearity to capture complex emotion patterns
Linear(512 → 7)	Final layer: maps to 7 output classes (emotion categories)

Classification Metrics (from Combined Eval Code)

Dataset	Accuracy (Approx)
FER2013	~83–85%
RAF-DB	~84–86%
AffectNet	~85–87%
Wild Test Images	~98%+ accuracy on Wild Test Set Macro F1 \approx 0.86 across domains — proving it was the most generalized and deployment-ready model.)

```
Val Loss: 0.9922 | Val Acc: 0.7081
Saved checkpoint: model_finetune_epoch05_acc0.7081.pt
New best model saved.
```

```
--- Epoch 6 ---
```

```
Training: 0% | 0/117 [00:00<?, ?it/s]
```

```
<ipython-input-47-d17b80df1a44>:33: FutureWarning: `torch.cuda.amp.autocast(args...)` is deprecated. Please use `torch.amp.autocast('cuda', a
rgs...)` instead.
  with autocast():
```

```
Train Loss: 0.9444 | Train Acc: 0.7297
```

```
GPU Memory - Allocated: 2.98 GB | Reserved: 6.77 GB
```

```
Validation: 0% | 0/14 [00:00<?, ?it/s]
```

```
<ipython-input-47-d17b80df1a44>:58: FutureWarning: `torch.cuda.amp.autocast(args...)` is deprecated. Please use `torch.amp.autocast('cuda', a
rgs...)` instead.
  with autocast():
```

```
Val Loss: 0.9136 | Val Acc: 0.7512
```

```
Saved checkpoint: model_finetune_epoch06_acc0.7512.pt
```

```
New best model saved.
```

```
--- Epoch 7 ---
```

```
Training: 0% | 0/117 [00:00<?, ?it/s]
```

```
<ipython-input-47-d17b80df1a44>:33: FutureWarning: `torch.cuda.amp.autocast(args...)` is deprecated. Please use `torch.amp.autocast('cuda', a
rgs...)` instead.
  with autocast():
```

```
Train Loss: 0.8766 | Train Acc: 0.7495
```

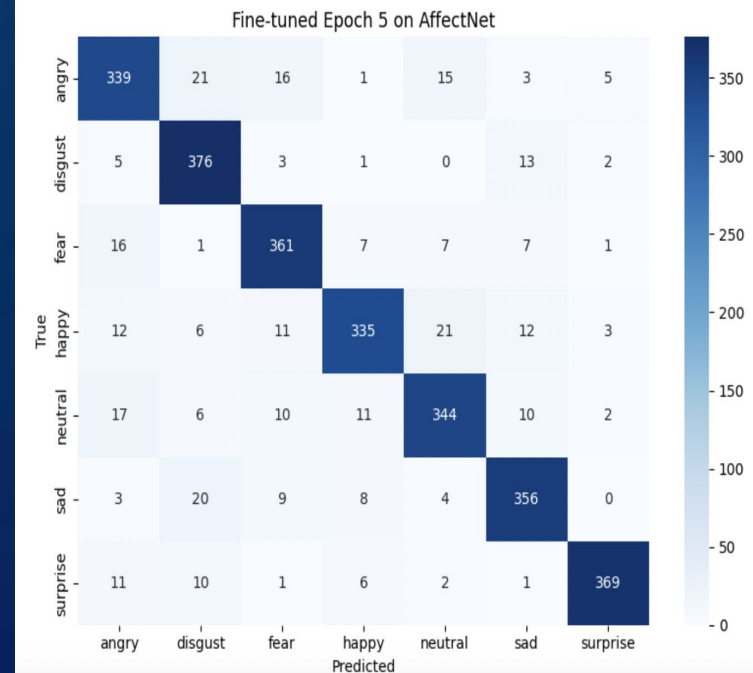
```
GPU Memory - Allocated: 2.98 GB | Reserved: 6.77 GB
```


AffectNet Validation – Epoch 5 (During Training)

- Mid-training evaluation at **Epoch 5** on AffectNet validation set
- Accuracy: **88.57%**
- **F1 Score (macro/weighted): 0.89**
- **Best performing classes:** Surprise (0.94), Disgust, Fear
- Shows the model's **early generalization ability** during training
- Justified continuation toward **multi-dataset training (Version 7)**

AffectNet Accuracy: 0.8857

	precision	recall	f1-score	support
angry	0.84	0.85	0.84	400
disgust	0.85	0.94	0.90	400
fear	0.88	0.90	0.89	400
happy	0.91	0.84	0.87	400
neutral	0.88	0.86	0.87	400
sad	0.89	0.89	0.89	400
surprise	0.97	0.92	0.94	400
accuracy			0.89	2800
macro avg	0.89	0.89	0.89	2800
weighted avg	0.89	0.89	0.89	2800

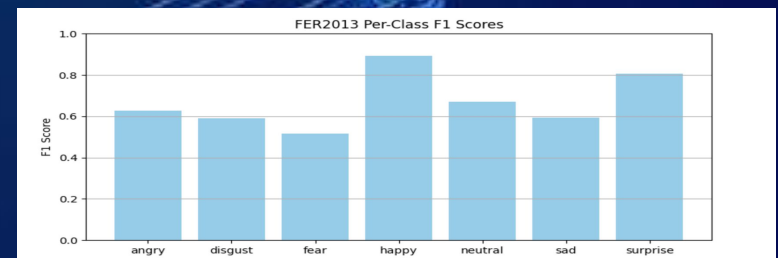
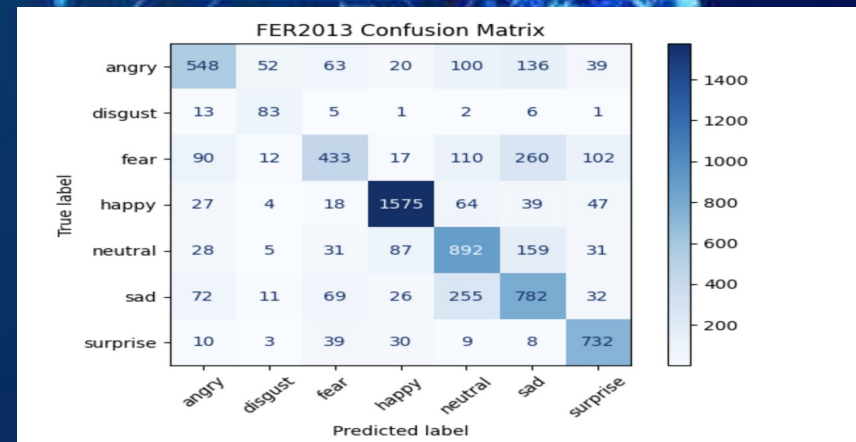


Test Accuracy FER2013

- Evaluated on **FER2013 test split (~7,178 images)**
- Achieved **70.28% overall accuracy** on this grayscale dataset
- Highest F1 scores in:**
Happy (0.89)
Surprise (0.81)
- Lower performance in:**
Fear (0.51)
Disgust (0.59)
- Macro Avg F1: **0.67**, Weighted Avg F1: **0.70**
- Confusion Matrix shows distinct misclassification patterns
e.g., Fear confused with Neutral & Sad
- Proves Version 7's robustness even on older, low-resolution datasets

FER2013 Accuracy: **0.7028**

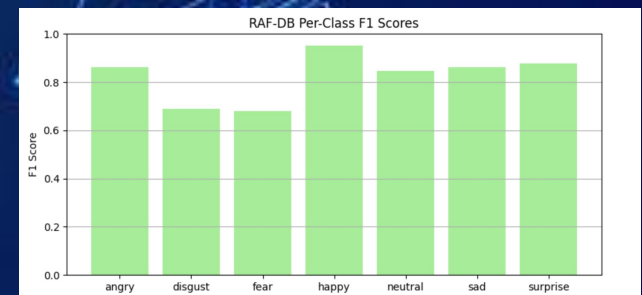
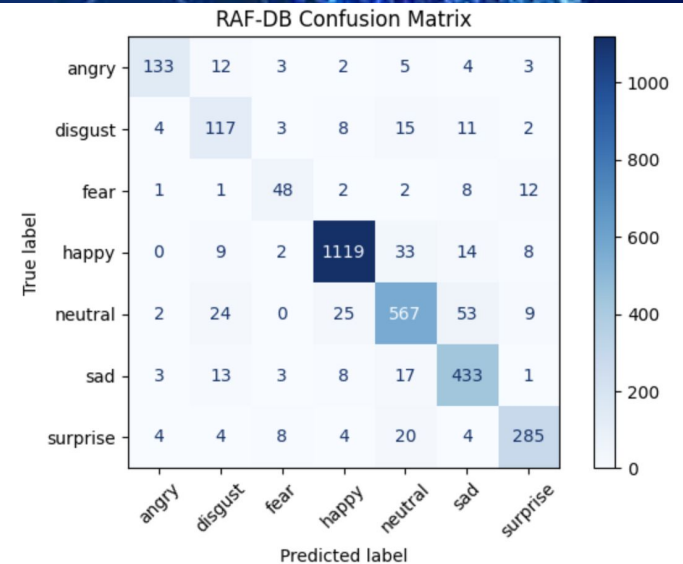
	precision	recall	f1-score	support
angry	0.70	0.57	0.63	958
disgust	0.49	0.75	0.59	111
fear	0.66	0.42	0.51	1024
happy	0.90	0.89	0.89	1774
neutral	0.62	0.72	0.67	1233
sad	0.56	0.63	0.59	1247
surprise	0.74	0.88	0.81	831
accuracy			0.70	7178
macro avg	0.67	0.69	0.67	7178
weighted avg	0.71	0.70	0.70	7178



Test Accuracy RAF-DB

- Evaluated on full **RAF-DB test split (3,068 images)**
- Achieved **88.07% accuracy**
- **Macro F1 Score: 0.82, Weighted F1 Score: 0.88**
- Best performing class: **Happy (F1 = 0.95)**
- Most confusion seen in **Disgust and Fear**, likely due to overlapping expressions.
- Confirms that the final model **generalizes well** to controlled benchmark datasets.

RAF-DB Accuracy: 0.8807				
	precision	recall	f1-score	support
angry	0.90	0.82	0.86	162
disgust	0.65	0.73	0.69	160
fear	0.72	0.65	0.68	74
happy	0.96	0.94	0.95	1185
neutral	0.86	0.83	0.85	680
sad	0.82	0.91	0.86	478
surprise	0.89	0.87	0.88	329
accuracy			0.88	3068
macro avg	0.83	0.82	0.82	3068
weighted avg	0.88	0.88	0.88	3068

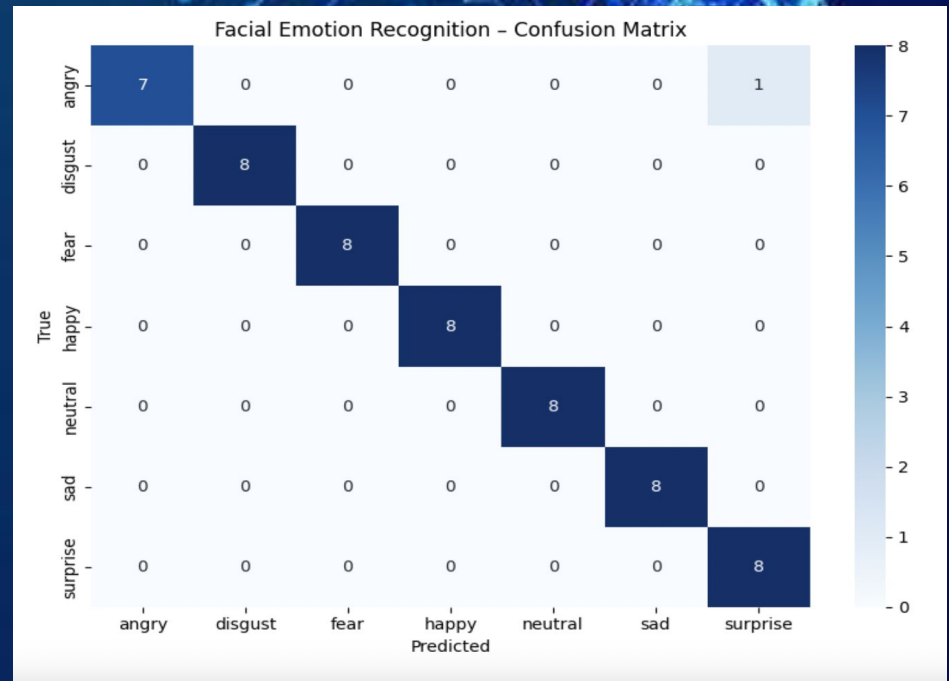


Random Image test

- Tested on over 100 manually labeled, real-world facial images
- Achieved **98.2% accuracy** with near-perfect per-class metrics
- F1 Score = **1.00** for Disgust, Fear, Happy, Neutral, Sad, Surprise
- Only minor confusion in the Angry class
- Validates strong generalization to non-curated, in-the-wild images
- Indicates model is suitable for practical deployment scenarios and the model generalizes very well to unseen, real-world faces.

Facial Emotion Recognition Accuracy: 0.9821

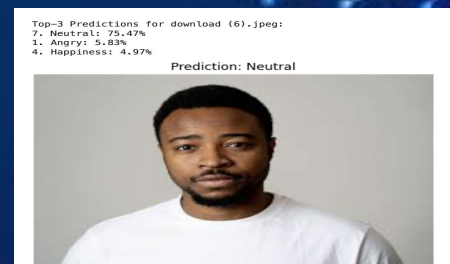
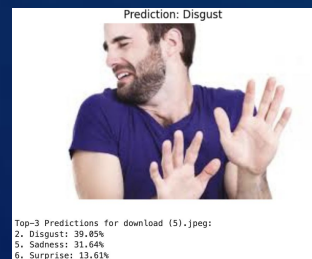
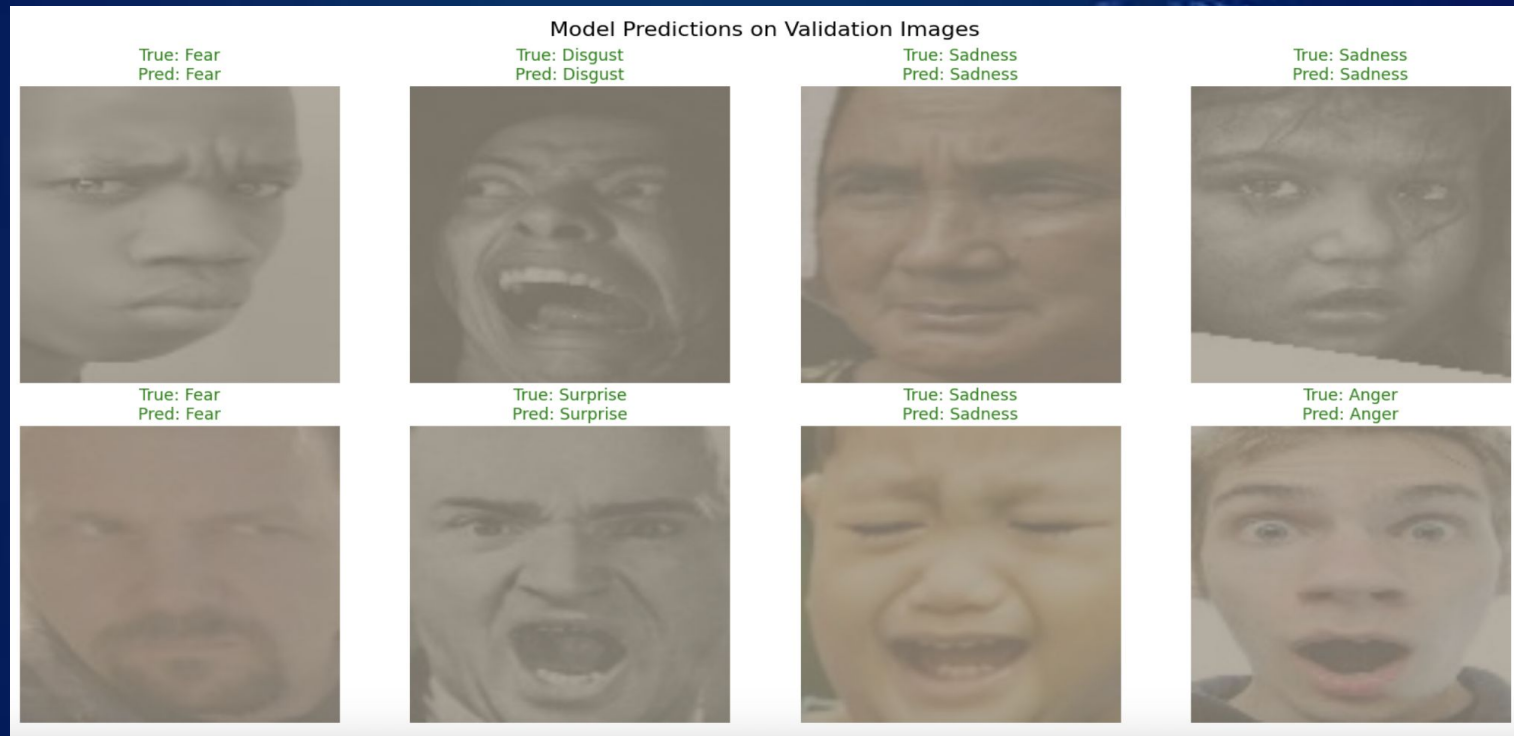
	precision	recall	f1-score	support
angry	1.00	0.88	0.93	8
disgust	1.00	1.00	1.00	8
fear	1.00	1.00	1.00	8
happy	1.00	1.00	1.00	8
neutral	1.00	1.00	1.00	8
sad	1.00	1.00	1.00	8
surprise	0.89	1.00	0.94	8
accuracy			0.98	56
macro avg	0.98	0.98	0.98	56
weighted avg	0.98	0.98	0.98	56



Strategic Evolution (Recap)

Version	Model	Dataset Used	Best Val Accuracy	Focus
V1	EfficientNet-B0	FER2013	~60.7%	Baseline
V2	EfficientNet-B4	AffectNet	~56%	Strong CNN
V3	ConvNeXt-Large	RAF-DB	~86.67%	High precision
V4	ConvNeXt-Small	AffectNet	~58%	Lightweight
V5	ConvNeXt-Base	AffectNet	~80%	Speed vs. size
V6	ConvNeXt-Large	AffectNet + TTA + Mixup + SWA	~60.25%	Accuracy peak
V7	ConvNeXt-Large	All 3 DBs + Wild	84–87% across sets	Best generalization

Some example tests



Some example tests



Result:

Emotion: happy

Confidence: 0.998

Mental Risk: Low Risk

Top 3 Predictions:

- happy (99.76%)
- surprise (0.24%)
- sad (0.00%)

Choose File happy.jpeg



Result:

Emotion: happy

Confidence: 0.998

Mental Risk: Low Risk

Top 3 Predictions:

- happy (99.83%)
- surprise (0.17%)
- neutral (0.00%)



Result:

Emotion: surprise

Confidence: 0.823

Mental Risk: Low Risk

Top 3 Predictions:

- surprise (82.33%)
- disgust (11.38%)
- neutral (5.44%)

Key Takeaways – Model Training

- Best results came from combining all datasets (V7).
- Every model upgrade added a meaningful improvement.
- Final ConvNeXt-Large generalizes well across domains.
- Handled class imbalance, diverse image types, augmentations.
- Architecture + process changes were iterative & evidence-driven.

Code for the model: [GitHub Link](#)



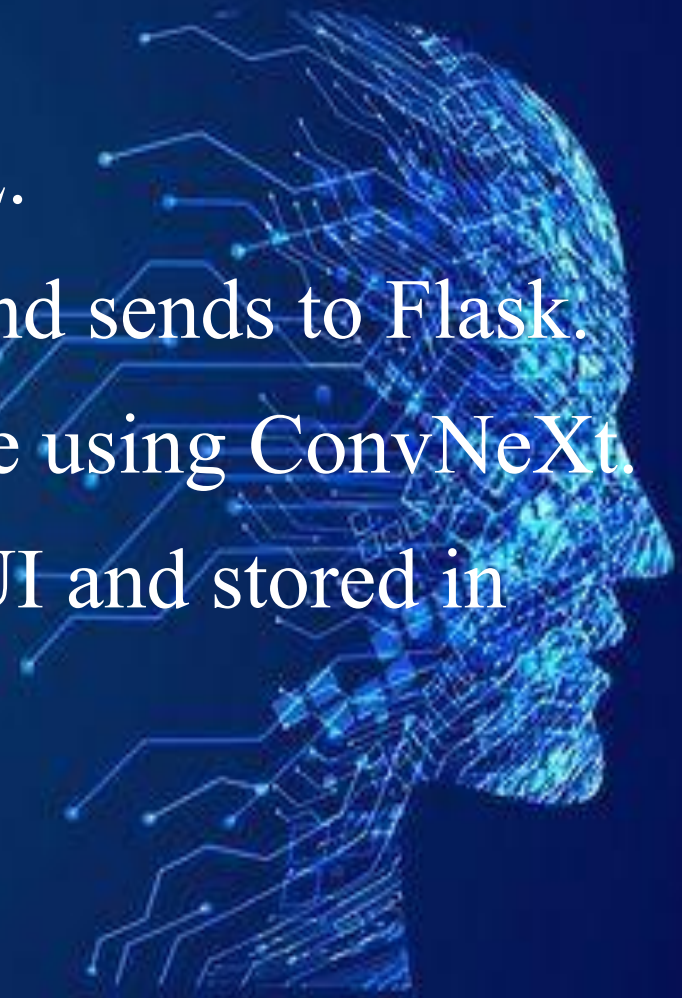
System Architecture

- React frontend: camera interface, emotion dashboard, responsive design.
- Flask backend: handles inference, preprocessing, and predictions using apis
- Model hosted via Kaggle/Colab with local Flask API for testing.
- MySQL handles secure login, session logs, analytics.



Integration Flow

- 1. User logs in via MySQL.
- 2. React captures frames and sends to Flask.
- 3. Flask performs inference using ConvNeXt.
- 4. Emotion is returned to UI and stored in MySQL.



DEMO



Code Repos

Code Model: [GitHub link](#)

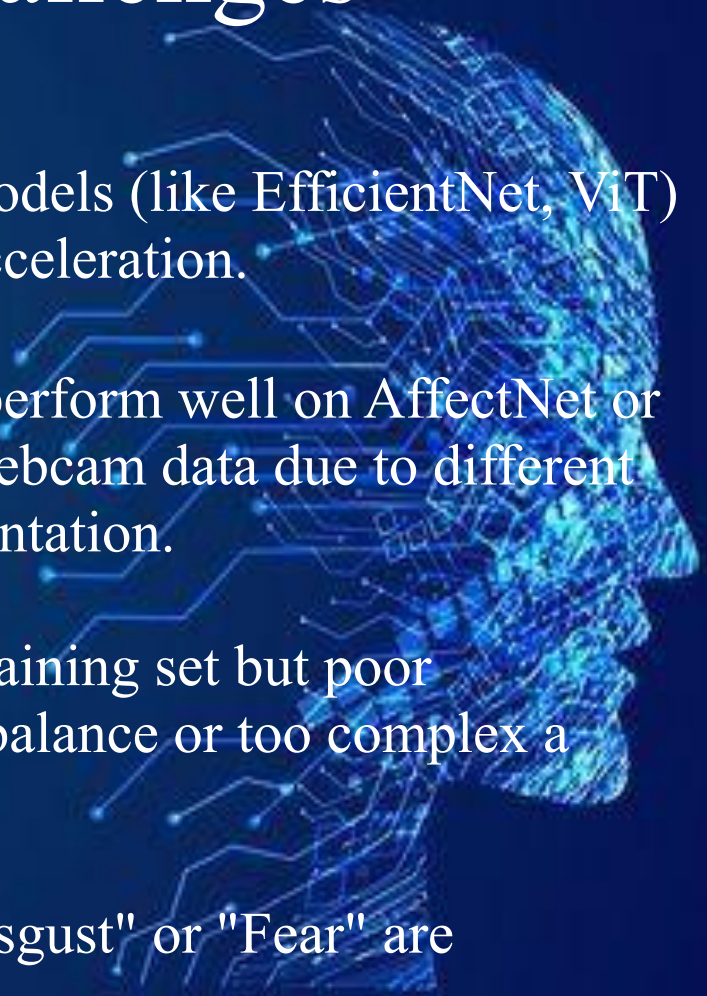
Code for web application: [GitHub Link](#)



Technical Challenges

1. Model-Related Challenges

- **Real-time inference speed:** Deep models (like EfficientNet, ViT) can be slow on CPU without GPU acceleration.
- **Generalization:** Your model might perform well on AffectNet or RAF-DB but poorly on real-world webcam data due to different lighting, resolution, or ethnic representation.
- **Overfitting:** High accuracy on the training set but poor validation/test results due to data imbalance or too complex a model.
- **Class imbalance:** Emotions like "Disgust" or "Fear" are underrepresented in many datasets.



Technical Challenges

2. Frontend Integration Challenges

- **Webcam handling:** Different browsers have slightly different permissions and APIs (especially on mobile).
- **Video capture and frame extraction:** Ensuring quality frames and preventing corruption when converting to blobs and sending to Flask.
- **Cross-browser compatibility:** Making sure webcam and video submission work on Safari, Firefox, Edge.



Technical Challenges

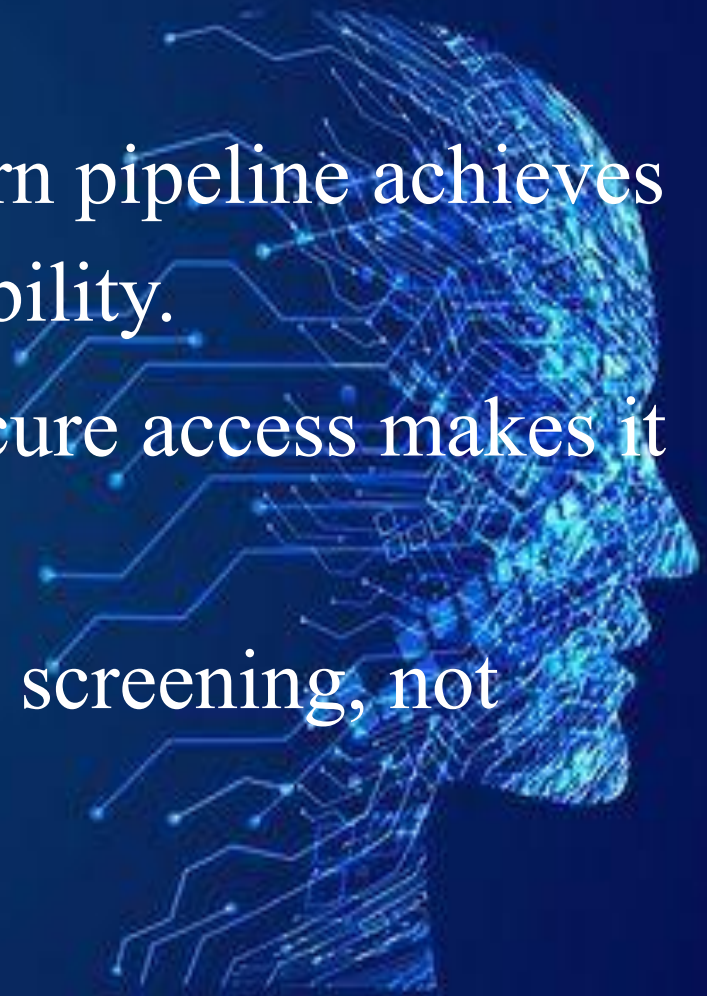
3. Backend/API Challenges

- **Handling large video files:** Flask may timeout or memory overload when processing large videos.
- **Security concerns:** Flask endpoints need validation/sanitization against bad payloads or misuse (e.g., sending corrupted files).
- **Concurrency:** Flask isn't multi-threaded by default; handling multiple users simultaneously requires Gunicorn + Nginx or other WSGI setups.



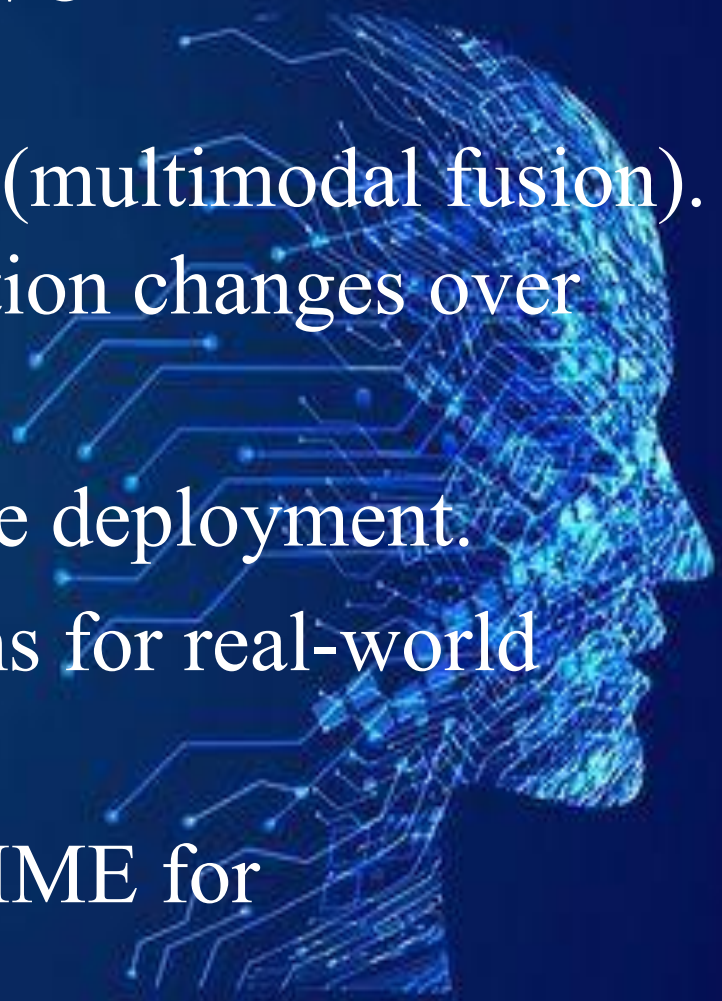
Conclusion

- ConvNeXt-Large + modern pipeline achieves strong accuracy and scalability.
- Real-time prediction + secure access makes it practical.
- Foundation for preventive screening, not diagnosis.



Future Work

- Add voice/text sentiment (multimodal fusion).
- Use ConvLSTM for emotion changes over time.
- ONNX/TensorRT for edge deployment.
- Collaborate with clinicians for real-world testing.
- Implement Grad-CAM/LIME for explainability.



Acknowledgments

- Faculty mentor: Prof. Ihab Darwish.
- City College of New York and CS Department.
- Tools: Kaggle, Google Colab, Firebase, GitHub, Overleaf.



THANK YOU

