

# 3D human interaction synthesis for action recognition data augmentation

Master Thesis







### **3D human interaction synthesis for action recognition data augmentation**

Master Thesis

June, 2024

By

Anders Bredgaard Thuesen

Copyright:      Reproduction of this publication in whole or in part must include the customary bibliographic citation, including author attribution, report title, etc.

Cover photo:    Vibeke Hempler, 2012

Published by:   DTU, Department of Applied Mathematics and Computer Science,  
Richard Petersens Plads, Building 324, 2800 Kgs. Lyngby Denmark  
[www.compute.dtu.dk](http://www.compute.dtu.dk)

ISSN:            [0000-0000] (electronic version)

ISBN:            [000-00-0000-000-0] (electronic version)

ISSN:            [0000-0000] (printed version)

ISBN:            [000-00-0000-000-0] (printed version)

## Approval

This thesis has been prepared over six months at the Section for Indoor Climate, Department of Civil Engineering, at the Technical University of Denmark, DTU, in partial fulfilment for the degree Master of Science in Engineering, MSc Eng.

It is assumed that the reader has a basic knowledge in the areas of statistics.

Anders Bredgaard Thuesen - s183926

.....  
*Signature*

.....  
*Date*

## **Abstract**

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

## Acknowledgements

**Anders Bredgaard Thuesen**, MSc Civil Engineering, DTU  
Creator of this thesis template.

**[Name]**, [Title], [affiliation]  
[text]

**[Name]**, [Title], [affiliation]  
[text]

# Contents

Preface . . . . .	ii
Abstract . . . . .	iii
Acknowledgements . . . . .	iv
<b>1 Introduction</b>	<b>1</b>
<b>2 Background</b>	<b>3</b>
2.1 Skinned Multi-Person Linear Model (SMPL) . . . . .	3
2.2 Human Mesh Reconstruction (HMR) . . . . .	3
2.3 Denoising Diffusion Probabilistic Models (DDPM) . . . . .	3
2.4 Classifier-free Guidance (CFG) . . . . .	5
2.5 Tranformer architecture . . . . .	5
2.6 Human Motion Generation . . . . .	5
<b>3 Data</b>	<b>7</b>
3.1 HumanML3D . . . . .	7
3.2 InterHuman . . . . .	7
3.3 Teton dataset . . . . .	7
<b>4 Methods</b>	<b>9</b>
4.1 Tracking & Matching . . . . .	9
4.2 Human pose sequence optimization . . . . .	9
4.3 3D scene reconstruction . . . . .	9
4.4 Motion and scene representation . . . . .	10
4.5 Diffusion model . . . . .	10
<b>5 Results</b>	<b>11</b>
<b>6 Discussion &amp; Conclusion</b>	<b>13</b>
<b>A Appendices</b>	<b>15</b>
A.1 Derivation of posterior distribution . . . . .	15





# 1 Introduction

In healthcare environments such as hospitals and care homes, data pertaining to critical incidents like falls are scarce due to their infrequent nature and the high privacy requirements surrounding such data. This scarcity poses significant challenges for training robust machine learning models, particularly in applications related to video classification where detailed understanding of such events is crucial. The primary goal of this thesis is to enhance the performance of video classification tasks by leveraging synthetic data, which is designed to closely mirror the underlying distribution of real-world incidents while enabling focused studies on specific, rare events.

To address the challenges inherent in collecting and utilizing real-world data from sensitive environments, this research proposes a novel approach using synthetic data generation. Synthetic data not only adheres to the distributional characteristics of genuine data but also provides flexibility to explore less common scenarios—specifically, those at the tail of the distribution which are typically underrepresented in available datasets.

This work introduces a sophisticated framework for generating synthetic data by explicitly modeling three-dimensional (3D) environments. This includes detailed interactions both among humans and between humans and their surroundings. By integrating these complex interactions as a strong inductive bias, the proposed generative diffusion model enhances the realism and applicability of the synthetic data.

To construct and train this model, we utilize publicly available datasets such as HumanML3D and InterHuman. These datasets include motion capture data of individuals and pairs interacting, each accompanied by textual descriptions. These are combined with 3D scene reconstructions derived from video captures in actual hospital and care home settings. This integration of human motion and scene specifics forms the foundation for our synthetic data generation process.

To effectively reconstruct 3D scenes from the captured videos, we employ state-of-the-art models such as ProHMR and Depth Anything. These models are instrumental in generating per-frame human pose estimations and scene depth labels. These outputs, along with 2D keypoint annotations, are fed into a joint optimization process. This process is critical as it unifies the coordinate systems of the human models and the environment, ensuring that the motion trajectories are smooth and coherent. The result is a highly accurate 3D representation of the scenes, which serves as a vital input for our synthetic data generation.



## 2 Background

### 2.1 Skinned Multi-Person Linear Model (SMPL)

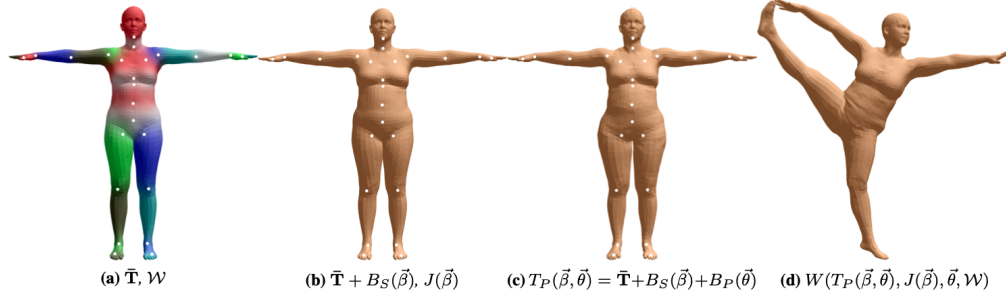


Figure 2.1: SMPL model

The Skinned Multi-Person Linear (SMPL) model is a parametric body shape model that accurately represents a wide range of human bodies and poses. It is built upon a foundation of linear blend skinning enhanced with corrective blend shapes, which are derived from a large dataset of body scans. The model captures the subtle deformations that occur with different body shapes and poses and can easily be rendered due to its compatibility with existing graphics pipelines. Since its publication, several extensions such as DMPL, incorporating dynamic soft-tissue deformation and SMPL-X, also modelling hands and facial expressions have been introduced. The model is parameterized by  $\vec{\beta}$ , capturing the variations from a mean body shape and  $\vec{\theta}$ , specifying the axis-angle rotation of 23 of the template skeleton joints. Mathematically, the model can be expressed as:

$$M(\vec{\beta}, \vec{\theta}) = W(T_P(\vec{\beta}, \vec{\theta}), J(\vec{\beta}), \vec{\theta}, \mathcal{W}) \quad (2.1)$$

where  $T_P(\vec{\beta}, \vec{\theta})$  returns the vertices of the rest pose, incorporating the deformations from the body shape and pose and is given by:

$$T_P(\vec{\beta}, \vec{\theta}) = \mathbf{T} + B_S(\vec{\beta}) + B_P(\vec{\theta}) \quad (2.2)$$

$J(\vec{\beta})$  returns the 3D joint locations from the shaped template vertices using a learned regression matrix  $\mathcal{J}$  and is given by:

$$J(\vec{\beta}) = \mathcal{J}(\mathbf{T} + B_S(\vec{\beta})) \quad (2.3)$$

$W$  is the skinning function (e.g. Linear Blend Skinning (LBS) or Dual-Quaternion Blend Skinning (DQBS)) and  $\mathcal{W}$  is the blend weights.

### 2.2 Human Mesh Reconstruction (HMR)

Recovering the mesh of humans have several applications...

### 2.3 Denoising Diffusion Probabilistic Models (DDPM)

In recent years, several types of generative models such as Variational Autoencoders (VAEs), Generative Adversarial Networks (GANs), autoregressive models and flow-based models have shown remarkable results in data generation of varying data modalities, such

as images, audio, videos and text. Most recently, Denoising Diffusion Probabilistic Models (DDPMs) have gained large popularity especially within the field of image generation due to several reasons such as high-quality data generation, versatility in several data domains as well as controllability, allowing one to steer the generation towards desired outputs. A DDPM is a parameterized ( $\theta$ ) Markov chain trained using variational inference to reverse a (forward) diffusion process,  $q(\mathbf{x}_{1:T}|\mathbf{x}_0)$  wherein the signal of the data,  $\mathbf{x}_0$ , is gradually destroyed by adding gaussian noise according to predefined noise schedule  $\{\beta_t \in (0, 1)\}_{t=1}^T$  giving rise to increasingly noisy samples,  $\mathbf{x}_1 \dots \mathbf{x}_T$ :

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(x_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}), \quad q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1}) \quad (2.4)$$

with  $T$  being the discretized number of diffusion steps before all original information is completely discarded. The goal of the inverse or backwards process then becomes to iteratively remove the noise, in order to arrive at the original data. More formally, the process is defined as:

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|x_t), \quad p_\theta(\mathbf{x}_{t-1}|x_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \sigma_\theta(\mathbf{x}_t, t)) \quad (2.5)$$

taking starting point in pure noise  $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$ , incrementally removing the noise through the learned functions,  $\mu_\theta(\mathbf{x}_t, t)$  and  $\sigma_\theta(\mathbf{x}_t, t)$  commonly parameterized by a deep neural network. Using the reparameterization trick, we are able to sample any noisy version of our data,  $\mathbf{x}_t$ , at time step  $t$  given our original data  $\mathbf{x}_0$ . Recall our forward transition probability function,  $q(\mathbf{x}_t|\mathbf{x}_{t-1})$ . Letting  $\alpha_t = 1 - \beta_t$ ,  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$  and using the reparameterization trick the expression can be rewritten as:

$$\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_{t-1} + \sqrt{1 - \alpha_t}\epsilon_{t-1} \quad (2.6)$$

where  $\epsilon_{t-1} \sim \mathcal{N}(0, 1)$ . Expanding the recursive definition then gives:

$$\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_{t-1} + \sqrt{1 - \alpha_t}\epsilon_{t-1} \quad (2.7)$$

$$= \sqrt{\alpha_t} \left( \sqrt{\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}}\epsilon_{t-2} \right) + \sqrt{1 - \alpha_t}\epsilon_{t-1} \quad (2.8)$$

$$= \sqrt{\alpha_t\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{\alpha_t(1 - \alpha_{t-1})}\epsilon_{t-2} + \sqrt{1 - \alpha_t}\epsilon_{t-1} \quad (2.9)$$

$$= \sqrt{\alpha_t\alpha_{t-1}}\mathbf{x}_{t-2} + \sqrt{1 - \alpha_t\alpha_{t-1}}\bar{\epsilon}_{t-2} \quad (2.10)$$

where  $\bar{\epsilon}_{t-2}$  merges the two independent Gaussians  $\epsilon_{t-1}$  and  $\epsilon_{t-2}$  into a single Gaussian with new variance as the sum of variances  $\alpha_t(1 - \alpha_{t-1}) + (1 - \alpha_t) = 1 - \alpha_t\alpha_{t-1}$ . Recursively applying the definition of  $\mathbf{x}_t$  and merging the gaussian noise terms results in the simplified expression:

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (2.11)$$

This result enables sampling of any noisy version of  $\mathbf{x}_t$  at time  $t$  given "clean" data  $\mathbf{x}_0$ . Conversely, given  $\mathbf{x}_0$  the reverse probability conditional probability

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\mu}_t, \tilde{\beta}_t)$$

becomes tractable to compute using Bayes rules with mean and variance given by (derivation in appendix A.1):

$$\tilde{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon \right) \quad (2.12)$$

$$\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \cdot \beta_t \quad (2.13)$$

The model is trained by minimizing the Evidence Lower Bound (ELBO) on the negative likelihood:

$$\mathcal{L}_{VLB} = \mathbb{E}_q \left[ -\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right] = \mathbb{E}_q \left[ -\log p(\mathbf{x}_T) - \sum_{t \geq 1} \log \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \right] \leq \mathbb{E} [-\log p_\theta(\mathbf{x}_0)] \quad (2.14)$$

however

## 2.4 Classifier-free Guidance (CFG)

## 2.5 Tranformer architecture

## 2.6 Human Motion Generation

Generating realistic human motion has several applications

- Human Motion Diffusion -



## 3 Data

### 3.1 HumanML3D

The HumanML3D dataset combines the HumanAct12 and AMASS datasets, integrating human motion captured using advanced motion capturing systems and converting the data to a unified parameterization. It covers a broad range of daily human activities, providing 14,616 motions in total, accompanied by 44,970 single-sentence descriptions. Each motion clip includes 3-4 descriptions, and the entire dataset amounts to 28.59 hours of recorded motion. Additionally, the data is augmented by mirroring all motions, with corresponding adjustments to descriptions, such as changing “clockwise” to “counterclockwise.”

### 3.2 InterHuman

The InterHuman dataset is a comprehensive, large-scale 3D dataset designed to capture human interactive motions involving two individuals. It includes approximately 107 million frames detailing a wide range of human interactions, from professional activities to daily social behaviors. Each motion sequence is paired with natural language annotations, totaling 23,337 descriptions, which provide context and detail for the captured interactions, enhancing the dataset’s utility for training and evaluating models.

### 3.3 Teton dataset

10s of video with 10 frames per second. Pseudo dataset of SMPL pose predictions from HMR2.0 for each frame





## 4 Methods

### 4.1 Tracking & Matching

We track the predicted staff and patients across multiple frames, indicated by  $t = 1 \dots T$ , by iteratively assigning the predictions,  $\{P_i^{(t)}\}_{i=1}^{M_t}$ , to the latest known tracks,  $\{Q_j^{(t)}\}_{j=1}^{N_t}$  by greedily picking the assignment with the minimum cost:

$$\underset{i,j}{\operatorname{argmin}} \mathcal{L}_{\text{track}}(P_i^{(t)}, Q_j^{(t-1)}), \quad (4.1)$$

until either all predictions or latest tracks have been assigned exactly once. In the case of any unassigned predictions, a new track is initialized. After tracks have been assigned, the latest track  $Q_j^{(t)}$  is updated with the assigned predictions. We choose the following composite loss function:

$$\mathcal{L}_{\text{track}}(P, Q) = \alpha \|P_{\text{3D kpts}} - Q_{\text{3D kpts}}\|_2 + \beta \|P_{\text{class}} - Q_{\text{class}}\|_\infty, \quad (4.2)$$

incorporating both the Euclidean distance between the 3D joint keypoints as well as the predicted person class (staff/patient), with  $\alpha$  and  $\beta$  weighting the influence of each term. The tracks,  $T_i$ , are then greedily matched to the ground truth annotations  $G_j$ , minimizing the loss:

$$\mathcal{L}_{\text{match}}(T, G) = \sum_t \begin{cases} \|a_{\text{track}, 2\text{D kpts}}^{(t)} - b_{\text{track}, 2\text{D kpts}}^{(t)}\|_2 & \text{if } a_{\text{track}}^{(t)} \in a_{\text{track}} \\ \gamma & \text{otherwise} \end{cases} \quad (4.3)$$

over the trajectory time horizon  $t = 1, 2, \dots, T$  where  $\gamma$  is the punishment for not detecting the person.

### 4.2 Human pose sequence optimization

### 4.3 3D scene reconstruction

#### Restoring depth scale and offsets

We utilize pseudo ground truth disparity maps generated by the Depth Anything model, inversely proportional to the scene depth, in order to reconstruct a point cloud of the scene. To recover the depth scale and offset we rasterize the predicted SMPL poses, extract the z-buffer, compute the inverse depth and regress the intersection with the normalized disparity maps using the Random Sample Consensus (RANSAC) algorithm robust to outliers.

$$\begin{bmatrix} px \\ y \\ 1/z \\ 1 \end{bmatrix} = \begin{bmatrix} f_x & 0 & 0 & p_x \\ 0 & f_y & 0 & p_y \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} \quad (4.4)$$

$$X = \frac{Z}{f_x} (x - p_x), \quad Y = \frac{Z}{f_y} (y - p_y) \quad (4.5)$$

where  $f_x, f_y$  are the horizontal and vertical focal lengths respectively and  $(p_x, p_y)$  is the principal point often chosen as the center  $(w/2, h/2)$  of the image.

#### **4.4 Motion and scene representation**

Previous work on single human motion generation uses a canonical representation

We use a 6D continuous representation of the joint angles as according to **Zhou\_2019\_CVPR**.

#### **4.5 Diffusion model**

## 5 Results



## **6 Discussion & Conclusion**





# A Appendices

## A.1 Derivation of posterior distribution

$$\begin{aligned}
 q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) &= q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} \\
 \Leftrightarrow \log q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) &= \log q(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) + \log q(\mathbf{x}_{t-1}|\mathbf{x}_0) - \log q(\mathbf{x}_t|\mathbf{x}_0) \\
 &= \frac{1}{2} \left( \frac{(\mathbf{x}_t - \sqrt{\alpha_t} \mathbf{x}_{t-1})^2}{\beta_t} + \frac{(\mathbf{x}_{t-1} - \sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0)^2}{1 - \bar{\alpha}_{t-1}} + \frac{(\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0)^2}{1 - \bar{\alpha}_t} \right) + K \\
 &= \frac{1}{2} \left( \frac{\mathbf{x}_t^2 - 2\sqrt{\alpha_t} \mathbf{x}_t \mathbf{x}_{t-1} + \alpha_t \mathbf{x}_{t-1}^2}{\beta_t} + \frac{\mathbf{x}_{t-1}^2 - 2\sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_{t-1} \mathbf{x}_0 + \bar{\alpha}_{t-1} \mathbf{x}_0^2}{1 - \bar{\alpha}_{t-1}} + C(x_t, x_0) \right) + K
 \end{aligned}$$

We collect all terms for  $\mathbf{x}_{t-1}$  and  $\mathbf{x}_{t-1}^2$ .

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Technical  
University of  
Denmark

Richard Petersens Plads, Building 324  
2800 Kgs. Lyngby  
Tlf. 4525 1700

[www.compute.dtu.dk](http://www.compute.dtu.dk)