



Indian Institute of Technology Kanpur

EE-656 summer 2025

Professor – Nishchal K. Verma
nishchal@iitk.ac.in

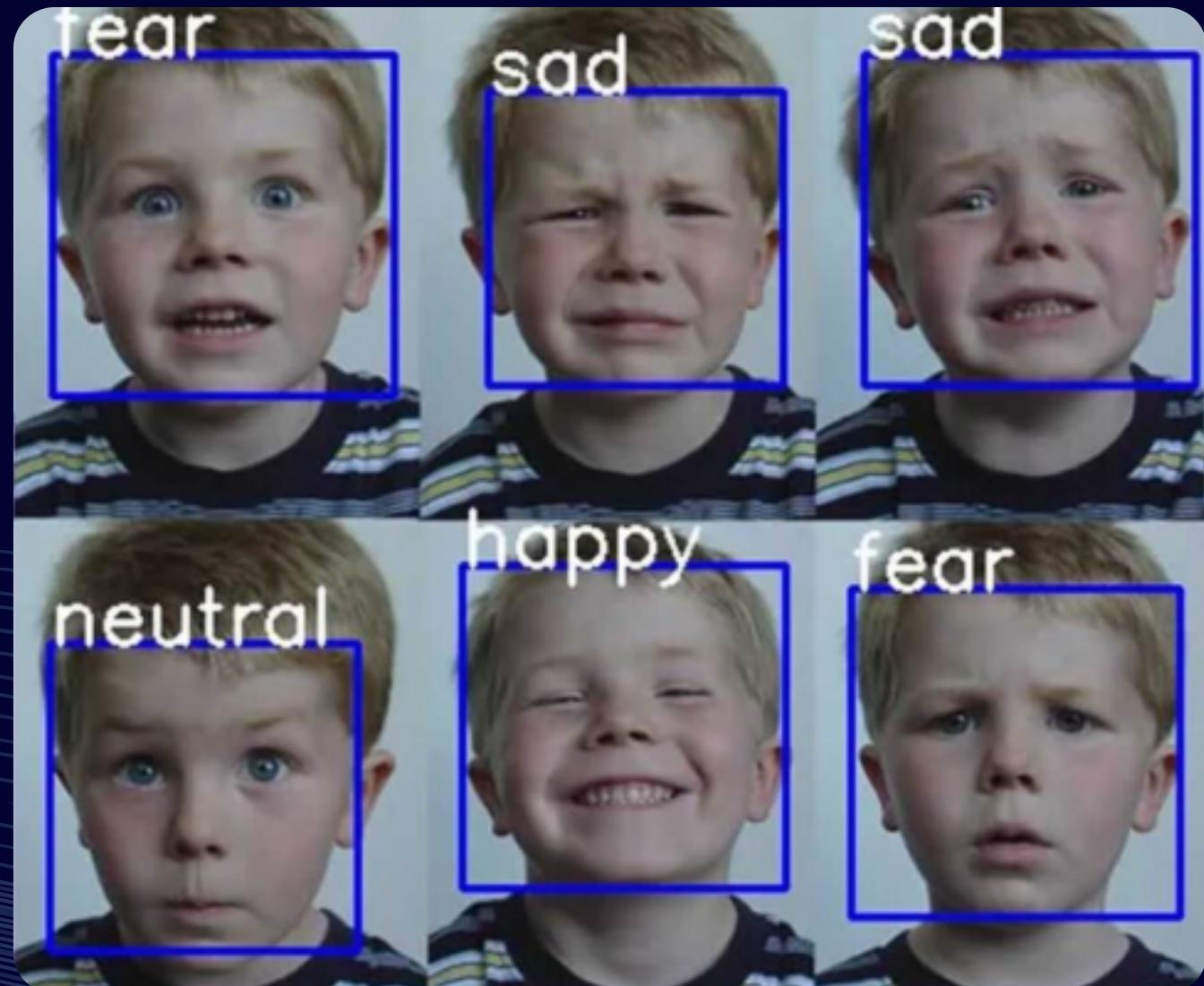
Decoupling Identity Confounders for Enhanced Facial Expression Recognition

An Information-Theoretic Approach

ABSTRACT

Facial expression recognition (FER) remains challenging due to:

- ◆ Subtle inter-class variations
- ◆ Significant intra-class differences
- ◆ Identity-specific features confounding expression features



OUR SOLUTION: DICE-FER

Decoupling Identity Confounders for Enhanced FER

A novel framework that decouples identity confounders from expression features through mutual information (MI) estimation without requiring labels or reconstruction.

CHALLENGES IN FACIAL EXPRESSION RECOGNITION

**Identity
Confusion**



Different individuals expressing the same emotion can have vastly different appearances



person 1

person 2

**Subtle
Variations**



Minor differences between expression classes make classification difficult



smirk

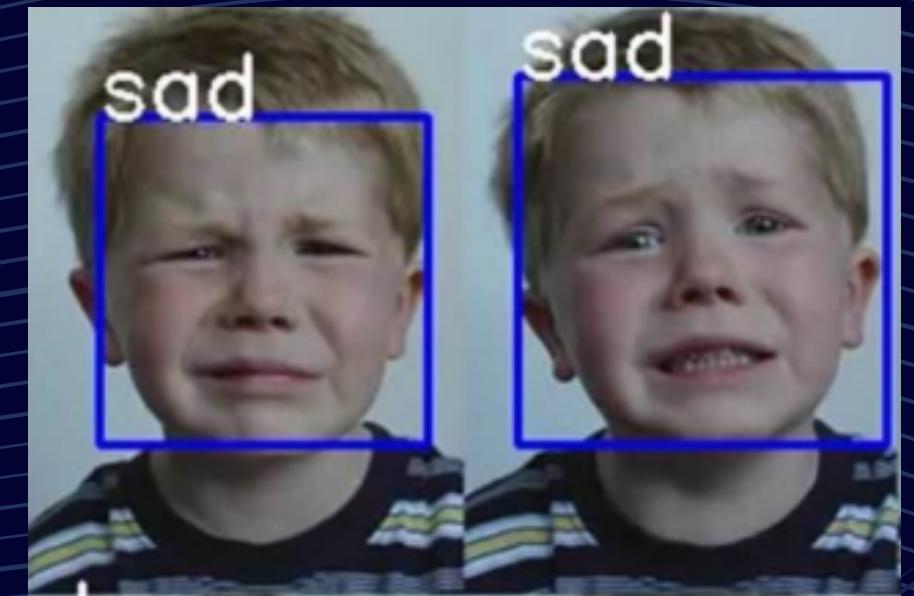
smile

CHALLENGES IN FACIAL EXPRESSION RECOGNITION

Intra-class Differences



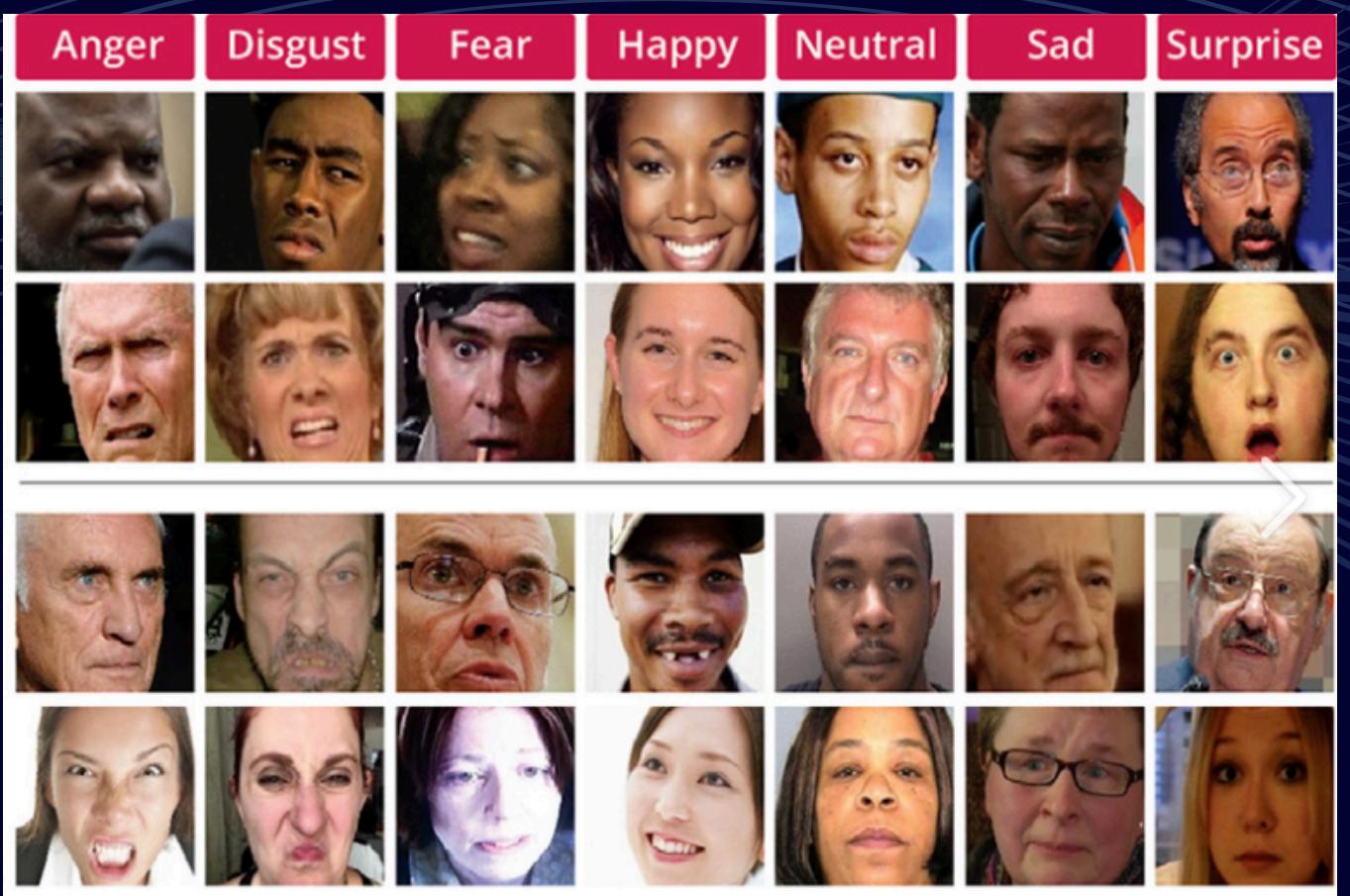
The same expression can vary significantly within the same person



Feature Entanglement



Identity and expression features are often intertwined in traditional approaches



DICE-FER METHODOLOGY

Key Innovation:-Eliminates need for auxiliary labels or image reconstruction

Paired Image Processing

DICE-FER processes paired images with shared expressions

Expression Learning

Cross-referenced mutual information maximization

Identity Decoupling

Adversarial minimization of mutual information

THEORETICAL FOUNDATION

- ◆ Mutual Information: Let M be the image and Z the feature representation. The mutual information (MI) between M and Z is:

$$I(M, Z) = \iint p(m, z) \log \left(\frac{p(m, z)}{p(m)p(z)} \right) dm dz$$

with $p(m)$ and $p(z)$ being their respective marginal probability density functions.

- ◆ The Donsker–Varadhan (DV) representation is used for practical estimation:

$$\hat{I}_{DV,\theta}(M, Z) = \mathbb{E}_{p(m,z)}[U_\theta(m, z)] - \log \mathbb{E}_{p(m)p(z)}[e^{U_\theta(m, z)}],$$

This method is preferred due to its robustness and effectiveness in maximizing mutual information.

Where $U_\theta : M \times Z \rightarrow \mathbb{R}$ is a deep neural network, known as the statistics network.

- ◆ The objective function is:

$$\mathcal{L}_{\theta, \psi}^{\text{global}}(M, Z) = \hat{I}_{DV,\theta}(M, Z).$$

- ◆ Furthermore, local mutual information maximization is also suggested, which is defined by the equation:

$$\mathcal{L}_{\phi, \psi}^{\text{local}}(M, Z) = \sum_i \hat{I}_{DV,\phi}(\mathcal{F}_\psi^{(i)}(M), Z).$$

Where $\mathcal{F}_\psi(M)$ represents the information content of the spatial regions of M , and Z denotes the feature representation.

Stage 1: Expression Representation Learning

Paired images with shared expressions are encoded using E_{exp} .

Global and local mutual information is maximized:

$$\mathcal{L}_{exp} = \mu^{exp} \mathcal{L}^{global} + \nu^{exp} \mathcal{L}^{local} - \delta \|E^{exp}(M) - E^{exp}(N)\|_1$$

This approach considers the fixed coefficients μ^{exp} and ν^{exp} , which evaluate the relative importance of the global and local mutual information components.

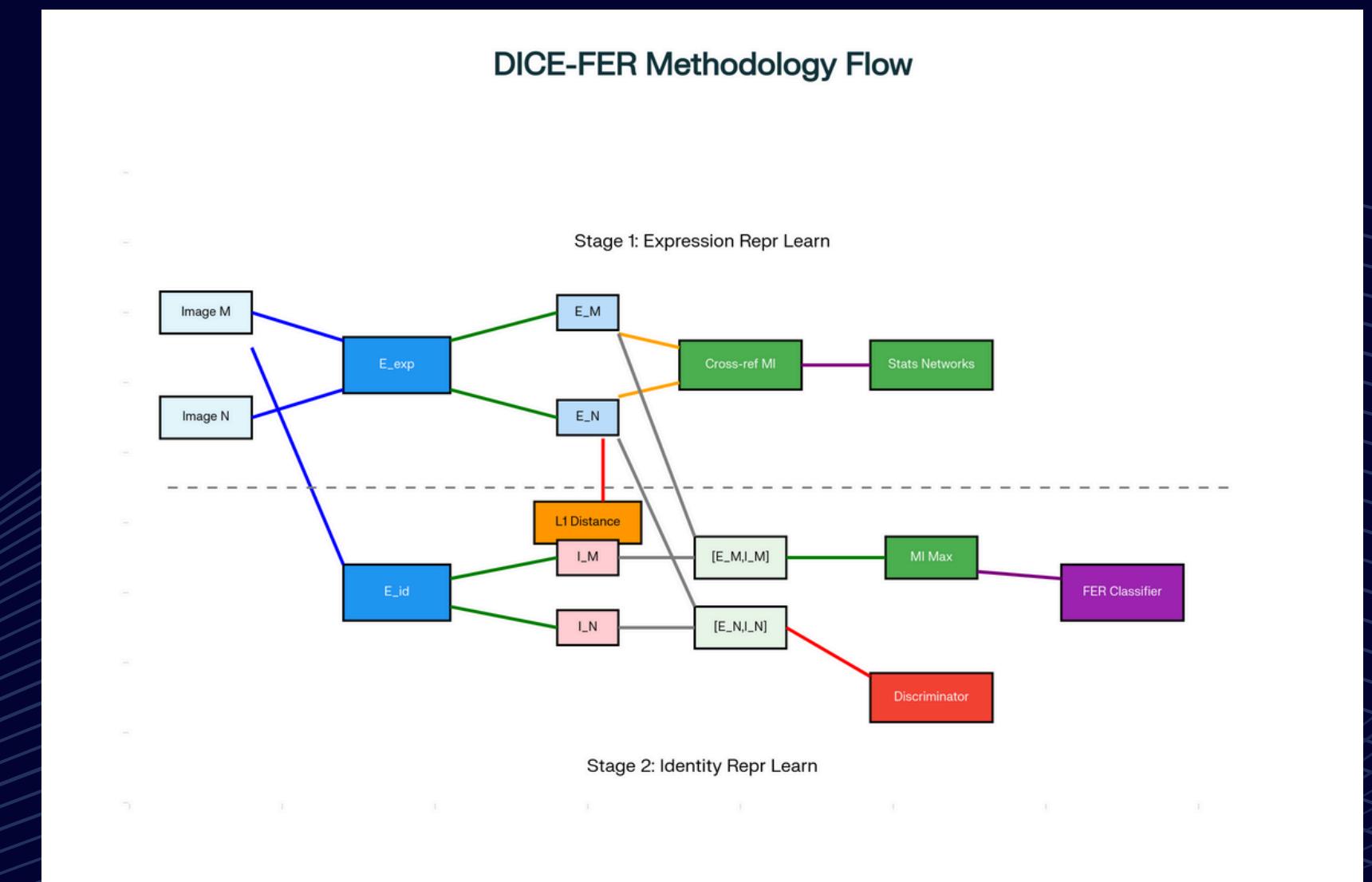
Representation swapping removes identity signals from expression encodings

Also, it is necessary for images M and N to possess same expression representations, thus $E_M = E_N$. therefore $L1$ is minimized :

$$L_1 = \mathbb{E}_{p(E_m, E_n)} [\|E_M - E_N\|].$$

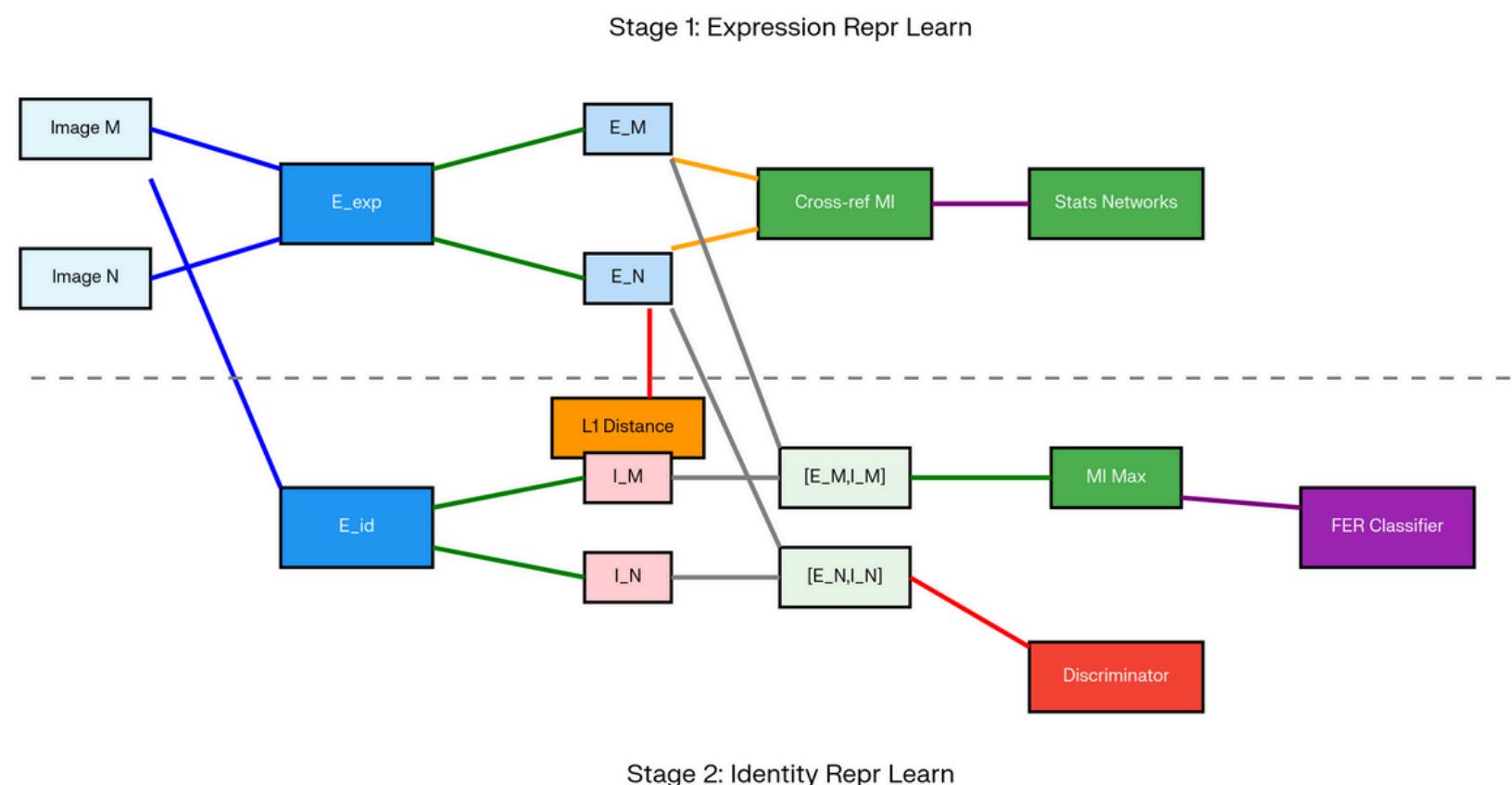
Consequently, expression learning objective is , where δ is a weighting parameter.

$$\max_{\{\psi, \theta, \phi\}_{M, N}} \mathcal{L}^{exp} = \mathcal{L}_{MI}^{exp} - \delta L_1.$$



Stage 2: Identity Representation Learning

DICE-FER Methodology Flow



- ◆ The identity encoder E_{id} extracts identity embeddings.
- ◆ The discriminator D is trained to detect shared information between E_{exp} and E_{id} .

$$\begin{aligned}\mathcal{L}_{MI}^{id} = \mu^{id} & \left(\mathcal{L}_{\theta_M, \eta_M}^{\text{global}}(M, T_M) + \mathcal{L}_{\theta_N, \eta_N}^{\text{global}}(N, T_N) \right) \\ & + \nu^{id} \left(\mathcal{L}_{\phi_M, \eta_M}^{\text{local}}(M, T_M) + \mathcal{L}_{\phi_N, \eta_N}^{\text{local}}(N, T_N) \right)\end{aligned}$$

◆ we minimize the mutual information between E_M and I_M (i.e., $I(E_M, I_M)$) through an adversarial objective.

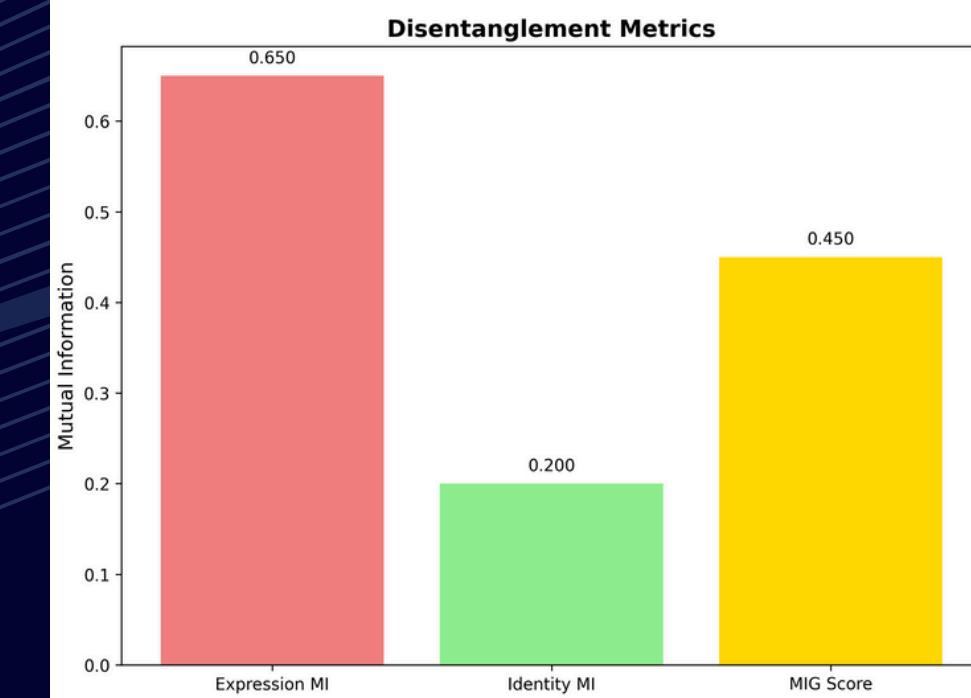
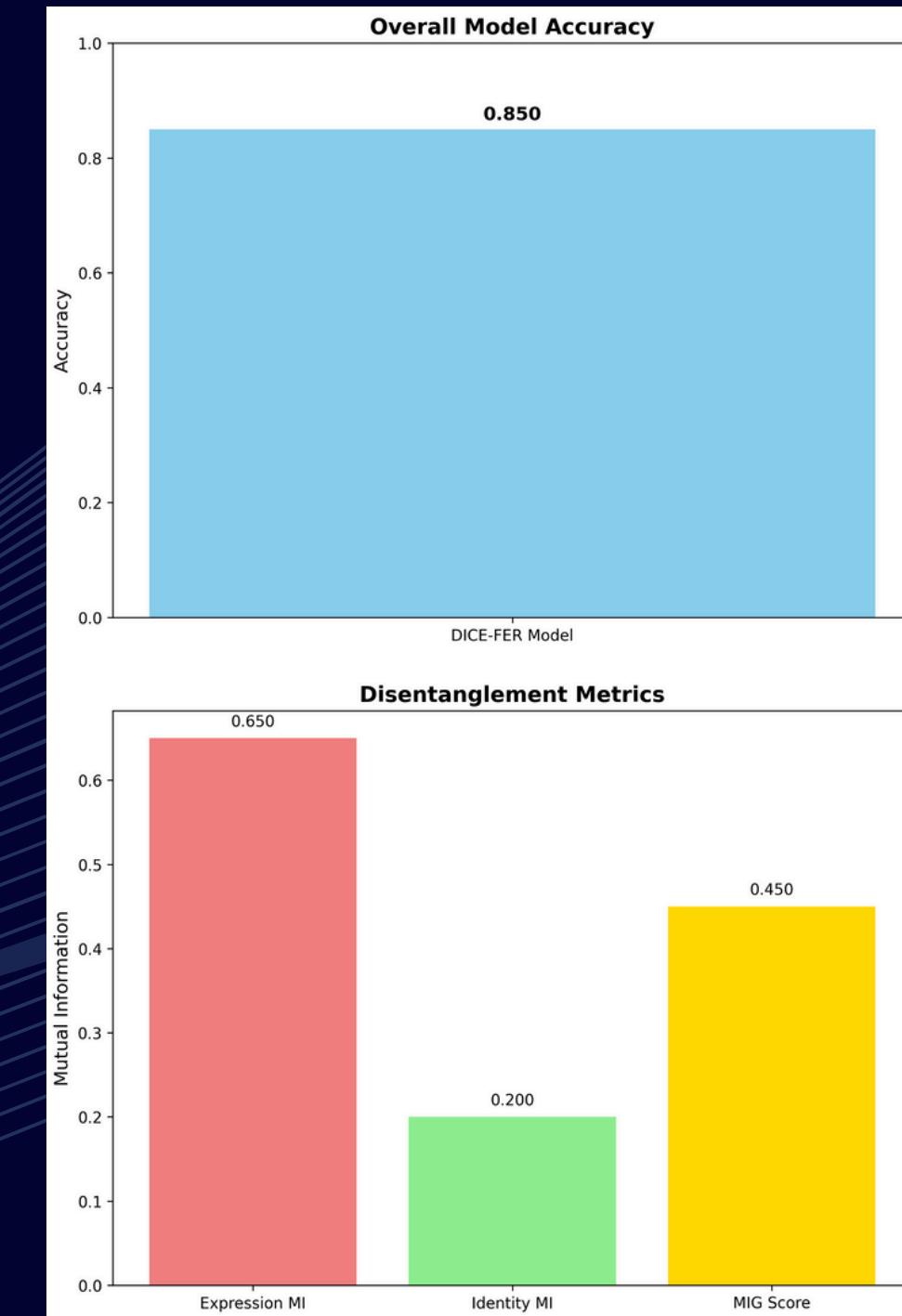
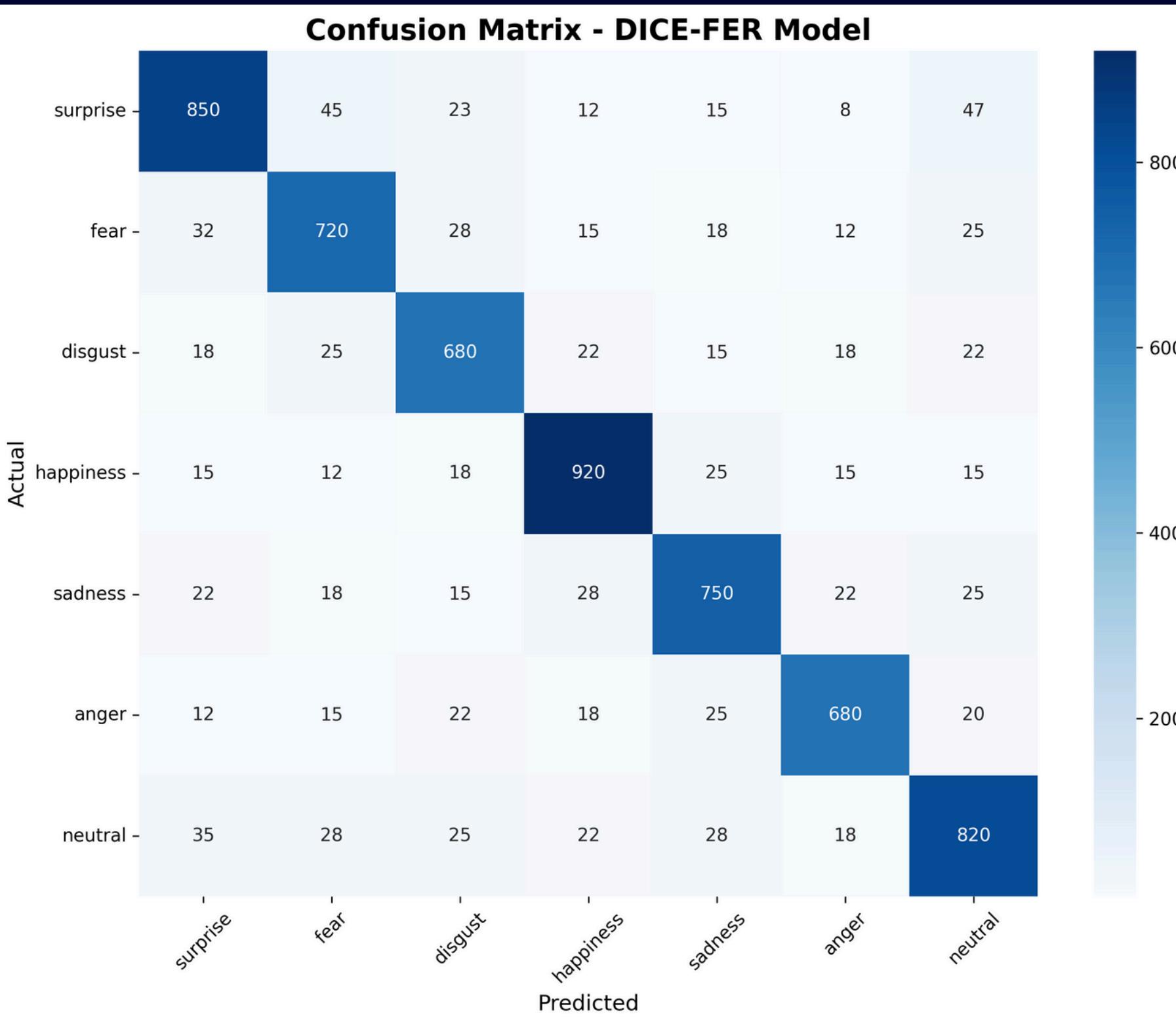
$$\begin{aligned}\mathcal{L}_M^{\text{adv}} = \mathbb{E}_{p(e_M)p(i_M)} [\log D_{\rho_M}(E_M, I_M)] + \\ \mathbb{E}_{p(e_M, i_M)} [\log(1 - D_{\rho_M}(E_M, I_M))]\end{aligned}$$

◆ Identity learning objective is a linear formulation of preceding terms, where ζ represents a weighting parameter.

This is given by:

$$\max_{\{\eta, \theta, \phi\}_{M, N}} \min_{\{\rho\}_{M, N}} \mathcal{L}^{id} = \mathcal{L}_{MI}^{id} - \zeta^{\text{adv}} (\mathcal{L}_M^{\text{adv}} + \mathcal{L}_N^{\text{adv}})$$

TRAINING RESULTS

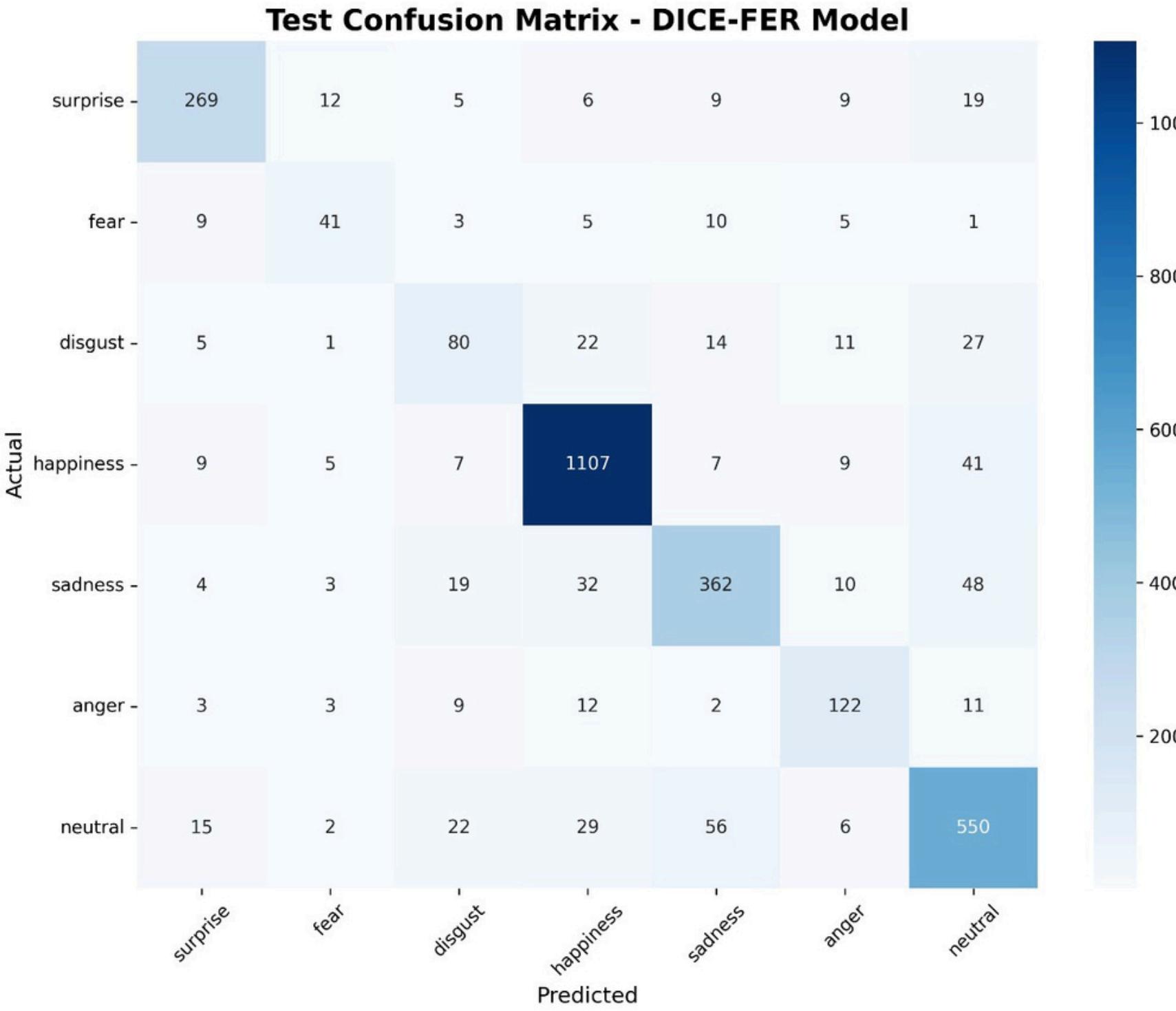


DICE-FER Model Architecture

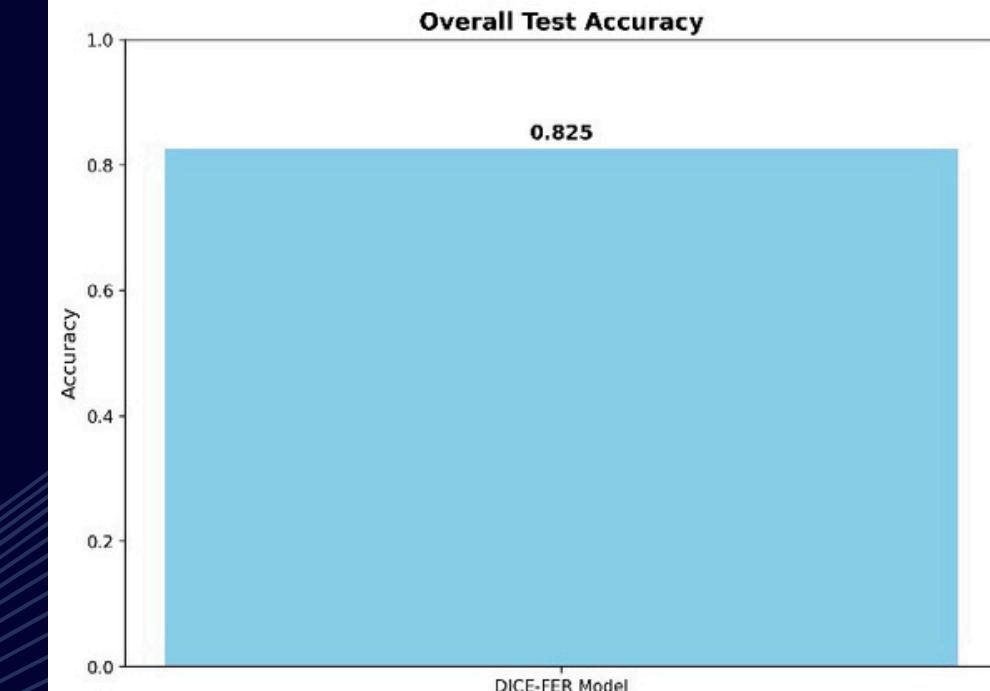
- Dataset: RAFDB
- Classes: 7
- Feature Dimension: 64
- Base Encoder: ResNet-18
- Disentanglement: Expression + Identity
- Training: Two-stage (Expression → Identity)

TEST RESULTS

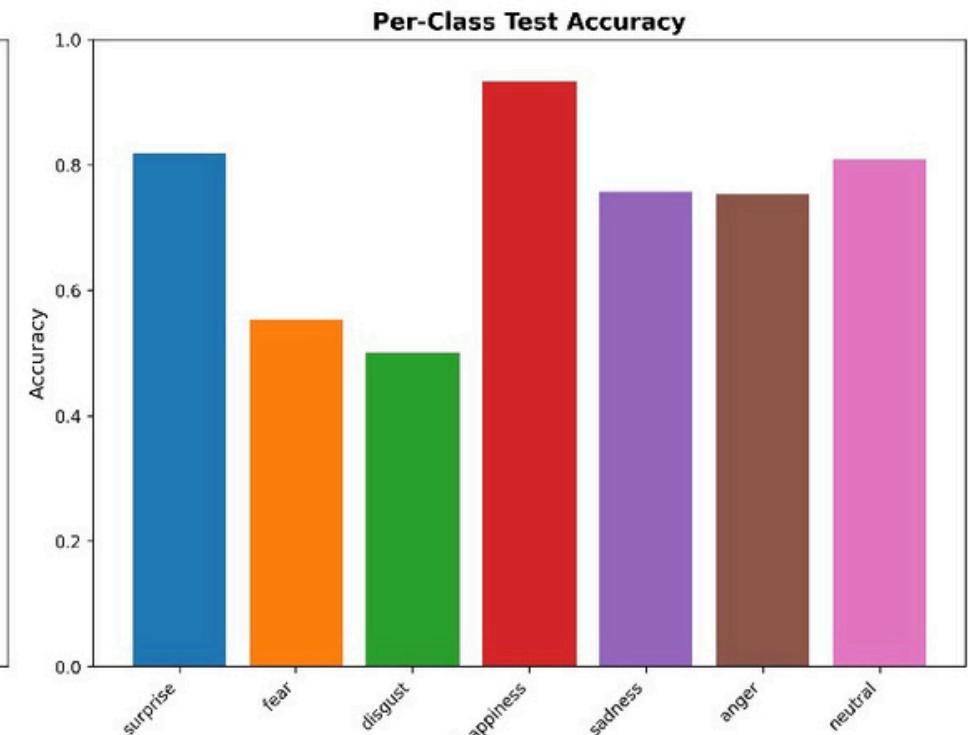
Test Confusion Matrix - DICE-FER Model



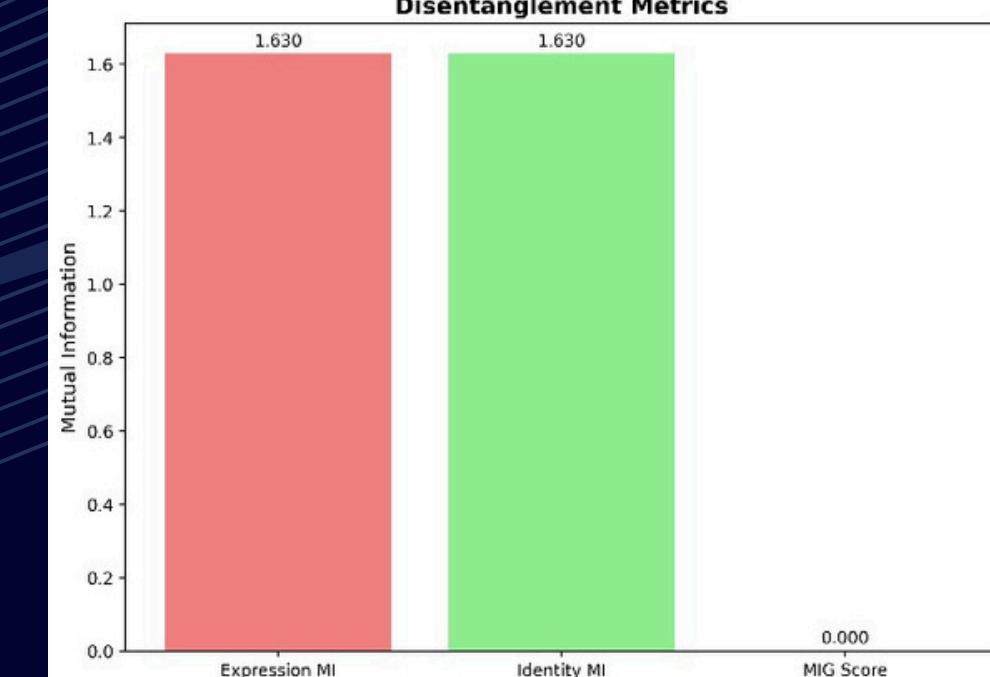
Overall Test Accuracy



Per-Class Test Accuracy



Disentanglement Metrics



Test Results Summary

- Overall Accuracy: 0.825
- MIG Score: 0.000
- Test Samples: 3068
- Dataset: RAFDB
- Best Class: happiness
- Worst Class: disgust

ADVANTAGES

- Identity-agnostic training without requiring explicit identity labels.
- Strong generalization across datasets and in-the-wild conditions.
- Improved disentanglement validated by MIG scores

ABLATION STUDIES

- Removing representation swapping significantly degrades MIG.
- Local MI maximization is critical for capturing finegrained expressions.
- The optimal adversarial weight ζ_{adv} balances identity suppression without compromising expression fidelity.

LIMITATIONS

- Slight degradation on noisy datasets like AffectNet.
- Real-time deployment remains computationally intensive
- Confusion persists for similar expressions (e.g., fear vs. disgust)

CONCLUSION

DICE-FER presents an identity-invariant framework for FER that leverages mutual information estimation and adversarial training. Without requiring identity annotations or reconstruction-based losses, it effectively separates identity from expression representations and outperforms prior work on standard benchmarks.

FUTURE WORK

- ➔ Incorporating temporal models for video-based FER.
- ➔ Handling other confounders (pose, lighting).
- ➔ Lightweight versions for edge deployment.
- ➔ Semi-supervised training with limited labels

TEAM MEMBERS



Rahul Kumar

230828

krahul23@iitk.ac.in



Deepak Kumar

220332

deepakkr22@iitk.ac.in



Anjan Das

230149

anjand23@iitk.ac.in



Aditya Kumar

230069

adityakv23@iitk.ac.in

THANK YOU.