# MagicPose: Realistic Human Poses and Facial Expressions Retargeting with Identity-aware Diffusion

Di Chang[1,2]     Yichun Shi[2]     Quankai Gao[1]     Jessica Fu[1]     Hongyi Xu[2]

Guoxian Song[2]     Qing Yan[2]     Xiao Yang[2]     Mohammad Soleymani[1]

[1] University of Southern California     [2] ByteDance Inc.

https://boese0601.github.io/magicdance/

{dichang,quankaig,fujessic,msoleyma}@usc.edu

{yichun.shi,hongyi.xu,guoxian.song,qing.yan,xiao.yang}@bytedance.com

International Conference on Machine Learning (**ICML**) 2024

Presenter: Di Chang

USC Viterbi
School of Engineering

University of Southern California

# Outline

- **Introduction & Problem Definition**

- **Recap Diffusion Model & ControlNet**

- **Motivation**

- **Method - Pipeline**

- **Results & User Study**

- **Conclusion & Future Work**

# Outline

- **Introduction & Problem Definition**

- Recap Diffusion Model & ControlNet

- Motivation

- Method - Pipeline

- Results & User Study

- Conclusion & Future Work

# Introduction

- AI Generated Content (AIGC) using **Diffusion Models** is the most popular topic in the recent Computer Vision Research Community.

- Some representative works include: **Stable Diffusion**(Text2Image), **Sora**(Text2Video), **EMO**(Audio2Video).

- Image & Video Generation is the most important application of Diffusion Model.



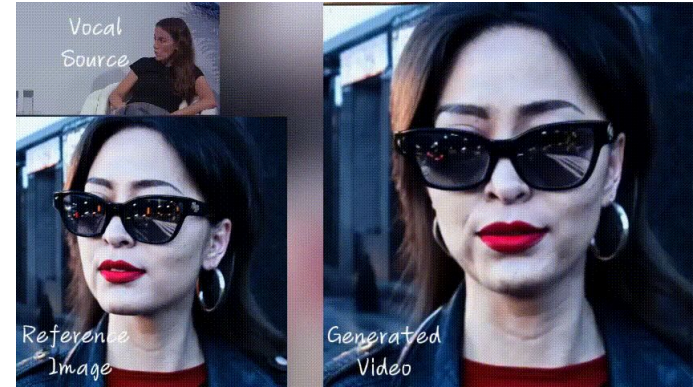A stylish woman wears sunglasses and red lipstick.

A stylish woman walks down a Tokyo street filled with warm glowing neon and animated city signage.

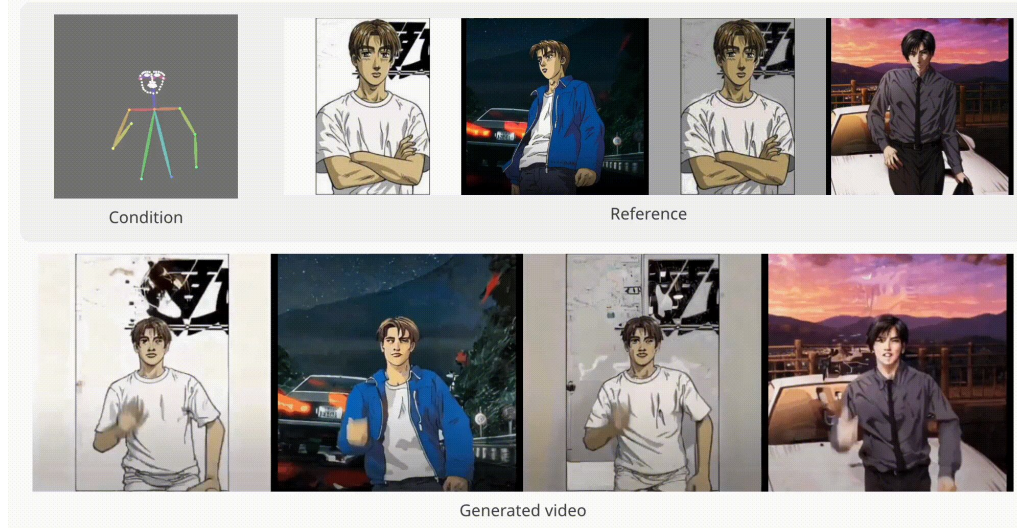**Stable Diffusion - LMU Munich**          **Sora - OpenAI**          **EMO - Alibaba**

# Problem Definition

- Human motion retargeting is the task of human image generation under the content control of a reference human subject **appearance** and the geometry control of **body motion and facial expression**.
- In this project, we propose **MagicPose**, an Image2Video model for motion retargeting with identity-aware diffusion.
- Such **Image2Image/Image2Video** technique can be applied to Virtual Reality, Creative Art Contents, Live Streaming, etc.
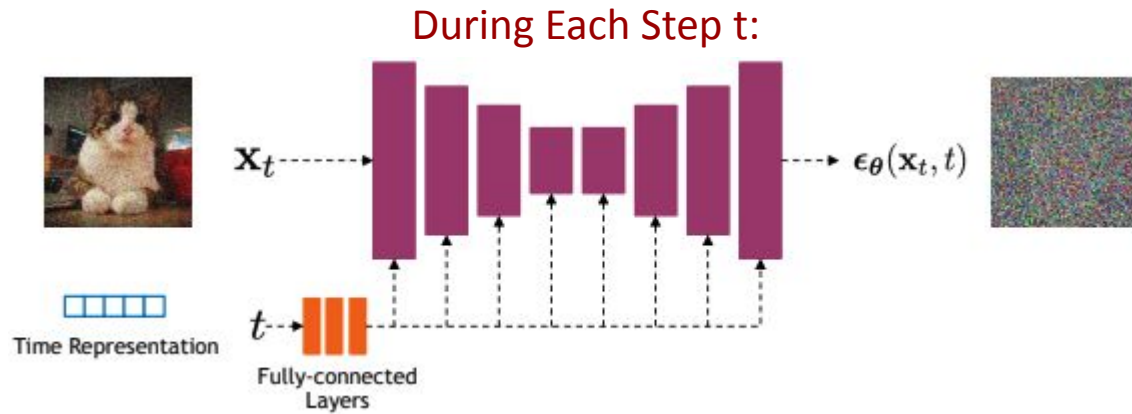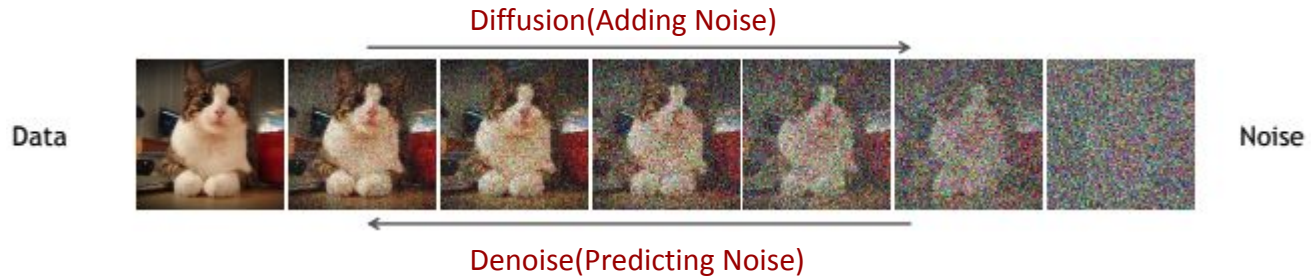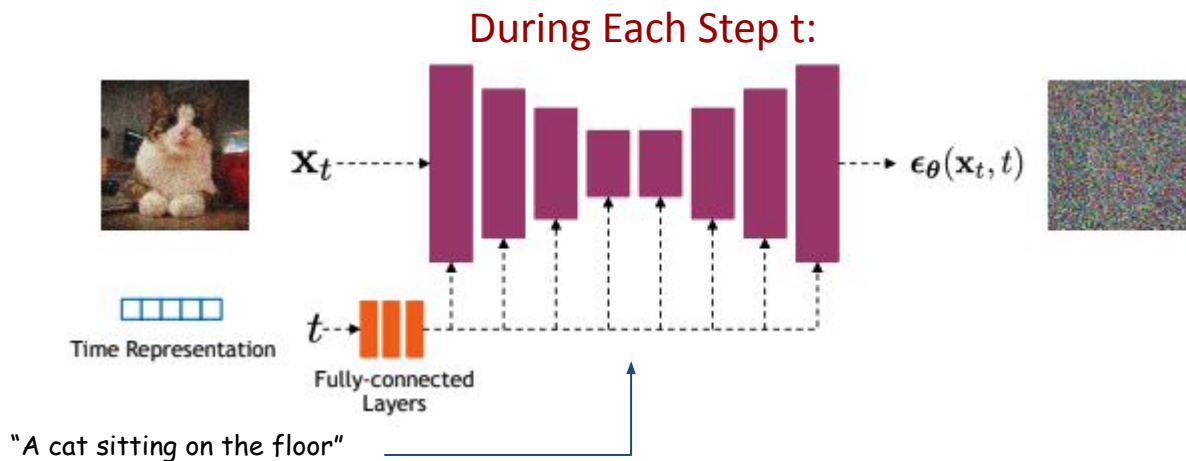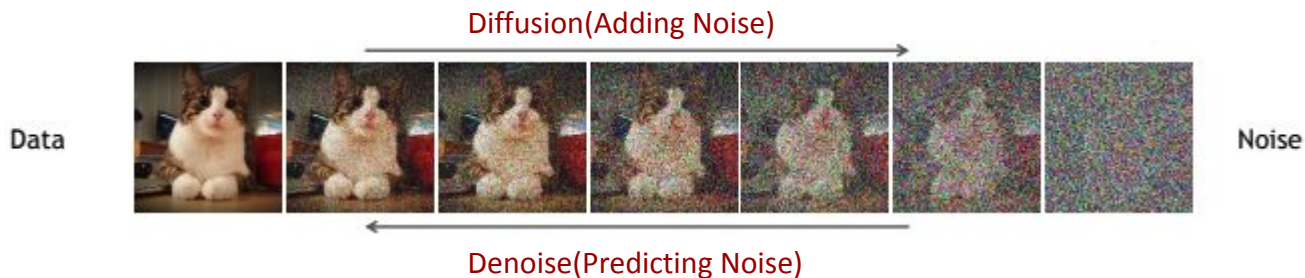


**MagicPose - USC**

# Outline

# Diffusion Model[1]

Diffusion(Adding Noise)

Data

Noise

Denoise(Predicting Noise)

## During Each Step t:

$\mathbf{x}_t$

$\epsilon_\theta(\mathbf{x}_t, t)$

Time Representation

$t$

Fully-connected Layers
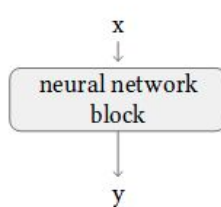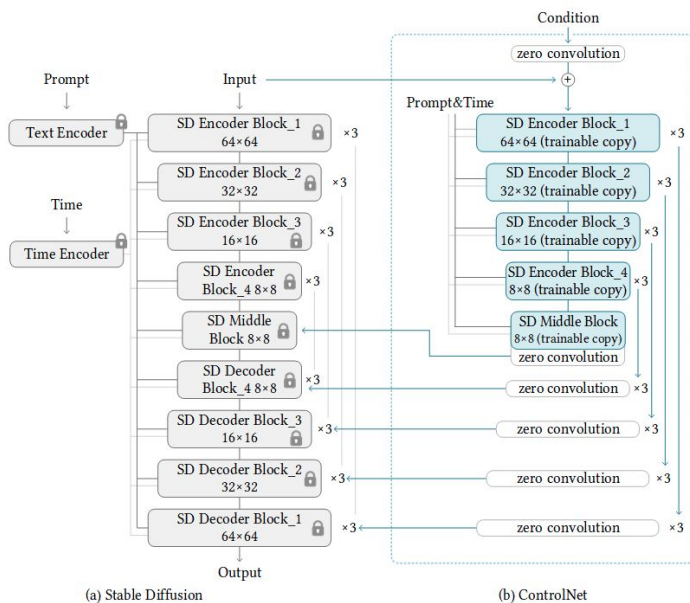
# Stable Diffusion[2] (Latent Diffusion Model-LDM)

# ControlNet[3]

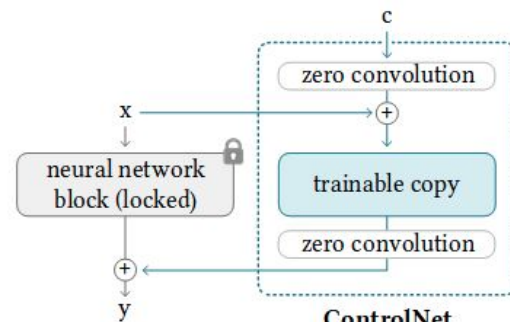- Text has already provided enough information for the content(appearance) of the object in the generation.
- Can we control the geometry of the object in the generation?
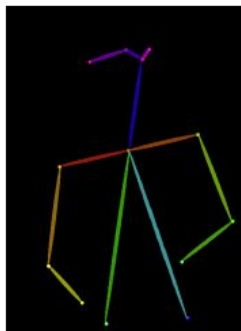- Can we freeze the weights of pretrained Stable Diffusion Model (keep the model safe)?

# Method - ControlNet[3]



Input human pose      Default      "chef in kitchen"      "Lincoln statue"

Input human pose      "man in suit"

# Outline

- Introduction & Problem Definition

- Recap Diffusion Model & ControlNet

- **Motivation**

- Method - Pipeline

- Results & User Study

- Conclusion & Future Work

# Motivation & Key challenges and limitations

- Existing works based on diffusion model cannot provide satisfactory identity/appearance preserving ability for real-human image generation.



Ground Truth          Pose          TPS          Disco          MagicPose

- Current works cannot generalize to unseen out-of-domain data after training on real-human data.

# Outline

- Introduction & Problem Definition

- Recap Diffusion Model & ControlNet

- Motivation

- **Method - Pipeline**

- Results & User Study

- Conclusion & Future Work

# Method - Exploration of Appearance Control Mechanism



Reference Image ControlNet Connected Attention Ours

ControlNet

# Method - Pipeline



Target Condition Map(s)

Noise

Reference

Pose ControlNet

Stable Diffusion UNet

Appearance Control Model

a) Appearance Control Pretraining

b) Appearance-Disentangled Pose Control

Multi-Source Self-Attention Module

$Q_1$ $K_1$ $K_2$ $V_1$ $V_2$

Self-attention

$K_2$ $V_2$ $Q_2$

Self-attention

: ResNetBlock

: TransformerBlock

: Zero Convolution

: Motion Module (optional)

: Appearance Control

: Pose Control

# Outline

- Introduction & Problem Definition

- Recap Diffusion Model & ControlNet

- Motivation

- Method - Pipeline

- **Results & User Study**
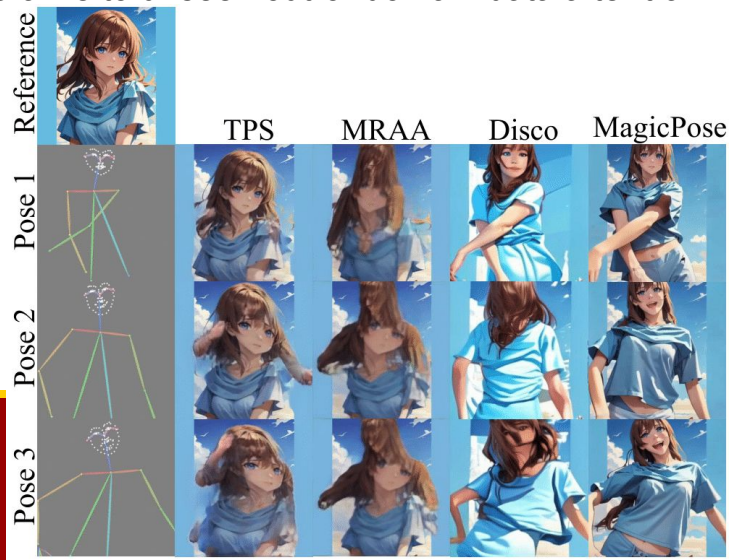
- Conclusion & Future Work

# Experiment - Dataset

TikTok Dataset

- Consists of 350 single-person **dance videos** (with video length of 10-15 seconds). Most of these videos contain the face and **upper-body** of a human.

EverybodyDanceNow Dataset

- Consists of full-body videos of five subjects. Experiments on this dataset aim to **test** our method's generalization ability to in-the-wild, **full-body out of domain motions**.

Self-collected Out-of-Domain Images

- Come from online resources. We use them to **test** our method's generalization ability to in-the-wild, **out of domain appearance.**

# Experiment - TikTok
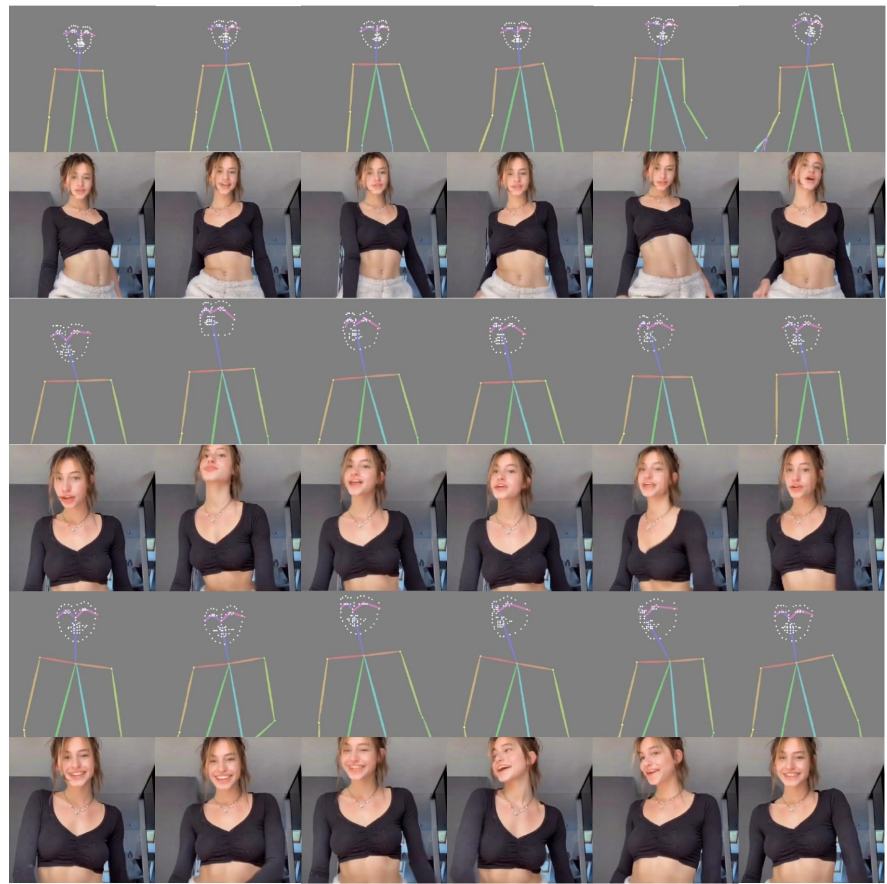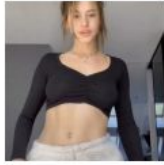


GT          Pose          TPS          Disco          MagicPose
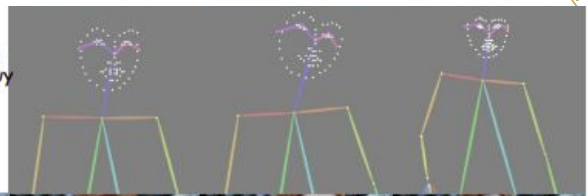
# Experiment - Out of Domain Appearance

# Comparison - Out of Domain Appearance



Reference      Pose      Animate Anyone      MagicPose

# Experiment - Out of Domain Motions

# Experiment - Quantitative Result

| Method | Image | | | | | | Video |
|---|---|---|---|---|---|---|---|
| | FID ↓ | SSIM ↑ | PSNR ↑ | LPIPS ↓ | L1 ↓ | Face-Cos ↑ | FID-VID ↓ |
| FOMM* (Siarohin et al., 2019a) | 85.03 | 0.648 | 29.01 | 0.335 | 3.61E-04 | 0.190 | 90.09 |
| MRAA* (Siarohin et al., 2021) | 54.47 | 0.672 | 29.39 | 0.296 | 3.21E-04 | 0.337 | 66.36 |
| TPS* (Zhao & Zhang, 2022) | 53.78 | 0.673 | 29.18 | 0.299 | 3.23E-04 | 0.280 | 72.55 |
| DisCo (Wang et al., 2023) | 50.68 | 0.648 | 28.81 | 0.309 | 4.27E-04 | - | 69.68 |
| DisCo† (Wang et al., 2023) | 30.75 | 0.668 | 29.03 | **0.292** | 3.78E-04 | 0.166 | 59.90 |
| MagicPose | **25.50** | **0.752** | **29.53** | **0.292** | **0.81E-04** | **0.426** | **46.30** |

**Note\***: **Face-Cos** is a novel metric which represents the cosine similarity of the extracted feature by AdaFace (Kim et al., 2022) of face area between generation and ground truth image.
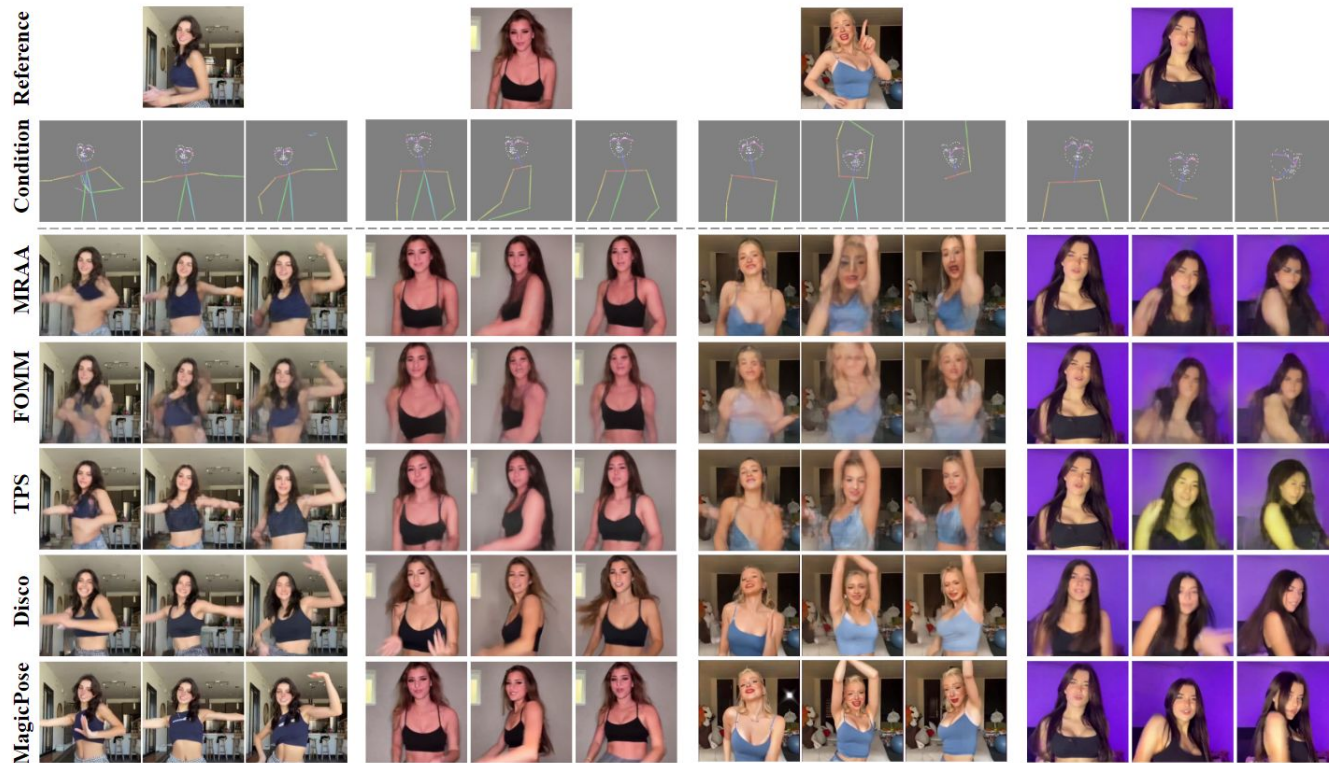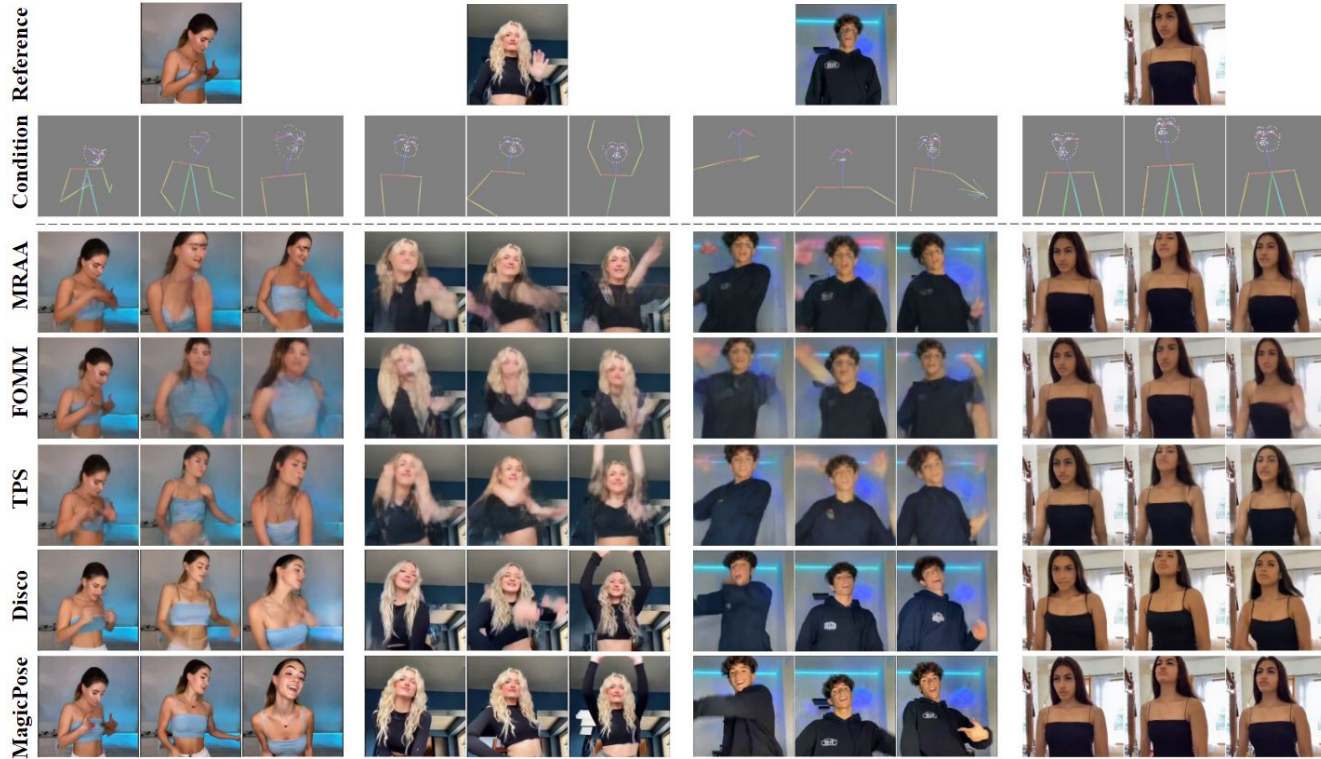
# User Study: Survey

- Participants: We use **Prolific**, an online platform designed to connect researchers with participants for academic studies and market research for all our user studies. The participants are English-speaking lay persons verified by this platform **without** any specific expertise in computer vision. We recruited **100** participants for this user study, and the hourly rate for this study is 72/hr. (~3 min per study)

- Data: We visualize 8 video sequences from the TikTok dataset and compare the performance of MagicPose to prior works. Examples are shown in the next slides.

- Procedure:
  a. For each of the 8 video sequences, we visualize different human poses and facial expressions.
  b. The methods are anonymized as A, B, C, D, E, and the order of the generated image from the corresponding method is randomized.
  c. We ask the Participants to choose **only one** best method.

- Tutorial & Example: See [here](here).

- Criteria for Judgment:
  a. 1) The appearance (Face, Clothes on the body, Background) of the generation should strictly match the given reference image input.
  b. 2) The motions and facial expressions of the generation should strictly match the given pose condition map input.

# User Study

# User Study

# User Study - Initial Votes from Users

# User Study - Table of Votes

Table 1. A large-scale user study with 100 participants. We collect the number of votes for eight video subjects from test set by five methods and report the percentage. Our MagicPose preserves the best identity information in pose and facial expression retargeting on all subjects.

| Method | Subject1 | Subject2 | Subject3 | Subject4 | Subject5 | Subject6 | Subject7 | Subject8 | Average |
|---|---|---|---|---|---|---|---|---|---|
| MRAA (Siarohin et al., 2019) | 8% | 6% | 0% | 5% | 2% | 2% | 8% | 4% | 4% |
| FOMO (Siarohin et al., 2021) | 3% | 1% | 3% | 1% | 1% | 0% | 5% | 8% | 3% |
| TPS (Zhao & Zhang, 2022) | 4% | 16% | 0% | 4% | 2% | 3% | 4% | 2% | 4% |
| Disco (Wang et al., 2023) | 12% | 3% | 9% | 18% | 5% | 20% | 33% | 27% | 16% |
| MagicPose | **73%** | **74%** | **88%** | **72%** | **90%** | **75%** | **50%** | **59%** | **73%** |

# Statistical Analysis - Chi Square Test

1. State the Null Hypothesis: There is no association between the video subjects and the choice of method. The distribution of votes for each method is the same across all groups, meaning any observed differences are due to chance.

2. Chi-square statistic: 116.02. p-value: $1.11 \times 10^{-12}$. Degrees of freedom: 28.

3. Conclusion and Discussion: Given the extremely small p-value (much less than 0.05), we can reject the Null Hypothesis. The differences in vote distribution are unlikely to have occurred by chance. The participants indeed prefer MagicPose more than other methods.

# Outline

- Introduction & Problem Definition

- Recap Diffusion Model & ControlNet

- Motivation

- Method - Pipeline

- Results & User Study
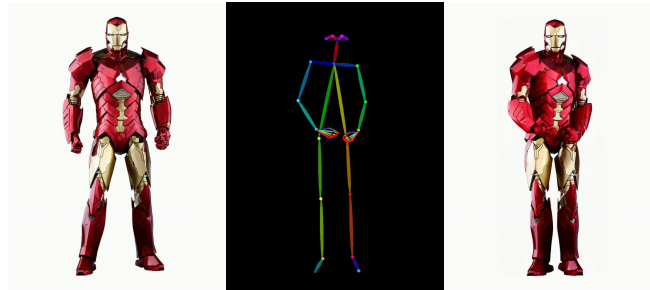
- **Conclusion & Future Work**

# Conclusion

- We propose an effective method - **MagicPose**, for human pose and expression retargeting.

- The proposed model can be used as a **plug-in** extension for Stable Diffusion.

- We introduce **Multi-Source Attention Module** that offers detailed appearance guidance.

- Experiment on out-of-domain data demonstrating strong **generalizability** of our model to diverse image styles and human motions.
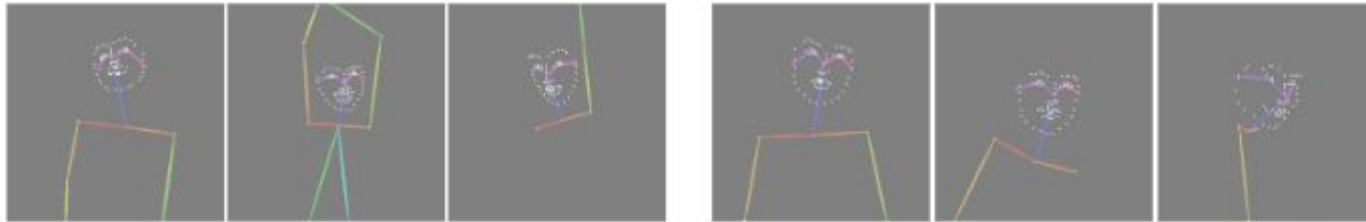
# Future Work

- Improve temporal consistency for video quality.



- Adopt more advanced pose detector for better motion representation. Current openpose detector is not stable and fails to detect complete skeleton.

# Social Impact

- Positive - Improving communication in digital or virtual environments
- Positive - Enhancing interactions in virtual meetings, online classrooms, and social networking platform
- Positive - Entertainment and media production, allowing for the creation of more lifelike and expressive characters in movies, video games, and animations


- Negative - Making fake animated videos of people which could be used in frauds
- Potential Solution - Digital watermarking and detection algorithms, enact and enforce stringent legal measures. Enhance public awareness and education on media literacy. Establish ethical guidelines within the tech industry.

# Reference

[1] Dhariwal, Prafulla, et al. "Diffusion Models Beat GANs on Image Synthesis." *Arxiv Report*. 2021

[2] Rombach, Robin, et al. "High-Resolution Image Synthesis with Latent Diffusion Models." *Proceedings of the IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)*. 2022

[3] Zhang, Lvmin, et al. "Adding Conditional Control to Text-to-Image Diffusion Models." *Proceedings of the IEEE/CVF Conference on International Conference on Computer Vision*. 2023

[4] Wang, Tan, et al. "DisCo: Disentangled Control for Realistic Human Dance Generation." *Proceedings of the IEEE / CVF Computer Vision and Pattern Recognition Conference (CVPR)*. 2024

# Thanks for listening!

Please scan the following QR-Code to check our project and more demos :))

Website:

Code: