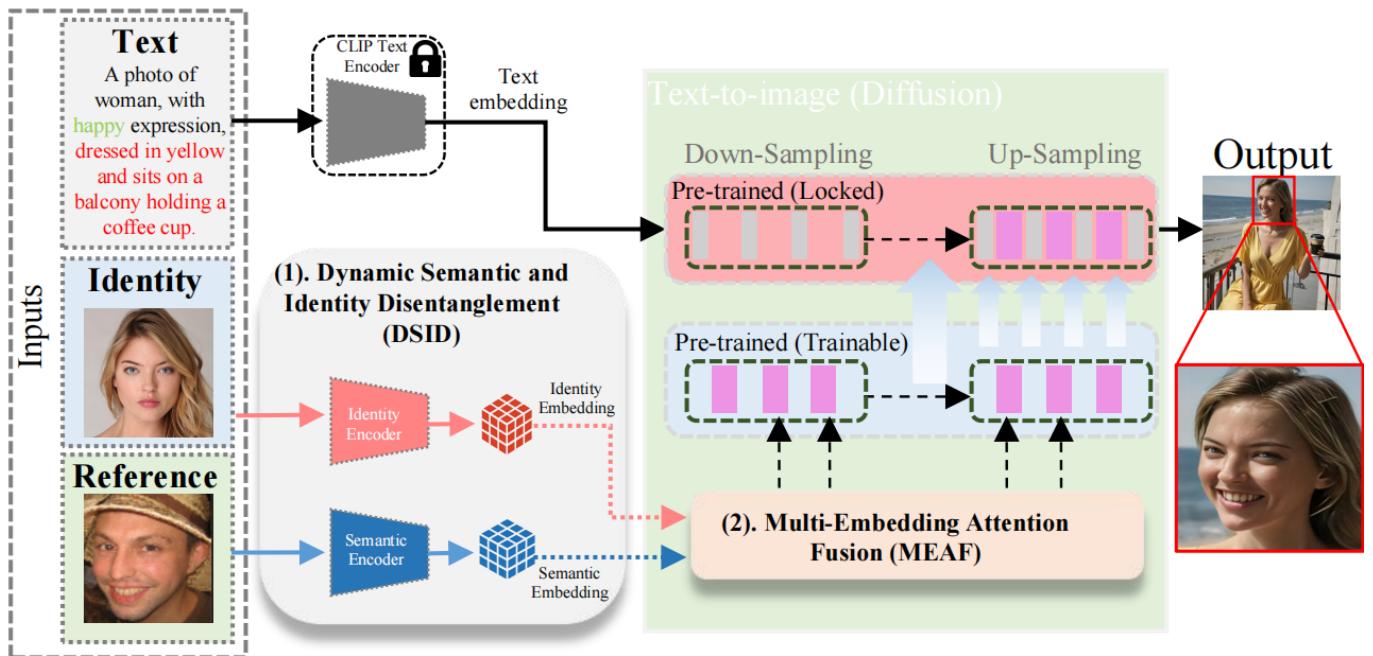


# DeepExPo

## Facial Expression and Pose Generation via Self-Supervised Disentangled Embeddings Fusion in Text-to-Image Diffusion Models

### 🔍 Overview

**DeepExPo** is a framework designed to provide fine-grained control over facial expression and head pose in text-to-image diffusion models, while maintaining identity preservation. Existing diffusion-based personalization methods typically fall short in expression control and pose alignment. DeepExPo overcomes these limitations through a self-supervised disentanglement strategy and adaptive embedding fusion with the U-Net of pre-trained text-to-image models.



**Description:** DeepExPo generates identity-preserving facial images with controllable expressions and head poses using three inputs: a text prompt for context, an identity image, and a reference image for semantic cues. Identity and semantic features are fused via the MEAF module into a pre-trained diffusion model.

Table 1: Capability comparison of different methods based on identity preservation (ID), expression control (EXP), pose manipulation (POSE), and context integration (CONTEXT). Ratings are derived from quantitative and qualitative metrics reported in the literature, as well as from our own empirical experiments.

Method	ID.	EXP.	POSE.	CONTEXT
GAN-based [7, 8]	High	High	High	None
Textual Inversion [9]	High	None	None	High
DreamBooth [10]	High	None	None	High
DisenBooth [11]	Moderate	Low	None	High
Celeb-Basis [12]	Moderate	Low	None	High
CIP-TIG [13]	Moderate	Low	None	High
MasterWeaver [14]	High	Moderate	None	High
EmojiDiff [15]	Moderate	Moderate	Low	High
DiffSFSR [16]	Moderate	Moderate	Moderate	High
<b>DeepExPo</b>	<b>High</b>	<b>High</b>	<b>High</b>	<b>High</b>

**Description:** Capability comparison of various methods based on identity preservation, expression control, pose manipulation, and context integration. Ratings are based on literature benchmarks and our own experimental analysis.

DeepExPo introduces a two-stage solution:

- **Dynamic Semantic and Identity Disentanglement (DSID) module**
- **Multi-Embedding Attention Fusion (MEAF) module**

This repository includes the full implementation of the **DeepExPo**.

## 🧠 Key Modules

### 1. Dynamic Semantic and Identity Disentanglement (DSID) Module

📦 `./modules/dsid/`

A self-supervised module designed to disentangle identity and dynamic semantic attributes (e.g., expressions and head pose) from single or paired image frames. This module:

- Leverages consecutive video frame pairs as implicit supervision to separate identity from both static and temporal variations
- Encodes identity and dynamic semantics into two independent latent spaces
- Enables composable conditioning on expressions and pose while preserving identity fidelity

### 2. Multi-Embedding Attention Fusion (MEAF) Module

The **MEAF** module employs a parallel multi-attention mechanism to fuse **semantic** and **identity embeddings** with intermediate U-Net features from a pre-trained diffusion model. This approach:

- Preserves **identity** and **facial fidelity**

- Avoids degradation typically caused by direct embedding injection
  - Enables effective conditioning without disrupting spatial coherence
- 

## Installation

To use **DeepExPo**, clone the repository and install the required dependencies.

### 1. Clone the Repository

```
git clone https://github.com/MSAfganUSTC/DeepExPo.git  
cd DeepExPo
```

After cloning, the folder contains the following structure:

### 2. Repository Structure

```
DeepExPo/  
|  
|   -- requir/  
|   |   -- config.yaml          # Configuration file  
|  
|   -- DeepExPo_DSID_Module/  
|   -- DeepExPo_MEAF_Module/  
|  
|   -- scripts/  
|   |   -- DeepExPo_Inference.ipynb    # Inference script  
|  
|   -- DeepExPo_Weights/  
|   |   -- DSID_Checkpoints/  
|   |   -- DeepExPo_MEAF_weights/      # Pretrained model weights  
|   |  
|   |   -- DSID_Checkpoints/          # Checkpoints for DSID module  
|   |   -- DeepExPo_MEAF_weights/     # Weights for MEAF module  
|  
|   -- Images/                      # Figures and illustrations used in  
the paper  
|   -- samples/  
|   |   -- subject.jpg            # Sample input image  
|  
|   -- LICENSE                      # License file  
|   -- README.md                    # Project documentation
```

### 3. Set Up the Environment

This project uses Conda for environment management. Make sure you have Conda installed.

Create the environment from the provided file:

```
conda env create -f configs/environment.yml  
conda activate DeepExPo
```

The dependencies include PyTorch, Hugging Face Transformers, and other necessary packages.

## 4. Usage

### Interactive Inference (Jupyter Notebook)

All the steps for loading a model, choosing an expression, and generating images are wrapped in an easy-to-run notebook.

#### 1. Activate the environment

```
conda activate DeepExPo    # or mamba activate DeepExPo
```

#### 2. Launch Jupyter and open the notebook

```
jupyter notebook scripts/DeepExPo_Inference.ipynb
```

#### 3. Run the cells from top to bottom

- The first cell lets you set paths like `MODEL_PATH` to load DSID checkpoints (both identity encoder and semantic encoder).
- The second cell asks for identity and reference images to extract embeddings, and lets you specify the `output_dir`.
- Later cells handle loading the model, running inference, and saving the results.
- Generated images will appear inside the notebook and in the output directory you specified. Below is a snapshot of the Jupyter Notebook interface used for inference:

This cell imports necessary libraries and sets the device (GPU or CPU). It also defines the file paths to the three pre-trained model checkpoints that will be loaded later.

```
[1]: import torch
from torchvision import transforms
from PIL import Image
import os

# --- Paths to model checkpoints ---
identity_encoder_path = r"DeepExPo/DeepExPo_Weights/DSID_Checkpoints/dsid-source-identity-cpk.pth"
semantic_encoder_path = r"DeepExPo/DeepExPo_Weights/DSID_Checkpoints/dsid-target-semtic-cpk.pth"
MEAF_model_path = r"DeepExPo/DeepExPo_Weights/DeepExPo_MEAF_weights.pth" # add extension if needed

# --- Device ---
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
print(f"Using device: {device}")

Using device: cuda
```

## Define or Import Model Architectures

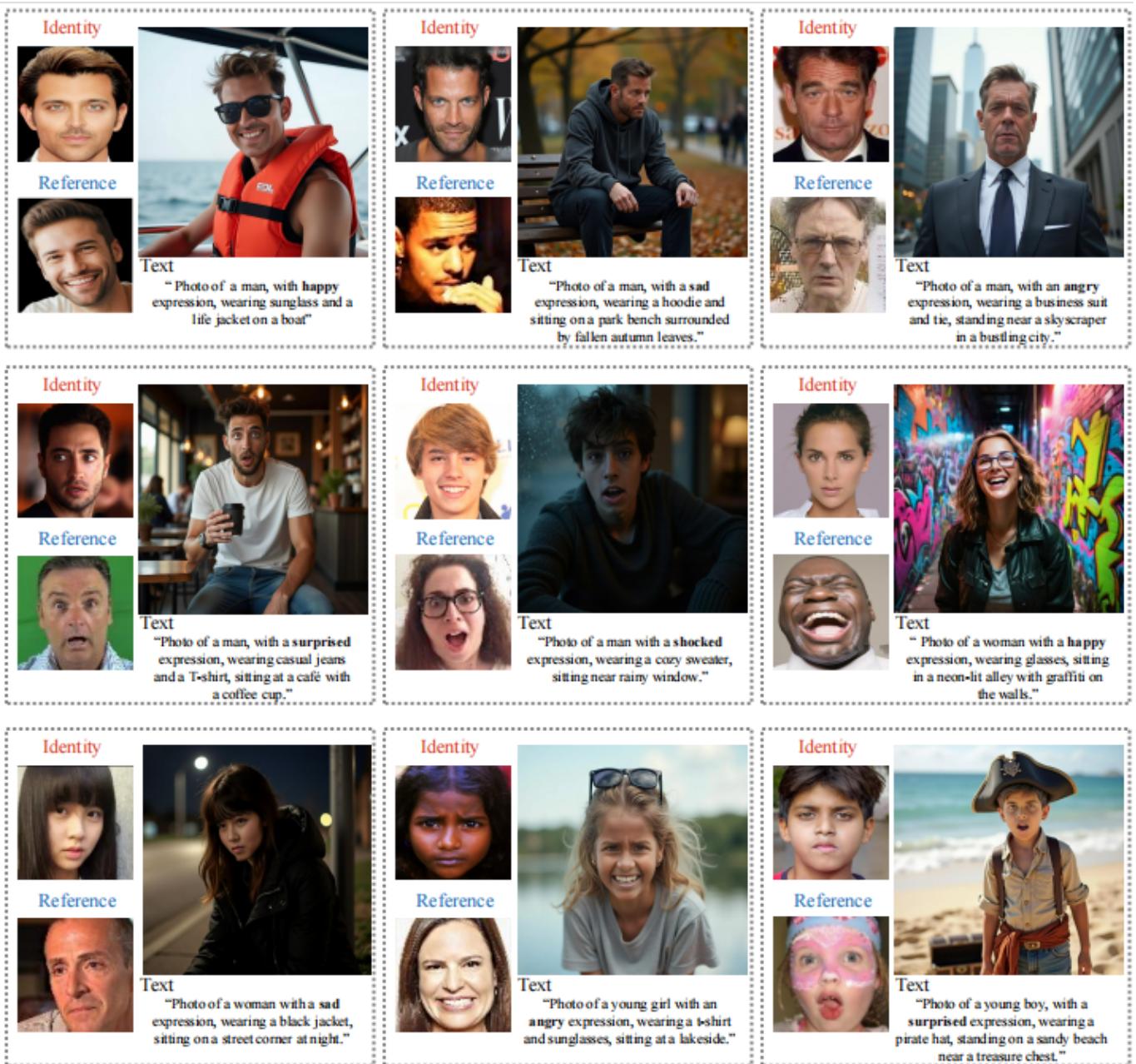
Here we define placeholder PyTorch model classes for the Identity Encoder, Semantic Encoder, and the MEAF Model. You should replace these with your actual model definitions or imports.

```
[ ]: class DSID(torch.nn.Module):
    def __init__(self):
        super().__init__()
        self.identity_encoder = ...
        self.semantic_encoder = ...

    def forward(self, x_id, x_sem):
        id_emb = self.identity_encoder(x_id)
        sem_emb = self.semantic_encoder(x_sem)
        return id_emb, sem_emb
```

## Load Pre-trained Model Weights

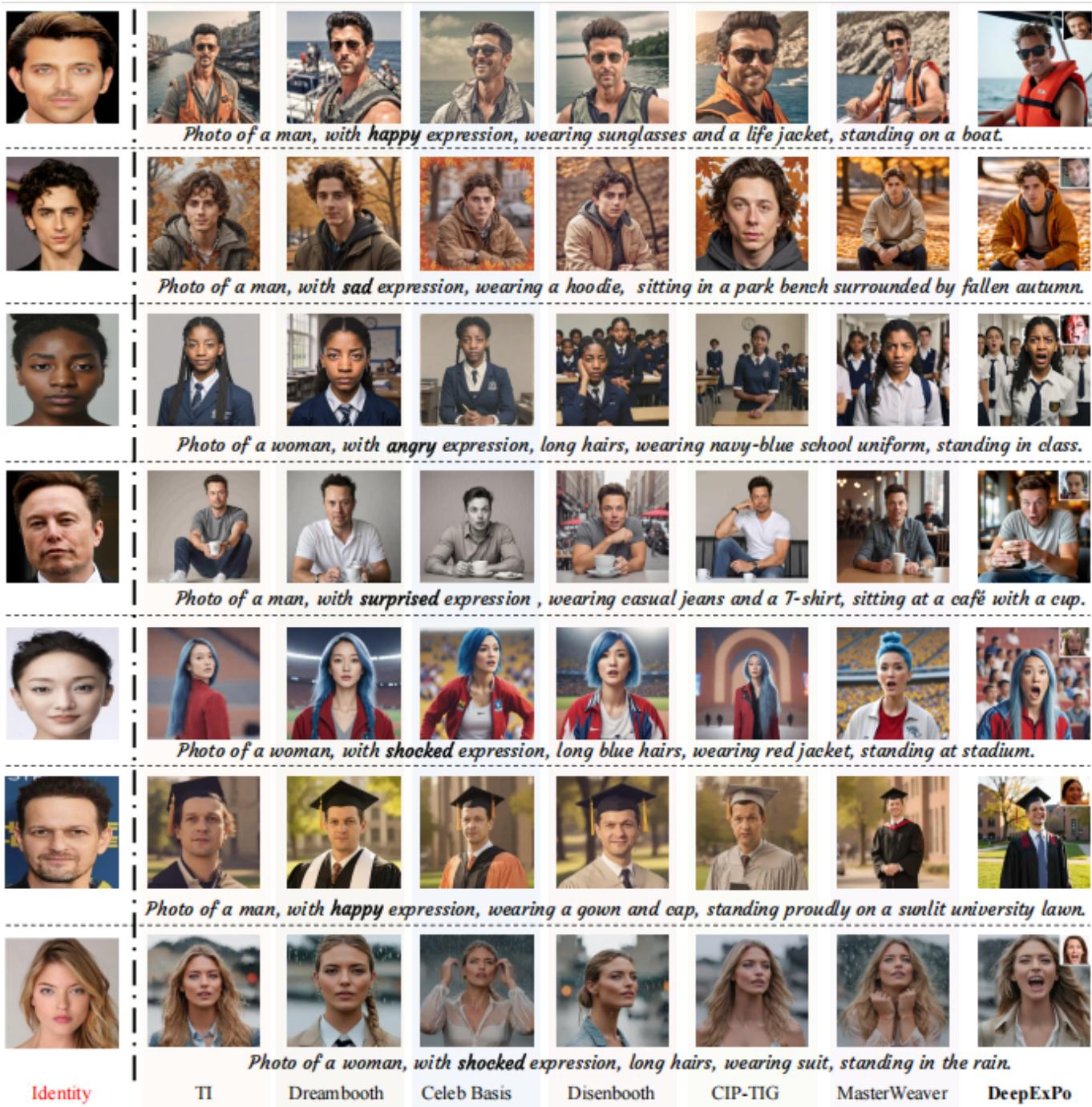
## Inference Results



**Description:** DeepExPo inference samples showing accurate expression and pose synthesis while maintaining subject identity and contextual fidelity.

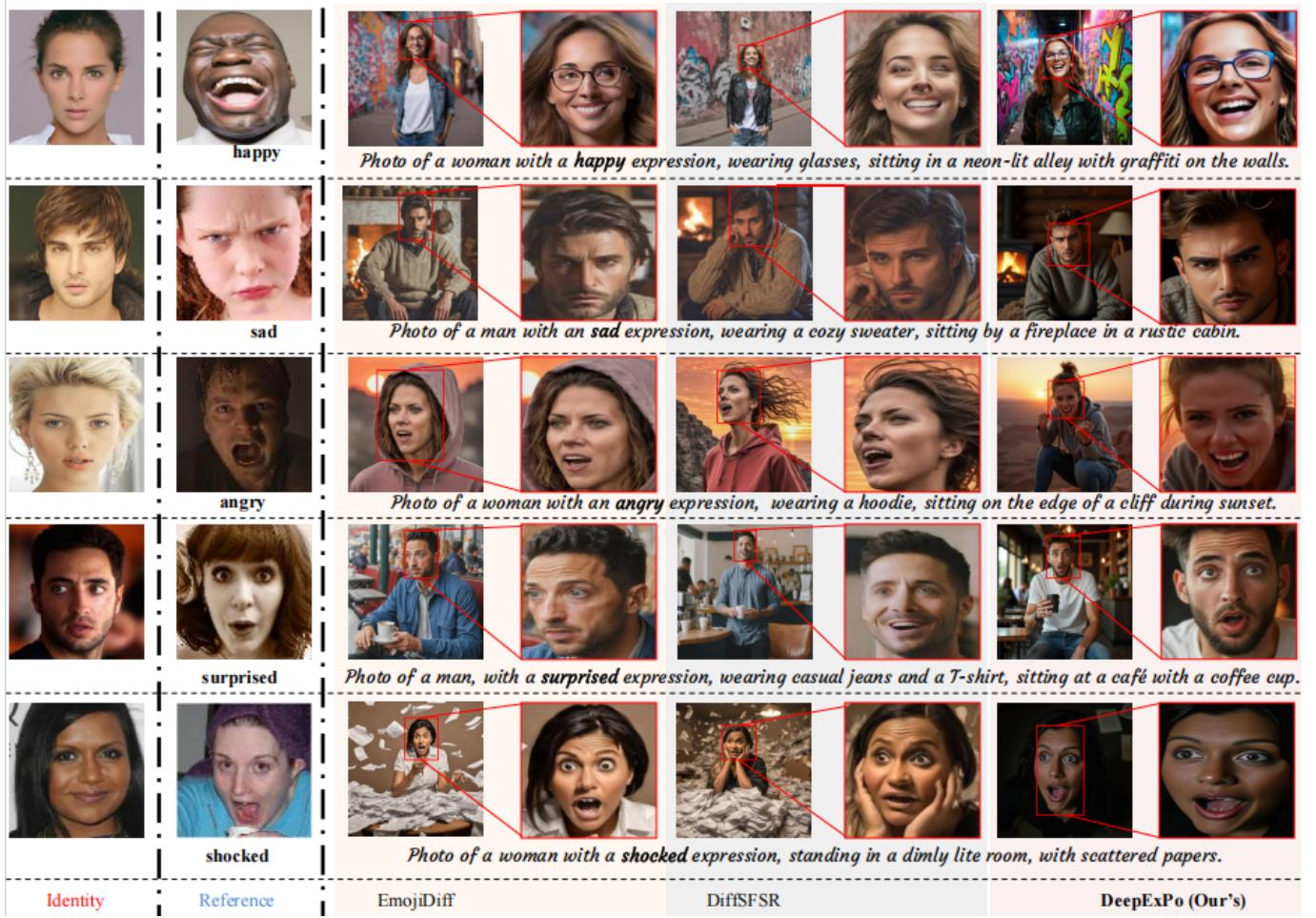
## Comparison with Baselines

### Personalized Generation Baselines



**Description:** Qualitative comparison against personalized generation models. **DeepExPo** maintains strong identity consistency while effectively transferring facial semantics and adhering to contextual prompts.

## 😊 Expression Generation Baselines



**Description:** DeepExPo achieves better identity preservation and precise semantic transfer (e.g., mouth and eye region alignment) compared to other expression generation techniques.

## Quantitative Evaluation

Table 3: Quantitative comparison against personalized and expression generation baselines. **Bold** indicates the best, and underline the second-best performance for each metric.

	Method	ID. $\downarrow$	SEMCHG. $\uparrow$	SIM. $\downarrow$	EXP.(%) $\uparrow$	POSE $\uparrow$	FID $\downarrow$	TIME $\downarrow$
Personalized	TI [9]	0.48	0.35	-	30	-	128.3	10s
	DreamBooth [10]	<b>0.25</b>	0.20	-	-	-	<u>119.5</u>	14s
	DisenBooth [11]	0.47	0.42	-	35	-	121.0	14s
	Celeb Basis [12]	0.56	0.40	-	38	-	122.5	20s
	CIP-TIG [13]	0.85	<u>0.70</u>	-	<u>80</u>	-	124.3	22s
	MW [14]	0.53	0.66	-	75	-	124.9	<u>8s</u>
Express	<b>DeepExPo</b>	<u>0.40</u>	<b>0.80</b>	<b>0.29</b>	<b>90</b>	<b>2.50</b>	<b>119.21</b>	<b>4s</b>
	EmojiDiff [15]	0.57	0.68	0.45	85	4.59	124.3	6s
	DiffSFSR [16]	<u>0.50</u>	<u>0.77</u>	<u>0.38</u>	<u>87</u>	<u>3.00</u>	<u>122.8</u>	<u>5s</u>
	<b>DeepExPo</b>	<b>0.40</b>	<b>0.80</b>	<b>0.29</b>	<b>90</b>	<b>2.50</b>	<b>119.21</b>	<b>4s</b>

**Description:** Quantitative evaluation highlights DeepExPo's superior performance in identity preservation and expression accuracy. Qualitative results emphasize expression realism, head pose accuracy, and facial fidelity.

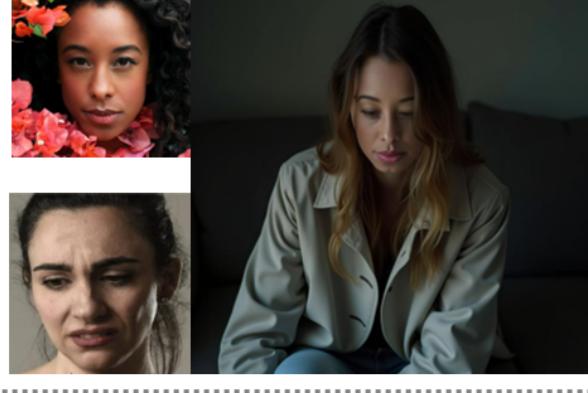
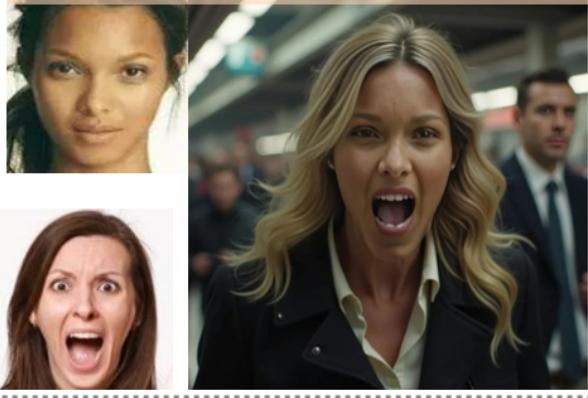
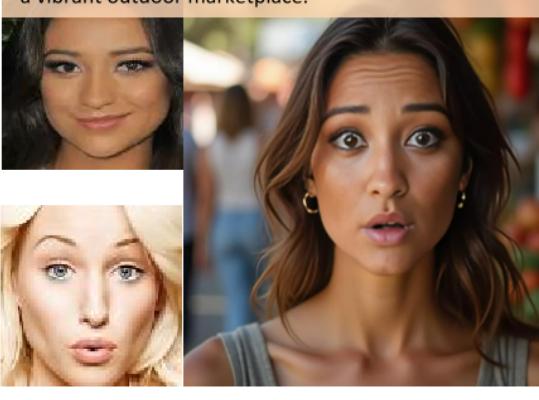
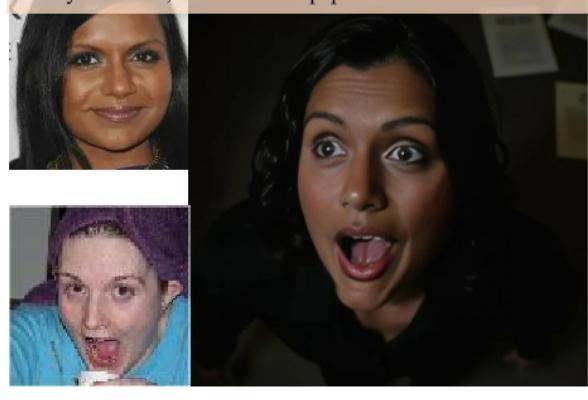
## 🎯 Additional Results

### 👤 Male Subjects

<p>Text Identity</p> <p>Photo of a man, with a <b>happy</b> expression, wearing sunglasses and a life jacket on a boat.</p>  	<p>Text Identity</p> <p>“ Photo of a man with <b>sad</b> expression, wearing gray hoodie driving a car”</p>  
<p>Text Identity</p> <p>Photo of a man, with an <b>angry</b> expression, wearing a business suit and tie, standing near a skyscraper in a bustling city.</p>  	<p>Text Identity</p> <p>Photo of a man, with an <b>angry</b> expression, wearing a business suit and tie, standing near a skyscraper in a bustling city.</p>  
<p>Text Identity</p> <p>Photo of a man with <b>surprised</b> expression, wearing a suit, —standing in studio”</p>  	<p>Text Identity</p> <p>Photo of a man, with a <b>shocked</b> expression, wearing a gray T-shirt and a hiking backpack, standing on a mountain trail.</p>  

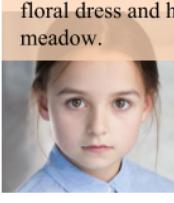
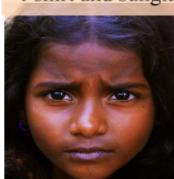
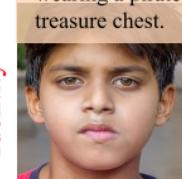
DeepExPo successfully generates realistic male facial expressions while preserving identity.

### 👤 Female Subjects

<p>Text Identity</p> <p>Photo of a woman, with a <b>happy</b> expression, wearing a cozy sweater , sitting at a café shop.</p> 	<p>Text Identity</p> <p>Photo of a woman, with a <b>sad</b> expression, dressed in a raincoat and sitting on sofa.</p> 
<p>Text Identity</p> <p>Photo of a woman with an <b>angry</b> expression, wearing a T-shirt and sunglasses, sitting at a lakeside.</p> 	<p>Text Identity</p> <p>Photo of a woman, with an <b>angry</b> expression, wearing a formal office dress, standing in the train.</p> 
<p>Text Identity</p> <p>Photo of a woman with a <b>surprised</b> expression, standing in a vibrant outdoor marketplace.</p> 	<p>Text Identity</p> <p>Photo of a woman with a <b>shocked</b> expression, standing in a dimly lite room, with scattered papers.</p> 
<p>Reference</p>	<p>Reference</p>

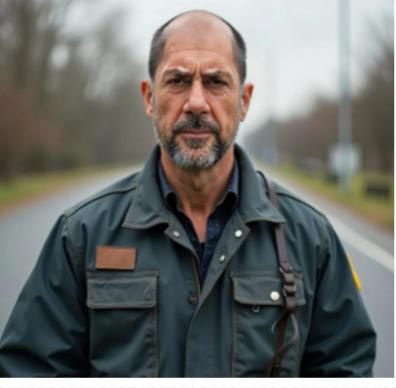
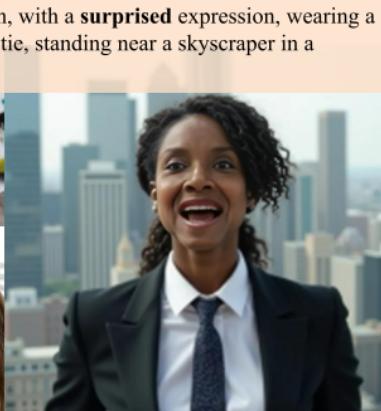
Expression synthesis for female subjects across five expression types with high visual fidelity.

## 👶 Children Across Ethnic Groups

<p>Text Identity</p>  	<p>Text Identity</p>  
<p>Text Identity</p>  	<p>Text Identity</p>  
<p>Text Identity</p>  	<p>Text Identity</p>  

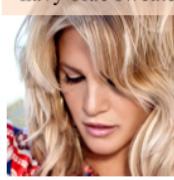
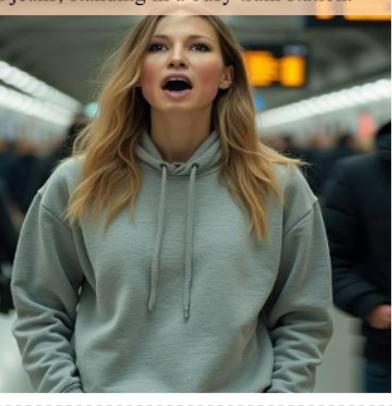
Results demonstrate effective generation for young boys and girls of Indian, African, and European descent.



	Text	Photo of a man, with a <b>happy</b> expression, , wearing sunglasses and a life jacket on a boat.				Text	Photo of a woman, with <b>sad</b> expression, wearing a gray T-shirt, sitting on the chair.		
	Reference	Identity				Reference	Identity		
	Text	Photo of a woman, with an <b>angry</b> expression, wearing a party hat, surrounded by birthday decorations and gifts.				Text	Photo of a man, with an <b>angry</b> expression, wearing a workshop suit, standing on road.		
	Reference	Identity				Reference	Identity		
	Text	Photo of a woman, with a <b>surprised</b> expression, wearing a business suit and tie, standing near a skyscraper in a bustling city.				Text	Photo of a man, with a <b>shocked</b> expression, wearing a hoodie and sitting on a park bench surrounded by fallen autumn leaves.		
	Reference	Identity				Reference	Identity		

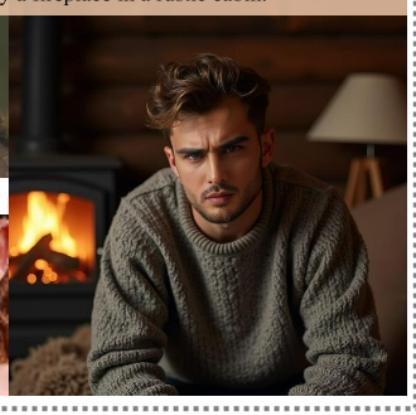
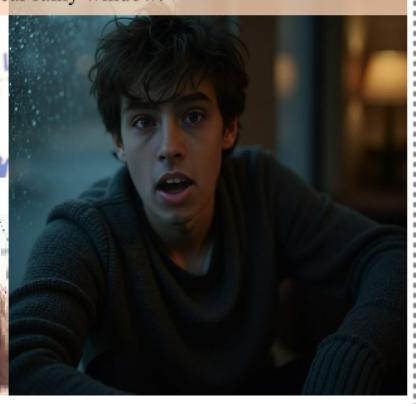
The model maintains identity and expression accuracy across various ethnic groups.

## Cross-Identity/Reference Combinations

<p>Text Identity</p>  	<p>Text Identity</p>  
<p>Text Identity</p>  	<p>Text Identity</p>  
<p>Text Identity</p>  	<p>Text Identity</p>  

Demonstrates flexibility in handling cross-gender and cross-ethnicity transformations while preserving identity.

## Extreme Orientations

<p>Text Identity</p> 	<p>Text Identity</p> 
<p>Text Identity</p> 	<p>Text Identity</p> 
<p>Text Identity</p> 	<p>Text Identity</p> 

Results show model limitations when both identity and reference inputs have extreme head poses, reflecting the boundaries of identity fidelity under such conditions.

## ✓ Conclusion

DeepExPo demonstrates robust performance in identity-preserving facial expression synthesis across diverse subjects, conditions, and contexts. Its ability to handle complex semantic cues and maintain realism positions it as a strong foundation for personalized human image generation in real-world applications.

## License

This project is licensed under the MIT License. See the [LICENSE](#) file for more details.

## Acknowledgments

- **Diffusion Models:** We used diffusion-based image generation techniques for creating high-quality images.
- **Open Source Tools:** The project leverages various open-source libraries such as PyTorch, Hugging Face Transformers, and others.

### Contact

For questions, collaboration, or code access, please contact:

**Muhammad Afgan**

 [msafgan@mail.ustc.edu.cn]