

HUMOS: Human Motion Model Conditioned on Body Shape

Shashank Tripathi^{1,2*} Omid Taheri¹ Christoph Lassner² Michael Black¹
Daniel Holden² Carsten Stoll²

¹Max Planck Institute for Intelligent Systems, Tübingen, Germany ²Epic Games
{stripathi, otaheri, black}@tue.mpg.de
{christoph.lassner, daniel.holden, carsten.stoll}@epicgames.com

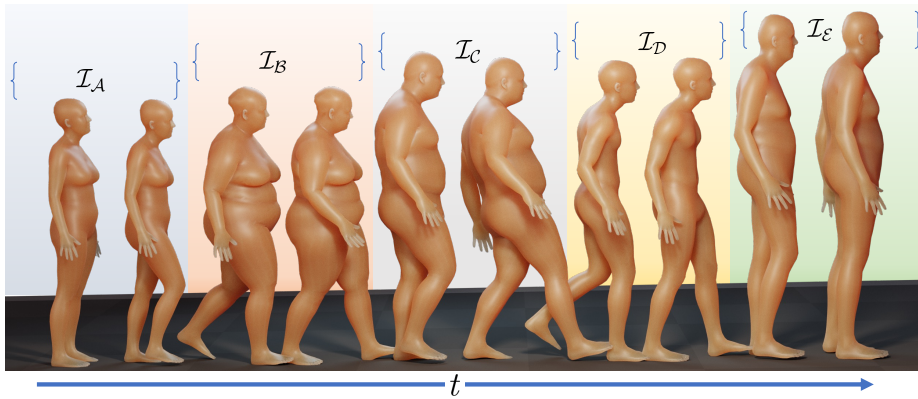


Fig. 1: People with different body shapes perform the same motion differently. Our method, HUMOS, generates *natural*, *physically plausible* and *dynamically stable* human motions conditioned on body shape. HUMOS uses a novel identity-preserving cycle consistency loss and differentiable dynamic stability and physics terms to learn an identity-conditioned manifold of human motions. Shown here is the same walk motion with a skip-step in the middle, generated by HUMOS for five different identities $\mathcal{I}_{A:\varepsilon}$. To demonstrate shape-conditioning, we visualize the same motion but successively change the identity after every 30 frames.

Abstract. Generating realistic human motion is crucial for many computer vision and graphics applications. The rich diversity of human body shapes and sizes significantly influences how people move. However, existing motion models typically overlook these differences, using a normalized, average body instead. This results in a homogenization of motion across human bodies, with motions not aligning with their physical attributes, thus limiting diversity. To address this, we propose a novel approach to learn a generative motion model conditioned on body shape. We demonstrate that it is possible to learn such a model from unpaired training data using cycle consistency, intuitive physics, and stability constraints that model the correlation between identity and movement. The

* work done during an internship at Epic Games

resulting model produces diverse, physically plausible, and dynamically stable human motions that are quantitatively and qualitatively more realistic than existing state of the art. More details are available on our project page <https://github.com/CarstenEpic/humos>.

1 Introduction

Modeling virtual humans that move and interact realistically with 3D environments is extremely important for interactive entertainment, AR/VR and simulation technology, with numerous applications in crowd simulation, gaming and robotics. There has been rapid progress in training models that generate human motion either unconditionally or conditioned on text or previous motions. Existing state-of-the-art human motion models [18, 55, 56, 69] are trained on datasets like AMASS [49], but they typically ignore body shape and proportions. However, variations in muscle mass distribution and body proportions contribute to a person’s distinct movement patterns. People with different body types will generally move differently when prompted to perform the same motion. We argue that to achieve physical realism and motion diversity, it is critical to condition generated human motions on body shape.

To address this problem we adopt a novel approach called HUMOS, that enhances traditional data-driven motion generation methods and uses a transformer-based conditional Variational Auto-Encoder (c-VAE) trained to generate human motion conditioned on *identity* features such as a subject’s body shape and sex. We take inspiration from a recent 3D human pose and shape estimation method, IPMAN [70], to propose new dynamic intuitive physics (IP) terms that are simple, fully differentiable, and compatible with parametric body models like SMPL [47]. Since IPMAN’s IP terms only apply to static 3D poses, they are not suited for dynamic human motion modeling. We go beyond this by proposing general IP terms that are effective for dynamic human motions involving a sequence of poses. We show that these dynamic IP terms are critical to effectively train our model without paired data of differently-shaped people performing the same action.

Specifically, we propose differentiable physics terms that improve the realism of generated motions by addressing common issues like foot sliding, ground penetration, and unrealistic floating effects. Our key contribution here is a dynamic stability term, that models the interaction between a body’s Center of Mass (CoM), Center of Pressure (CoP), and the Zero Moment Point (ZMP). Dynamic stability is a biomechanical concept, frequently employed for ensuring balance in humanoid robots [38], but has also been shown to hold true for human gait [57]. This approach ensures our generated motions are not only visually convincing but also more closely adhere to principles of biomechanics, making them suitable for a wide range of applications in realistic motion generation.

One of the key applications our model enables is retargeting of animation between characters with different identities. Existing methods typically generate human motions for a canonical body and then use a second character retargeting step to transfer the generated motions to the target body. Since classical

retargeting methods rely on simple heuristics and ignore body shape, they tend to fail for extreme poses and complex motions involving significant body-on-ground and self-contact. In contrast, HUMOS effectively learns how people with different body shapes and proportions perform the same motion; see Fig. 1.

Given an input motion for a particular identity, the HUMOS encoder outputs an encoding in an identity-agnostic latent space. The decoder receives this encoding along with a target identity and outputs a motion that resembles the input motion but as performed by the new identity. We leverage solutions from unpaired image-to-image translation literature [87] to design a self-supervised loss that leverages cyclic consistency within the encoder-decoder step. The cycle-consistent formulation results in realistic motions given a target identity.

We observe that the cycle-consistency alone is not enough for this task as the network may learn trivial solutions that ignore the target identity and output the same target motion as the source, while still satisfying the cycle consistency constraint. For example, merely copying the source motion to the target body will result in significant foot-sliding, ground penetration, floating, and dynamic instabilities. To prevent this, our key insight is to incorporate our IP and dynamic stability terms as training losses on the generated target body motions. This ensures the generated motions remain physically-plausible and dynamically stable, and encourages the network to use the conditioning body shape. This makes HUMOS the first data-driven human motion model that generates motions that are not only realistic but also physically plausible, dynamically stable, and tailored to the input body shape.

2 Related Work

We categorize related work into several broad areas: previous works which transfer motion between different skeleton proportions or topologies, works which use physics simulation as a prior to constrain human motion generation, and works which synthesize or generate motion, often conditioned on various different desired parameters.

Motion Transfer: Most industry methods of transferring motion between two characters assume that characters have either the same skeletal topology (but may differ in bone lengths) or a manual mapping between the two is provided. Motion is then largely transferred directly, without taking into account the further shape or identity difference between the characters. Simple heuristics such as inverse kinematics are used to remove any artifacts [20, 51, 62]. There have been several attempts to apply human motion data to entirely different characters or creatures [1, 76, 85], yet they tend to require paired data to function effectively which can be difficult to obtain at a large scale. Recent methods using Deep Learning have been developed which can transfer motion between different topologies [2, 43] without paired data, or even between entirely new mesh shapes [17]. Similarly, techniques have been developed which can re-target motion from other sources or data representations such as 2D videos [3], or can take into account physical constraints such as floor contacts [13, 72]. However, these

methods generally do not take into account the character’s shape or identity beyond their skeleton proportions. While some attempts have been made to build retargeting systems which take into account character identities [29, 58, 74, 83], these works are limited in scale, and only work on poses or very short windows of motion or on a small number of body types.

Physics-based Motion Modeling: Physics-based motion generation, particularly through the use of reinforcement learning (RL) within physics engines, has emerged as a prominent method for creating physically plausible humanoid motions. This approach, leveraging RL, navigates the complex solution space of human motion, aiming to produce motions that adhere to physical laws. Common approaches in this direction include the development of locomotion skills and user-controllable policies for character animation through deep RL [14, 50, 52–54, 64, 65, 75]. Despite the principled framework RL offers, it comes with its limitations. The extensive training required due to the high-dimensional space of human motions, the reliance on reward functions over data for motion generation, and the computational expense of physics simulators present significant challenges. Also, these engines are typically non-differentiable black boxes, making them incompatible with data-driven learning frameworks [24, 70]. To overcome these challenges, physics-based trajectory optimization and motion imitation have been applied for 3D human pose estimation, further highlighting the significance of physics in capturing human dynamics [14, 47, 66, 67, 77–79, 81, 82]. These simulators often utilize simplistic, non-differentiable models that fail to capture the intricacies of skin-level contact or muscle activations, leading to motions that, while physically plausible, lack the naturalness, diversity, and expressiveness found in data-driven approaches. We take inspiration from [70] and biomechanical physics terms, and propose to combine the physical plausibility with our data-driven method to generate more accurate and lifelike human motions while considering body shape and physics.

Data-Driven Human Motion Generation: Early efforts in human motion generation utilized deterministic models, producing single motion outcomes and failing to capture the stochastic nature of human motions [6, 11, 14, 23, 27, 30–33, 36]. The shift towards deep generative models like GANs and VAEs marked a significant advancement, enabling the modeling of human motions’ probabilistic nature conditioned on various inputs such as past motions [7, 12, 19, 21, 22, 34, 35, 68], music or speech [15, 28, 42, 44–46, 84], and text or action labels [4, 5, 16, 26, 32, 33, 55, 61]. The introduction of denoising diffusion models [10, 18, 69] represents a leap forward, merging the strengths of traditional generative models to achieve state-of-the-art performance in motion generation. Despite these advancements, a key challenge in shape conditioning remains: the lack of paired training data between body shape and motions. This leads most data-driven methods to forgo shape conditioning altogether. These approaches typically normalize training data, mainly from the AMASS [49] dataset, to a canonical skeleton or mean body shape [8, 9, 18, 55, 56, 60, 69], successfully learning the manifold of realistic motions for a canonical body, but often at the expense of physical and biomechanical realism. Such simplifications result in pronounced inconsistencies like foot sliding

and ground penetration, undermining the realism necessary for applications in animation and virtual reality [37, 59].

3 Method

Our goal is to learn an identity-conditioned motion model capable of generating 1) *realistic*, 2) *physically-plausible* and 3) *dynamically-stable* human motions. Specifically, we represent a 3D human motion as a sequence of poses $\mathbf{P}_{1:T} = \mathbf{P}_1, \dots, \mathbf{P}_T$, where T denotes the number of frames. We follow prior work [8, 9, 55, 56] and represent each pose \mathbf{P} by the 3D SMPL [47] vertex mesh, $\mathbf{V}(\theta, \beta, \mathcal{G})$. We choose a mesh-based representation as any physical analysis of human motion requires accurate modeling of skin-level surface contact. SMPL is a convenient choice as it parameterizes the 3D mesh into disentangled pose, θ , body shape, β and gender, \mathcal{G} parameters, allowing explicit and independent control over the gender, body shape and pose. For ease of notation, we combine the body shape and gender into a single *identity* parameter, $\mathcal{I} = (\beta, \mathcal{G})$, and use it as the conditioning signal in our motion model. Given the target identity, \mathcal{I}_t and an arbitrary duration T , we generate a sequence consisting of the global root joint position, $\mathbf{x}_t \in \mathbb{R}^3$ and the root-relative joint rotations, $\mathbf{r}_t \in \mathbb{R}^{J \times 6}$ in the 6D rotation representation [86], where $J = 23$ is the number of SMPL joints and one global rotation. Although SMPL uses parent-relative rotations defined on the kinematic chain, we empirically observe that using root-relative joint rotations results in more stable gradients and faster convergence. Consequently, we convert the SMPL parent-relative pose parameters, to root-relative rotations to construct our motion features. Additionally, we process all sequences by removing the z-component of the first-frame root orientation \mathbf{r}_1^z and the horizontal root translation, \mathbf{x}_1^x and \mathbf{x}_1^y . Doing this canonicalizes all sequences to start at the origin with the same forward facing direction and makes network training easier. Please refer to Sup. Mat. for a detailed description of our motion representation.

3.1 HUMOS model architecture

HUMOS is designed as a conditional Variation Auto-Encoder (c-VAE) network that generates sequential motion features in a non-autoregressive manner where we output motion features for T consecutive frame in one-shot. Our choice of non-autoregressive training is driven by the observation that while auto-regressive approaches are effective in generating simple motions like walking, running, *etc.*, one-shot approaches yield better motion diversity [60]. Consequently, most text-to-motion approaches employ non-autoregressive generative modeling since they focus on generating diverse motions conditioned on text. Following this trend, we use a non-autoregressive training paradigm to output diverse motions conditioned on body shape. We use a Transformer [71] architecture to obtain spatio-temporal embeddings from the input motion features. Next, we describe our motion encoder followed by the motion decoder (see Fig. 2 for an overview).

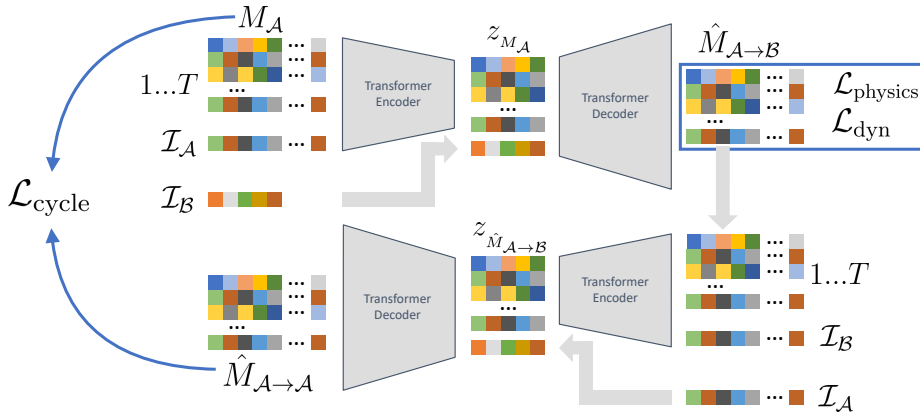


Fig. 2: Visual representation of our architecture. The *Encoder* takes as input a motion M_A and its associated identity \mathcal{I}_A , and outputs a latent (identity invariant) encoding of the motion z_{M_A} . The *Decoder* takes as input the latent encoding of the motion z_{M_A} , along with a different identity \mathcal{I}_B , and produces a retargeted motion appropriate for the given identity $\hat{M}_{A \rightarrow B}$. The same *Encoder* and *Decoder* are used with the original identity to produce a cycle loss $\mathcal{L}_{\text{cycle}}$, while a physics loss $\mathcal{L}_{\text{physics}}$ ensures the retargeted motion $\hat{M}_{A \rightarrow B}$ is realistic with respect to the given identity \mathcal{I}_B and prevents the cycle consistency loss from collapsing to a trivial solution.

Motion encoder: We extend the Transformer-based VAE motion encoder from TEMOS [56] by conditioning it on identity features, \mathcal{I} . Given a motion sequence M_A of arbitrary length, T , and the conditioning signal, \mathcal{I}_A of the source identity, our encoder \mathcal{E} outputs distribution tokens μ and Σ for the shape-conditioned motion latent space. Using the reparameterization trick [40], we sample the latent vector $z_{M_A} \in R^d$, embedding the input features into a d -dimensional latent space. To represent the temporal ordering in the input sequence, we use positional encodings in the form of sinusoidal functions and concatenate them with the input features [71].

Motion decoder: The decoder takes as input the latent code z_{M_A} and the target identity, \mathcal{I}_B and generates a sequence of 3D poses $\hat{\mathbf{P}}_{1:T}$ representing the source 3D motion as performed by an individual with the target body shape and gender. Our motion decoder \mathcal{D} is built on a Transformer architecture, which incorporates time information through T sinusoidal positional embeddings as queries. A concatenated combination of latent vectors and identity features serves as the key-value pair. Our decoder architecture mirrors the motion encoder, except for the input and output layers.

3.2 Self-supervised shape-conditioned training

Training a shape-conditioned c-VAE in a fully supervised manner requires paired identity and motion data. While datasets like AMASS [49] contain a large collection of 3D motion capture sequences, often containing diverse motions per-

formed by the same person, we seldom find different identities performing the same motion. Such pairs are essential for any motion model to effectively disambiguate motion from body type. This lack of *paired* data impedes training a shape-conditioned model. In fact, most existing methods capable of generating human motions ignore the body shape, generating motions only for the canonical skeleton and mean body shape using SMPL’s *neutral* gender body model.

To circumvent the lack of paired data in AMASS, we draw inspiration from the success in image-to-image translation [87] and use a novel self-supervised training strategy that leverages cycle-consistency in the shape-motion space. As shown in Fig. 2, we randomly sample a motion capture sequence for identity \mathcal{A} and extract motion features $M_{\mathcal{A}}$ comprising the global root-joint translation $\mathbf{x}_{\mathcal{A}}$ and the root-relative joint rotations $\mathbf{r}_{\mathcal{A}}$. We project the identity features, $\mathcal{I}_{\mathcal{A}} = (\beta_{\mathcal{A}}, \mathcal{G}_{\mathcal{A}})$ using a linear layer. The identity features are concatenated with the motion features and fed to the motion encoder which embeds them in a shape-agnostic latent code, $z_{M_{\mathcal{A}}}$. The HUMOS decoder takes the latent code $z_{M_{\mathcal{A}}} = \mathcal{E}(M_{\mathcal{A}}, \mathcal{I}_{\mathcal{A}})$ as input and a randomly sampled target identity, \mathcal{B} , with the projected $\mathcal{I}_{\mathcal{B}}$ as the new conditioning input to the decoder. The task of the decoder is to generate the root-joint translations $\hat{\mathbf{x}}_{\mathcal{A} \rightarrow \mathcal{B}}$ and the joint rotations $\hat{\mathbf{r}}_{\mathcal{A} \rightarrow \mathcal{B}}$ representing $\hat{M}_{\mathcal{A} \rightarrow \mathcal{B}} = \mathcal{D}(z_{M_{\mathcal{A}}}, \mathcal{I}_{\mathcal{B}})$, *i.e.* the motion $M_{\mathcal{A}}$ in the style of the identity of \mathcal{B} .

Since we lack any explicit ground truth to supervise $\hat{M}_{\mathcal{A} \rightarrow \mathcal{B}}$, as illustrated in Fig. 2, we employ cycle-consistency by reversing the forward step, this time using $\hat{M}_{\mathcal{A} \rightarrow \mathcal{B}}$ as the source motion and \mathcal{A} as the target identity. We input $\hat{M}_{\mathcal{A} \rightarrow \mathcal{B}}$ and $\mathcal{I}_{\mathcal{B}}$ to the HUMOS encoder and extract latent code, $z_{\hat{M}_{\mathcal{A} \rightarrow \mathcal{B}}}$. The latent code, along with projected $\mathcal{I}_{\mathcal{A}}$ are fed to the decoder resulting in $\hat{M}_{\mathcal{A} \rightarrow \mathcal{A}} = \mathcal{D}(z_{\hat{M}_{\mathcal{A} \rightarrow \mathcal{B}}}, \mathcal{I}_{\mathcal{A}})$. In the full cycle, since the motion style remains the same even as identities change, the same identity features should result in the same motion. Consequently, the reconstructed $\hat{M}_{\mathcal{A} \rightarrow \mathcal{A}}$ should match the source motion $M_{\mathcal{A}}$ and we define the cycle-consistency loss as $\mathcal{L}_{\text{cycle}} = \mathcal{L}_{\text{rot}} + \mathcal{L}_{\text{pos}}$ where \mathcal{L}_{rot} computes the geodesic distance in the rotational space by converting 6D joint rotations, $\mathbf{r}_{t_{\mathcal{A}}}$ and $\hat{\mathbf{r}}_{t_{\mathcal{A} \rightarrow \mathcal{A}}}$ to rotation matrices $R_{t_{\mathcal{A}}}$ and $\hat{R}_{t_{\mathcal{A} \rightarrow \mathcal{A}}}$ using the Gram-Schmidt process [86]. \mathcal{L}_{pos} is the smooth L1 loss between source and reconstructed root joint positions, $\mathbf{x}_{t_{\mathcal{A}}}$ and $\hat{\mathbf{x}}_{t_{\mathcal{A} \rightarrow \mathcal{A}}}$. Specifically,

$$\mathcal{L}_{\text{cycle}} = \mathcal{L}_{\text{rot}} + \mathcal{L}_{\text{pos}} \quad (1)$$

$$\mathcal{L}_{\text{rot}} = \sum_{t=1}^T \arccos \frac{\text{Tr}(R_{t_{\mathcal{A}}}(\hat{R}_{t_{\mathcal{A} \rightarrow \mathcal{A}}})^{-1}) - 1}{2}, \quad \mathcal{L}_{\text{pos}} = \sum_{t=1}^T \|\mathbf{x}_{t_{\mathcal{A}}} - \hat{\mathbf{x}}_{t_{\mathcal{A} \rightarrow \mathcal{A}}}\|_1. \quad (2)$$

3.3 Intuitive-physics (IP) terms

Our motion encoder aggregates spatio-temporal features over successive frames to learn a shape-agnostic latent embedding by disentangling the motion “style”

from identity-specific attributes. The decoder, in turn, leverages the shape-agnostic latent code and maps the motion style to a new target body. While intuitive, $\mathcal{L}_{\text{cycle}}$, however, is not enough as training with only $\mathcal{L}_{\text{cycle}}$ is prone to trivial solutions. Without special care, the encoder-decoder architecture could learn to generate identical motion $M_{\mathcal{A}} \approx \hat{M}_{\mathcal{A} \rightarrow \mathcal{B}} \approx \hat{M}_{\mathcal{A} \rightarrow \mathcal{A}}$ at intermediate steps, ignoring the identity conditioning while naïvely minimizing $\mathcal{L}_{\text{cycle}}$. To alleviate this, as shown in Fig. 2, we incorporate intuitive physics terms $\mathcal{L}_{\text{physics}}$ on $\hat{M}_{\mathcal{A} \rightarrow \mathcal{B}}$ that address physical inconsistencies such as ground penetration, floating meshes and foot sliding. If HUMOS naïvely copies the same source motion $M_{\mathcal{A}}$ on the target body \mathcal{B} , it would result in motions $\hat{M}_{\mathcal{A} \rightarrow \mathcal{B}}$ that have significant ground penetration, floating meshes and foot sliding.

Our intuitive physics terms are fast, simple and fully-differentiable. Following [80], we design IP terms to address penetration, float, and foot sliding individually. Our $\mathcal{L}_{\text{penetrate}}$ minimizes the per-frame distance of the lowest vertex *below* the ground from the ground plane. $\mathcal{L}_{\text{float}}$ minimizes the per-frame distance of the lowest vertex *above* the ground from the ground plane. The foot sliding loss, $\mathcal{L}_{\text{slide}}$, minimizes the horizontal x-y component of the foot joint velocities if they are determined to be in ground contact using a distance threshold from the ground. We collate them together as $\mathcal{L}_{\text{physics}} = \mathcal{L}_{\text{penetrate}} + \mathcal{L}_{\text{float}} + \mathcal{L}_{\text{slide}}$.

3.4 Dynamic stability term

In the real world, motion is the result of internal muscular forces and external forces acting on the body and the surrounding scene. Human bodies are typically *stable*, i.e. they have the ability to control their body position and momentum during movement without falling over.

Tripathi *et al.* [70] successfully use the notion of static stability in 3D human pose and shape estimation to output physically-plausible and biomechanically stable poses from RGB images. In static poses, a body is considered stable if the gravity-projection of the center of mass (CoM) onto the ground is within the base of support (BoS). The base of support is defined as the convex hull of all points in contact with the ground. Since the base of support requires a convex hull computation that is not easily differentiable, [70] minimize the distance between an estimated center of pressure (CoP) and the projected CoM instead, and use it as a proxy static stability loss or energy term which is minimized during training and optimization. However, this *static* treatment of stability is only applicable to static poses. Humans are highly dynamic by nature and we need a general treatment of stability analysis that extends to all scenarios involving human movement and locomotion.

Dynamic stability extends this concept to bodies in motion. We follow the concept of zero-moment point (ZMP) [73], which has been widely used in robotics and biomechanics [41, 57]. Assuming flat ground, the ZMP is the point on the ground where the horizontal component of the moment of ground reaction force is zero. If this point lies within the base of support, the ZMP is equivalent to the center of pressure and the motion is considered dynamically stable (see Sup. Mat. video for an example in human gait).

The ZMP is defined as a function of the CoM’s acceleration and the net moment torques along the CoM and can be computed in closed form in a fully differentiable manner:

$$\mathcal{Z} = \mathcal{C}_m - \frac{n \times \mathcal{M}_{\mathcal{C}_m}^{gi}}{\mathcal{F}^{gi} \cdot n} \quad (3)$$

where \mathcal{C}_m is the projection of the center of mass onto the ground plane, and n is the normal to the ground plane. \mathcal{F}^{gi} is force of inertia calculated as

$$\mathcal{F}^{gi} = mg - ma_G \quad (4)$$

with m being the total mass of the body, g the acceleration of gravity, and a_G the acceleration of the center of mass \mathcal{G} . $\mathcal{M}_{\mathcal{C}_m}^{gi}$ is the moment around the projected center of mass \mathcal{C}_m

$$\mathcal{M}_{\mathcal{C}_m}^{gi} = \overrightarrow{\mathcal{C}_m \mathcal{G}} \times mg - \overrightarrow{\mathcal{C}_m \mathcal{G}} \times ma_G - \dot{\mathcal{H}}_G \quad (5)$$

where $\dot{\mathcal{H}}_G$ is the rate of change of angular momentum or torque at the center of mass.

We calculate the total body mass m by using the volume of the SMPL mesh as a proxy for total weight. To calculate center of mass \mathcal{G} and its acceleration as well as the moment $\mathcal{M}_{\mathcal{C}_m}^{gi}$, we distribute the total mass m to point masses at the vertices of the body mesh proportional to the volume of the body part they are part of. The accelerations are calculated numerically using the central differences. Please refer to Sup. Mat. for detailed derivations of the formulas.

Similar to Tripathi *et al.* [70] we use an estimation of the center of pressure as proxy for calculating the distance to the base of support. The center of pressure \mathcal{C}_P is calculated as weighted average of all vertices close to ground plane in a frame.

For dynamically stable motions, the ZMP and CoP should coincide. We, therefore, minimize the distance between ZMP and CoP and define our dynamic stability loss as

$$\mathcal{L}_{\text{dyn}} = \rho(\|\mathcal{C}_P - \mathcal{Z}\|_2) \quad (6)$$

where ρ is the Geman-McClure penalty function [25] which stabilizes training by making \mathcal{L}_{dyn} robust to noisy ZMP estimates.

Dynamic stability computation requires ground support. For sequences where the human is not supported by the ground, *i.e.* lowest vertex >25 cm, we disable the dynamic stability term during training. Thus, although the dynamic stability term is designed to help grounded sequences, it does not hurt non-grounded ones.

3.5 Latent embedding losses

To enable motion generation at inference time, we regularize the distributions of the latent embedding spaces, $z_{M_A} = \mathcal{N}(\mu_A, \Sigma_A)$ and $z_{M_{A \rightarrow B}} = \mathcal{N}(\mu_{A \rightarrow B}, \Sigma_{A \rightarrow B})$ to be similar to the normal distribution $\psi = \mathcal{N}(0, I)$ by minimizing the Kullback-Leibler (KL) divergence via

$$\mathcal{L}_{\text{KL}} = \text{KL}(z_{M_A}, \psi) + \text{KL}(z_{M_{A \rightarrow B}}, \psi) \quad (7)$$

Since the HUMOS latent embeddings encode motion style rather than identity-specific attributes, we also encourage the embeddings $z_{M_{\mathcal{A}}} \sim \mathcal{N}(\mu_{\mathcal{A}}, \Sigma_{\mathcal{A}})$ and $z_{M_{\mathcal{A} \rightarrow \mathcal{B}}} \sim \mathcal{N}(\mu_{\mathcal{A} \rightarrow \mathcal{B}}, \Sigma_{\mathcal{A} \rightarrow \mathcal{B}})$ to be as close as possible to each other via the the following cycle-consistent L1 loss:

$$\mathcal{L}_{\text{E}} = \|z_{M_{\mathcal{A}}} - z_{M_{\mathcal{A} \rightarrow \mathcal{B}}}\|_1 \quad (8)$$

The resulting total loss in HUMOS training is the weighted sum of all the individual loss terms:

$$\mathcal{L} = \lambda_{\text{cycle}} \mathcal{L}_{\text{cycle}} + \lambda_{\text{physics}} \mathcal{L}_{\text{physics}} + \lambda_{\text{dyn}} \mathcal{L}_{\text{dyn}} + \lambda_{\text{KL}} \mathcal{L}_{\text{KL}} + \lambda_{\text{E}} \mathcal{L}_{\text{E}} \quad (9)$$

The loss weights are determined empirically and set to $\lambda_{\text{cycle}} = 1$, $\lambda_{\text{physics}} = 1$, $\lambda_{\text{dyn}} = 0.0001$, $\lambda_{\text{KL}} = 10^{-5}$ and $\lambda_{\text{E}} = 10^{-2}$.

4 Experiments

We first discuss our data processing and implementation details in (Sec. 4.1). Next, we introduce baselines and the evaluation metrics (Sec. 4.2) used in our comparisons. Then, we discuss quantitative, perceptual and qualitative comparisons of our method with baselines (Sec. 4.3) and present an ablation study (Sec. 4.4).

4.1 Data and implementation details

For training, we use the AMASS dataset which contains 480 unique gender identities out of which 274 are male and 206 are female with diverse body shapes and sizes. Please refer to Sup. Mat. for a full analysis on the diversity and distribution of body shape β parameters in AMASS. We first subsample raw SMPL-H sequences from AMASS [49] to 20 fps following Guo *et al.* [32]. We augment the data by mirroring sequences left-to-right. We exclude sequences where the feet are more than 20 cm above the ground and where the lowest vertex across all frames is not grounded. This is to ensure ground support as it is an essential component for dynamic stability. We observed that normalizing sequences for consistent facing and start position in the first frame helped with training. We then extract input features by converting the root orientation and joint rotations to 6D form [86] and concatenate root translation, betas and gender. In addition, we augment the AMASS dataset by applying left-right flip augmentation to the pose parameters, effectively doubling the amount of training data. We show a step-by-step visualization of our data-processing pipeline in Sup. Mat.¹

We train our models for 1300 epochs with the AdamW [39, 48] optimizer using a fixed learning rate 10^{-5} and a batch size of 60. Both our encoder and decoder consist of 6 transformer layers. We train with sequence length $T = 200$ frames on

¹ All datasets were obtained and used only by the authors affiliated with academic institutions.

arbitrary length clips sampled from AMASS. For clips longer or shorter than T frames, we extend or clip the sampling interval by either including neighboring frames or dropping boundary frames. For short AMASS videos with $< T$ frames, we repeat the last frame. Please refer to Sup. Mat. for more details.

4.2 Baselines and evaluation metrics

We focus on the task of shape-conditioned motion reconstruction for evaluating the performance of HUMOS. Since no existing baseline directly addresses shape-conditioned motion reconstruction, we create new baselines by combining a state-of-the-art motion generation model, TEMOS [56] and retargeting its output motions to a target body shape using 1) naïve retargeting and 2) using the commercial retargeting system, Rokoko [62]. For these experiments, we reconstruct the same motions from the AMASS test-split for both the TEMOS baselines and HUMOS.

TEMOS generates motions for a canonicalized mean-shape SMPL body by directly regressing the pose and global root-joint translation. For a fair comparison, we use its “unconditional” variant which does not require any text inputs. For obtaining motions for a target body, we do simple retargeting where we randomly sample identities from AMASS and naïvely *copy* the target $\beta_{\mathcal{B}}$ and gender $\mathcal{G}_{\mathcal{B}}$ parameters to the motions obtained from TEMOS. We call this baseline, “TEMOS-Simple”. Intuitively, naïvely copying a new identity to a neutral mean-shape body motion will result in ground penetration, floating and foot sliding artefacts. We address ground penetration by translating the whole motion sequence such that the lowest vertex in the sequence is on the ground. We refer to this baseline as “TEMOS-Simple-G”. For another strong baseline, we use the Rokoko retargeting system to retarget TEMOS generated motions to the target body. We call this baseline “TEMOS-Rokoko”, with “TEMOS-Rokoko-G” being the variant where we ground the Rokoko output sequences as described above. For consistency in evaluation, given an input motion, we sample the same target identities across our method and all baselines.

Evaluation Metrics. For evaluating the physical plausibility of generated motions, we use the physics-based metrics suggested by Yuan *et al.* [80]. The *Penetrate* metric measures ground penetration by computing the distance (in cm) between the ground and the lowest body mesh vertex *below* the ground. *Float* measuring the amount of unsupported floating by computing the distance (in cm) between the ground and the lowest body mesh vertex *above* the ground. *Skate* measures foot skating by computing the percentage of adjacent frames where the foot joints in contact with the ground have a non-zero average velocity. We also report two metrics for measuring dynamic stability. *Dyn. Stability* computes the percentage of frames where the ZMP is outside the base of support. The *BoSDist* metric measures the distance of the ZMP to the closest edge of the BoS convex hull, if the ZMP lies outside the BoS, indicating the pose is dynamically unstable. For details, please refer to Sup. Mat.

4.3 Comparison to baselines

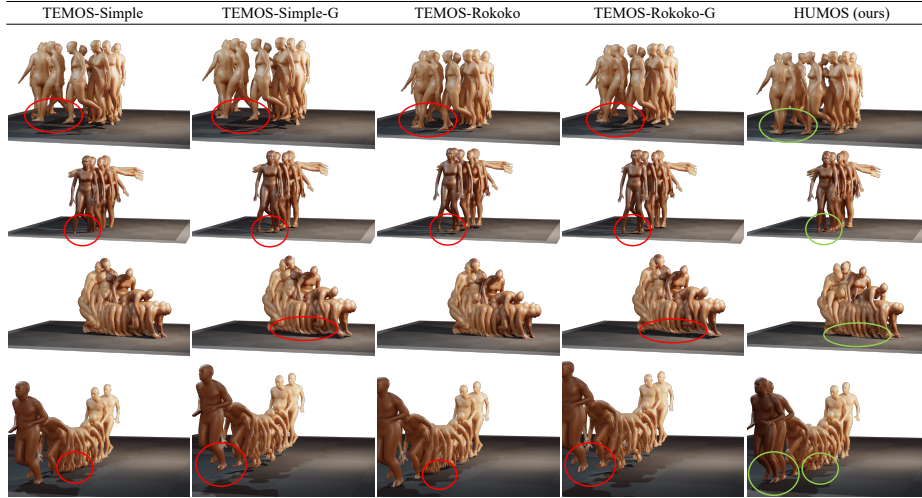


Fig. 3: Qualitative comparison of shape-conditioned motion generation. Each row represents generations across different methods for a unique body shape and gender. HUMOS generated motions are more realistic, physically plausible and dynamically stable compared to baselines. The red circles on the baseline methods highlight issues such as floating, penetrations, and foot skating, compared to more realistic results on highlighted in green with HUMOS. **Q Zoom in.**

Quantitative. We summarize our main results in Tab. 1. As we lack ground-truth motions for the target body shape, we rely on physics and stability metrics to compare our method with baselines. Our method substantially outperforms on all metrics except on *Penetrate*. As expected, the lowest ground penetration is observed for TEMOS-Simple-G and TEMOS-Rokoko-G as both were specifically altered to ground the sequence. However, this comes at the cost of increasing the *Float* metric. In contrast, HUMOS simultaneously improves both *Penetrate* and *Float* indicating that the network learns to modify body pose (*e.g.* foot tilt) in addition to learning the correct global translation for grounding the motion. HUMOS also improves over baselines on foot skating, achieving a $\sim 7.3\%$ *Skate* compared to 20% and 27% for the TEMOS-Rokoko and TEMOS-Simple baselines. HUMOS’s motions are also dynamically stable in 71.9% of all frames, a significant improvement of 16% over the closest baseline. Even for dynamically unstable poses, when the ZMP is outside the BoS, it is close to the BoS edge as indicated by the low *BoSDict* metric for our method.

Qualitative. We provide additional qualitative comparisons with baselines in Fig. 3. Each row represents the same pair of source motion and target body across all methods. We highlight physical plausibility issues such as foot-skate, ground penetration and floating in red. The green highlighted region points to

the improvement in HUMOS’s results over baselines. HUMOS motions showcase realistic ground support. To better evaluate our performance on shape-conditioned motion generation, foot-sliding and dynamic stability, please watch our Sup. Mat. video.

Table 1: Comparison of HUMOS with baselines on the shape-conditioned motion reconstruction task.

Method	Penetrate (cm) ↓	Float (cm) ↓	Skate (%) ↓	Dyn. Stability (%) ↑	BoS Dist (cm) ↓
TEMOS [56]-Simple	6.82	6.55	27.07	45.85	16.94
TEMOS [56]-Simple-G	0.75	4.39	27.07	45.85	16.94
TEMOS [56]-Rokoko [63]	4.14	3.85	20.05	55.92	16.58
TEMOS [56]-Rokoko [63]-G	0.75	4.44	20.05	55.92	16.58
HUMOS	1.23	1.04	7.37	71.9	14.62

Perceptual Study. To evaluate realism of the generated motions given the target body shape, we perform a human study on Amazon Mechanical Turk (AMT). We randomly select 30 videos generated from our method and the next closest baselines, “TEMOS-Rokoko” and “TEMOS-Rokoko-G”. The participants are shown a single video and after watching the whole video at least once, are allowed to select their response to the question “How realistic is the motion given this body shape?” on a Likert scale of scores between 5 (completely realistic) to 1 (completely unrealistic). Each rating task was completed by 25 participants. In the study, we added 4 catch trials, 2 containing ground-truth AMASS motions and 2 containing significant ground penetration. 17 out of 75 participants who failed the catch trials were excluded from our study. As shown in Tab. 2 participants prefer HUMOS motions and give it an average rating of 3.64 out of 5, compared to 3.25 for TEMOS-Rokoko and 3.19 for TEMOS-Rokoko-G. Curiously, between the two baselines, participants preferred the motions without a heuristics-based grounding indicating that the use of heuristics for motion retargeting struggles with generalization. We include more details about the study and layout in Sup. Mat.

Table 2: Perceptual study comparing HUMOS with two closest baselines. Given a video of a generated motion, participants select 5-point ratings for the question “How realistic is the motion given this body shape?”

Method	Average Rating ↑	Std. Dev. ↓
TEMOS-Rokoko	3.25	1.26
TEMOS-Rokoko-G	3.19	1.27
HUMOS	3.64	1.11

4.4 Ablation Study

We evaluate the importance of our key contributions, $\mathcal{L}_{\text{cycle}}$, $\mathcal{L}_{\text{physics}}$ and \mathcal{L}_{dyn} in Tab. 3. As shown, $\mathcal{L}_{\text{cycle}}$ alone achieves a significant $\sim 33\%$ improvement in

Penetrate, $\sim 32\%$ in *Float*, $\sim 25\%$ in *Skate* over the TEMOS-Rokoko baseline indicating that our cycle-consistent training paradigm is effective in training HUMOS for the shape-conditioned motion generation task. Adding $\mathcal{L}_{\text{physics}}$ further improves physical plausibility, resulting in the biggest improvement in foot skating ($\sim 47\%$). While both $\mathcal{L}_{\text{cycle}}$ and $\mathcal{L}_{\text{physics}}$ help, adding \mathcal{L}_{dyn} results in the best HUMOS configuration across all metrics. With all losses active, HUMOS motions are dynamically stable 71.9% the times.

Table 3: Ablation study comparing the improvements from cycle-consistent training ($\mathcal{L}_{\text{cycle}}$), physics losses ($\mathcal{L}_{\text{physics}}$) and the dynamic stability term (\mathcal{L}_{dyn}).

Method	Penetrate (cm) ↓	Float (cm) ↓	Skate (%) ↓	Dyn. Stability (%) ↑	BoS Dist (cm) ↓
TEMOS-Rokoko	4.14	3.85	20.05	55.92	16.58
$\mathcal{L}_{\text{cycle}}$	2.74	2.62	15.04	64.00	16.96
$\mathcal{L}_{\text{cycle}} + \mathcal{L}_{\text{physics}}$	1.55	1.44	7.93	67.82	16.41
$\mathcal{L}_{\text{cycle}} + \mathcal{L}_{\text{physics}} + \mathcal{L}_{\text{dyn}}$	1.23	1.04	7.37	71.9	14.62

5 Conclusion

In this paper we presented a method for shape-conditioned motion generation that used a set of physically inspired constraints to allow for self-supervised disentanglement of character motion and identity. This allows for motion generation and retargetting of a higher quality than previous methods both qualitatively and quantitatively.

In terms of limitations, although our method represents an improvement over previous work there are still motion artefacts introduced by the model. Additionally, the differences in the style of motion produced by characters of very different body shapes remain subtle. This may be due to the limited shape diversity in the training set. In the future it would be interesting to examine how this data distribution affects the diversity and generalization capabilities of the model. We also do not take into account self-penetrations that may arise during shape-conditioned motion generation. Addressing this would be another promising direction for future work. While human motion is influenced by both body shape and individual motion style, we only consider body shape. Motion style includes factors like emotional state, physiological impediments, societal influences, and personal biases, which are not annotated in existing mocap datasets. With style-specific annotations, it would be useful to extend HUMOS to include style attributes as additional conditioning signals.

Acknowledgements. We sincerely thank Tsvetelina Alexiadis, Alpar Cseke, Tomasz Niewiadomski, and Taylor McConnell for facilitating the perceptual study, and Giorgio Becherini for his help with the Rokoko baseline. We are grateful to Iain Matthews, Brian Karis, Nikos Athanasiou, Markos Diomataris, and Mathis Petrovich for valuable discussions and advice. Their invaluable contributions enriched this research significantly.

References

1. Abdul-Massih, M., Yoo, I., Benes, B.: Motion style retargeting to characters with different morphologies. *Computer Graphics Forum* **36**(6), 86–99 (2017). <https://doi.org/https://doi.org/10.1111/cgf.12860>, <https://onlinelibrary.wiley.com/doi/abs/10.1111/cgf.12860> 3
2. Aberman, K., Li, P., Lischinski, D., Sorkine-Hornung, O., Cohen-Or, D., Chen, B.: Skeleton-aware networks for deep motion retargeting. *ACM Trans. Graph.* **39**(4) (aug 2020). <https://doi.org/10.1145/3386569.3392462>, <https://doi.org/10.1145/3386569.3392462> 3
3. Aberman, K., Wu, R., Lischinski, D., Chen, B., Cohen-Or, D.: Learning character-agnostic motion for motion retargeting in 2d. *ACM Transactions on Graphics* **38**(4), 1–14 (Jul 2019). <https://doi.org/10.1145/3306346.3322999>, <http://dx.doi.org/10.1145/3306346.3322999> 3
4. Ahn, H., Ha, T., Choi, Y., Yoo, H., Oh, S.: Text2Action: Generative adversarial synthesis from language to action. In: *International Conference on Robotics and Automation (ICRA)* (2018) 4
5. Ahuja, C., Morency, L.P.: Language2pose: Natural language grounded pose forecasting. In: *2019 International Conference on 3D Vision (3DV)*. pp. 719–728. *IEEE* (2019) 4
6. Aksan, E., Kaufmann, M., Hilliges, O.: Structured prediction helps 3d human motion modelling. In: *Proceedings of the IEEE International Conference on Computer Vision*. pp. 7144–7153 (2019) 4
7. Aliakbarian, S., Saleh, F.S., Salzmann, M., Petersson, L., Gould, S.: A stochastic conditioning scheme for diverse human motion prediction. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 5223–5232 (2020) 4
8. Athanasiou, N., Petrovich, M., Black, M.J., Varol, G.: TEACH: temporal action composition for 3D humans. In: *3DV*. pp. 414–423. *IEEE* (2022) 4, 5
9. Athanasiou, N., Petrovich, M., Black, M.J., Varol, G.: SINC: Spatial composition of 3D human motions for simultaneous action generation. In: *Proc. International Conference on Computer Vision (ICCV)*. pp. 9984–9995 (Oct 2023) 4, 5
10. Bao, F., Li, C., Sun, J., Zhu, J., Zhang, B.: Estimating the optimal covariance with imperfect mean in diffusion probabilistic models. In: *International Conference on Machine Learning* (2022) 4
11. Bao, F., Li, C., Zhu, J., Zhang, B.: Analytic-DPM: An analytic estimate of the optimal reverse variance in diffusion probabilistic models. In: *International Conference on Learning Representations* (2022) 4
12. Barsoum, E., Kender, J., Liu, Z.: Hp-gan: Probabilistic 3d human motion prediction via gan. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*. pp. 1418–1427 (2018) 4
13. Basset, J., Wuhrer, S., Boyer, E., Multon, F.: Contact preserving shape transfer for rigging-free motion retargeting. In: *Proceedings of the 12th ACM SIGGRAPH Conference on Motion, Interaction and Games. MIG '19, Association for Computing Machinery, New York, NY, USA* (2019). <https://doi.org/10.1145/3359566.3360075>, <https://doi.org/10.1145/3359566.3360075> 3
14. Bergamin, K., Clavet, S., Holden, D., Forbes, J.R.: Drecon: data-driven responsive control of physics-based characters. *ACM Transactions on Graphics (TOG)* **38**(6), 1–11 (2019) 4

15. Bhattacharya, U., Childs, E., Rewkowski, N., Manocha, D.: Speech2affectivegestures: Synthesizing co-speech gestures with generative adversarial affective expression learning. In: Proceedings of the 29th ACM International Conference on Multimedia. pp. 2027–2036 (2021) 4
16. Bhattacharya, U., Rewkowski, N., Banerjee, A., Guhan, P., Bera, A., Manocha, D.: Text2gestures: A transformer-based network for generating emotive body gestures for virtual agents. In: 2021 IEEE Virtual Reality and 3D User Interfaces (VR). pp. 1–10. IEEE (2021) 4
17. Celikkan, U., Yaz, I.O., Capin, T.: Example-based retargeting of human motion to arbitrary mesh models. *Computer Graphics Forum* **34**(1), 216–227 (2015). <https://doi.org/https://doi.org/10.1111/cgf.12507>, <https://onlinelibrary.wiley.com/doi/abs/10.1111/cgf.12507> 3
18. Chen, X., Jiang, B., Liu, W., Huang, Z., Fu, B., Chen, T., Yu, G.: Executing your commands via motion diffusion in latent space. In: CVPR. pp. 18000–18010. IEEE (2023) 2, 4
19. Choi, J., Kim, S., Jeong, Y., Gwon, Y., Yoon, S.: ILVR: Conditioning method for denoising diffusion probabilistic models. arXiv preprint arXiv:2108.02938 (2021) 4
20. Choi, K.J., Ko, H.S.: On-line motion retargetting. In: Proceedings. Seventh Pacific Conference on Computer Graphics and Applications (Cat. No.PR00293). pp. 32–42 (1999). <https://doi.org/10.1109/PCCGA.1999.803346> 3
21. Dhariwal, P., Nichol, A.Q.: Diffusion models beat GANs on image synthesis. *Advances in Neural Information Processing Systems* (2021) 4
22. Dockhorn, T., Vahdat, A., Kreis, K.: GENIE: Higher-order denoising diffusion solvers. *Advances in Neural Information Processing Systems* (2022) 4
23. Fragkiadaki, K., Levine, S., Felsen, P., Malik, J.: Recurrent network models for human dynamics. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 4346–4354 (2015) 4
24. Fussell, L., Bergamin, K., Holden, D.: Supertrack: motion tracking for physically simulated characters using supervised learning. *ACM Trans. Graph.* **40**(6) (dec 2021). <https://doi.org/10.1145/3478513.3480527>, <https://doi.org/10.1145/3478513.3480527> 4
25. Geman, S.: Statistical methods for tomographic image restoration. *Bull. Internat. Statist. Inst.* **52**, 5–21 (1987) 9, 22
26. Ghosh, A., Cheema, N., Oguz, C., Theobalt, C., Slusallek, P.: Synthesis of compositional animations from textual descriptions. In: International Conference on Computer Vision (ICCV) (2021) 4
27. Ghosh, P., Song, J., Aksan, E., Hilliges, O.: Learning human motion models for long-term predictions. In: 2017 International Conference on 3D Vision (3DV). pp. 458–466. IEEE (2017) 4
28. Ginosar, S., Bar, A., Kohavi, G., Chan, C., Owens, A., Malik, J.: Learning individual styles of conversational gesture. In: Computer Vision and Pattern Recognition (CVPR) (2019) 4
29. Gomes, T., Martins, R., Ferreira, J., Azevedo, R., Torres, G., Nascimento, E.: A Shape-Aware Retargeting Approach to Transfer Human Motion and Appearance in Monocular Videos. *International Journal of Computer Vision* (Apr 2021). <https://doi.org/10.1007/s11263-021-01471-x>, <https://inria.hal.science/hal-03257490>, 19 pages, 13 figures 4
30. Gopalakrishnan, A., Mali, A., Kifer, D., Giles, L., Ororbias, A.G.: A neural temporal model for human motion prediction. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 12116–12125 (2019) 4

31. Grenander, U., Miller, M.I.: Representations of knowledge in complex systems. *Journal of the Royal Statistical Society: Series B (Methodological)* **56**(4), 549–581 (1994) [4](#)
32. Guo, C., Zou, S., Zuo, X., Wang, S., Ji, W., Li, X., Cheng, L.: Generating diverse and natural 3D human motions from text. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 5152–5161 (June 2022) [4](#), [10](#)
33. Guo, C., Zuo, X., Wang, S., Zou, S., Sun, Q., Deng, A., Gong, M., Cheng, L.: Action2motion: Conditioned generation of 3d human motions. In: *Proceedings of the 28th ACM International Conference on Multimedia*. pp. 2021–2029 (2020) [4](#)
34. Habibie, I., Holden, D., Schwarz, J., Yearsley, J., Komura, T.: A recurrent variational autoencoder for human motion synthesis. In: *British Machine Vision Conference (BMVC)* (2017) [4](#)
35. He, C., Saito, J., Zachary, J., Rushmeier, H.E., Zhou, Y.: NeMF: Neural motion fields for kinematic animation. In: *NeurIPS* (2022) [4](#), [24](#)
36. Holden, D., Saito, J., Komura, T.: A deep learning framework for character motion synthesis and editing. *ACM Transactions on Graphics (TOG)* **35**(4), 1–11 (2016) [4](#)
37. Hoyet, L., McDonnell, R., O’Sullivan, C.: Push it real: perceiving causality in virtual interactions. *ACM Trans. Graph.* **31**(4), 90:1–90:9 (2012) [5](#)
38. Kang, H.j., Hashimoto, K., Kondo, H., Hattori, K., Nishikawa, K., Hama, Y., Lim, H.o., Takanishi, A., Suga, K., Kato, K.: Realization of biped walking on uneven terrain by new foot mechanism capable of detecting ground surface. In: *2010 IEEE International Conference on Robotics and Automation*. pp. 5167–5172 (2010). <https://doi.org/10.1109/ROBOT.2010.5509348> [2](#)
39. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: *ICLR* (2015) [10](#)
40. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. In: *ICLR* (2014) [6](#)
41. Kondak, K., Hommel, G.: Control and online computation of stable movement for biped robots. *IEEE/RSJ International Conference on Intelligent Robots and Systems* **1**, 874–879 (2003) [8](#)
42. Lee, H., Yang, X., Liu, M., Wang, T., Lu, Y., Yang, M., Kautz, J.: Dancing to music. In: *Neural Information Processing Systems (NeurIPS)* (2019) [4](#)
43. Lee, S., Kang, T., Park, J., Lee, J., Won, J.: Same: Skeleton-agnostic motion embedding for character animation. In: *SIGGRAPH Asia 2023 Conference Papers*. SA ’23, Association for Computing Machinery, New York, NY, USA (2023). <https://doi.org/10.1145/3610548.3618206>, <https://doi.org/10.1145/3610548.3618206> [3](#)
44. Li, B., Zhao, Y., Zhelun, S., Sheng, L.: Danceformer: Music conditioned 3d dance generation with parametric motion transformer. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. vol. 36, pp. 1272–1279 (2022) [4](#)
45. Li, J., Yin, Y., Chu, H., Zhou, Y., Wang, T., Fidler, S., Li, H.: Learning to generate diverse dance motions with transformer. *arXiv preprint arXiv:2008.08171* (2020) [4](#)
46. Li, R., Yang, S., Ross, D.A., Kanazawa, A.: Ai choreographer: Music conditioned 3d dance generation with aist++. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 13401–13412 (2021) [4](#)
47. Loper, M., Mahmood, N., Romero, J., Pons-Moll, G., Black, M.J.: SMPL: A skinned multi-person linear model. *Transactions on Graphics (TOG)* **34**(6), 248:1–248:16 (2015) [2](#), [4](#), [5](#)
48. Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: *ICLR* (2017), <https://api.semanticscholar.org/CorpusID:53592270> [10](#)

49. Mahmood, N., Ghorbani, N., F. Troje, N., Pons-Moll, G., Black, M.J.: AMASS: Archive of motion capture as surface shapes. In: International Conference on Computer Vision (ICCV). pp. 5441–5450 (2019) [2](#), [4](#), [6](#), [10](#), [23](#)
50. Makoviychuk, V., Wawrzyniak, L., Guo, Y., Lu, M., Storey, K., Macklin, M., Hoeller, D., Rudin, N., Allshire, A., Handa, A., State, G.: Isaac gym: High performance GPU based physics simulation for robot learning. In: Vanschoren, J., Yeung, S. (eds.) Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1, NeurIPS Datasets and Benchmarks 2021, December 2021, virtual (2021), <https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/hash/28dd2c7955ce926456240b2ff0100bde-Abstract-round2.html> [4](#)
51. Motion builder. <https://www.autodesk.com/products/motionbuilder/overview> [3](#)
52. Peng, X.B., Abbeel, P., Levine, S., van de Panne, M.: Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *ACM Transactions on Graphics (TOG)* **37**(4), 1–14 (2018) [4](#)
53. Peng, X.B., Kanazawa, A., Malik, J., Abbeel, P., Levine, S.: Sfv: Reinforcement learning of physical skills from videos. *ACM Transactions on Graphics (TOG)* **37**(6), 1–14 (2018) [4](#)
54. Peng, X.B., van de Panne, M.: Learning locomotion skills using deeprl: Does the choice of action space matter? In: Proceedings of the ACM SIGGRAPH/Eurographics Symposium on Computer Animation. pp. 1–13 (2017) [4](#)
55. Petrovich, M., Black, M.J., Varol, G.: Action-conditioned 3D human motion synthesis with transformer VAE. In: ICCV. pp. 10965–10975. IEEE (2021) [2](#), [4](#), [5](#)
56. Petrovich, M., Black, M.J., Varol, G.: TEMOS: generating diverse human motions from textual descriptions. In: ECCV (22). Lecture Notes in Computer Science, vol. 13682, pp. 480–497. Springer (2022) [2](#), [4](#), [5](#), [6](#), [11](#), [13](#)
57. Popovic, M.B., Goswami, A., Herr, H.: Ground reference points in legged locomotion: Definitions, biological trajectories and control implications. *International Journal of Robotics Research* **24**(10) (2005) [2](#), [8](#)
58. Regateiro, J., Boyer, E.: Temporal shape transfer network for 3d human motion. In: 2022 International Conference on 3D Vision (3DV). pp. 424–432 (2022). <https://doi.org/10.1109/3DV57658.2022.00054> [4](#)
59. Reitsma, P.S.A., Pollard, N.S.: Perceptual metrics for character animation: sensitivity to errors in ballistic motion. *ACM Trans. Graph.* **22**(3), 537–542 (2003) [5](#)
60. Rempe, D., Birdal, T., Hertzmann, A., Yang, J., Sridhar, S., Guibas, L.J.: HuMoR: 3D human motion model for robust pose estimation. In: International Conference on Computer Vision (ICCV). pp. 11468–11479. IEEE (2021) [4](#), [5](#)
61. Ren, Z., Pan, Z., Zhou, X., Kang, L.: Diffusion motion: Generate text-guided 3d human motion by diffusion model. arXiv preprint arXiv:2210.12315 (2022) [4](#)
62. Rokoko. <https://www.rokoko.com/> [3](#), [11](#)
63. Rokoko: Rokoko studio live plugin for blender. <https://github.com/Rokoko/rokoko-studio-live-blender> (2023) [13](#)
64. Schulman, J., Moritz, P., Levine, S., Jordan, M.I., Abbeel, P.: High-dimensional continuous control using generalized advantage estimation. In: Bengio, Y., LeCun, Y. (eds.) 4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2–4, 2016, Conference Track Proceedings (2016) [4](#)
65. Schulman, J., Wolski, F., Dhariwal, P., Radford, A., Klimov, O.: Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347 (2017) [4](#)

66. Shimada, S., Golyanik, V., Xu, W., Pérez, P., Theobalt, C.: Neural monocular 3d human motion capture with physical awareness. *ACM Transactions on Graphics (ToG)* **40**(4), 1–15 (2021) [4](#)
67. Shimada, S., Golyanik, V., Xu, W., Theobalt, C.: Physcap: Physically plausible monocular 3d motion capture in real time. *ACM Transactions on Graphics (TOG)* **39**(6) (2020) [4](#)
68. Taheri, O., Choutas, V., Black, M.J., Tzionas, D.: GOAL: Generating 4D whole-body motion for hand-object grasping. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 13253–13263 (2022) [4](#)
69. Tevet, G., Raab, S., Gordon, B., Shafir, Y., Cohen-Or, D., Bermano, A.H.: Human motion diffusion model. In: *ICLR. OpenReview.net* (2023) [2](#), [4](#)
70. Tripathi, S., Müller, L., Huang, C.H.P., Omid, T., Black, M.J., Tzionas, D.: 3D human pose estimation via intuitive physics. In: *Computer Vision and Pattern Recognition (CVPR)*. pp. 4713–4725 (2023), <https://ipman.is.tue.mpg.de> [2](#), [4](#), [8](#), [9](#), [21](#)
71. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. In: *NeurIPS*. vol. 30 (2017) [5](#), [6](#)
72. Villegas, R., Ceylan, D., Hertzmann, A., Yang, J., Saito, J.: Contact-aware re-targeting of skinned motion. In: *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*. pp. 9700–9709 (2021). <https://doi.org/10.1109/ICCV48922.2021.00958> [3](#)
73. Vukobratović, M., Borovac, B.: Zero-moment point—thirty five years of its life. In: *International Journal of Humanoid Robotics*. pp. 157–173 (2004) [8](#)
74. Wang, J., Wen, C., Fu, Y., Lin, H., Zou, T., Xue, X., Zhang, Y.: Neural pose transfer by spatially adaptive instance normalization. *CoRR* **abs/2003.07254** (2020), <https://arxiv.org/abs/2003.07254> [4](#)
75. Won, J., Gopinath, D., Hodgins, J.: A scalable approach to control diverse behaviors for physically simulated characters. *ACM Transactions on Graphics (TOG)* **39**(4), 33–1 (2020) [4](#)
76. Yamane, K., Ariki, Y., Hodgins, J.: Animating non-humanoid characters with human motion data. In: *Proceedings of the 2010 ACM SIGGRAPH/Eurographics Symposium on Computer Animation*. p. 169–178. SCA '10, Eurographics Association, Goslar, DEU (2010) [3](#)
77. Yi, X., Zhou, Y., Habermann, M., Shimada, S., Golyanik, V., Theobalt, C., Xu, F.: Physical inertial poser (pip): Physics-aware real-time human motion tracking from sparse inertial sensors. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 13167–13178 (2022) [4](#)
78. Yuan, Y., Kitani, K.: Dlow: Diversifying latent flows for diverse human motion prediction. In: *Proceedings of the European Conference on Computer Vision (ECCV)*. pp. 346–364. Springer (2020) [4](#)
79. Yuan, Y., Kitani, K.: Residual force control for agile human behavior imitation and extended motion synthesis. *Advances in Neural Information Processing Systems* (2020) [4](#)
80. Yuan, Y., Song, J., Iqbal, U., Vahdat, A., Kautz, J.: Physdiff: Physics-guided human motion diffusion model. In: *ICCV*. pp. 15964–15975. IEEE (2023) [8](#), [11](#)
81. Yuan, Y., Wei, S.E., Simon, T., Kitani, K., Saragih, J.: Simpoe: Simulated character control for 3d human pose estimation. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2021) [4](#)
82. Zell, P., Wandt, B., Rosenhahn, B.: Joint 3d human motion capture and physical analysis from monocular videos. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*. pp. 17–26 (2017) [4](#)

83. Zhang, J., Weng, J., Kang, D., Zhao, F., Huang, S., Zhe, X., Bao, L., Shan, Y., Wang, J., Tu, Z.: Skinned motion retargeting with residual perception of motion semantics & geometry. In: 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 13864–13872 (2023). <https://doi.org/10.1109/CVPR52729.2023.01332> 4
84. Zhang, M., Cai, Z., Pan, L., Hong, F., Guo, X., Yang, L., Liu, Z.: Motiondiffuse: Text-driven human motion generation with diffusion model. arXiv preprint arXiv:2208.15001 (2022) 4
85. Zhou, K., Bhatnagar, B.L., Pons-Moll, G.: Unsupervised shape and pose disentanglement for 3d meshes. CoRR **abs/2007.11341** (2020), <https://arxiv.org/abs/2007.11341> 3
86. Zhou, Y., Barnes, C., Lu, J., Yang, J., Li, H.: On the continuity of rotation representations in neural networks. In: CVPR. pp. 5745–5753. Computer Vision Foundation / IEEE (2019) 5, 7, 10, 24
87. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: ICCV. pp. 2223–2232 (2017) 3, 7

Appendix

In the supplementary materials, we provide a detailed derivation of the ZMP and the dynamic stability term (Appendix A), analyze the effect of body shape on motion (Appendix B), provide additional qualitative results (Appendix C), ablations for the latent embedding losses (Appendix D), a discussion on AMASS shape diversity (Appendix E), and finally additional implementation details (Appendix F).

Video. Our research focuses on humans in motion with diverse body shapes and sizes, making motion a critical aspect of our results. Given the difficulty of conveying motion quality through a static document, we strongly recommend that readers view the provided supplemental video for an in-depth overview of our methodology and findings.

A Detailed derivation for the Zero Moment Point (ZMP) and the dynamic stability term

Before we compute the ZMP, we first compute the body Center of Mass (CoM) by adapting the CoM formulation of Tripathi *et al.* [70] to dynamic humans. For every sequence, we use their body part segmentation and the differentiable “*close-translate-fill*” [70] to compute per-part volumes \mathcal{V}^{P_i} by splitting the mesh in the first frame into 10 parts. Using the per-part volumes, the CoM is calculated for time instance t , as a volume weighted-average of $N_U = 6890$ mesh vertex points.

$$\mathcal{G}_t = \frac{\sum_{i=1}^{N_U} \mathcal{V}^{P_{v_i}} v_{i_t}}{\sum_{i=1}^{N_U} \mathcal{V}^{P_{v_i}}}, \quad (\text{S.1})$$

The acceleration of the CoM, $a_{\mathcal{G}}$, is obtained using the central difference as,

$$a_{\mathcal{G}_t} = \frac{\mathcal{G}_{t+1} - 2\mathcal{G}_t + \mathcal{G}_{t-1}}{\Delta t^2} \quad (\text{S.2})$$

With $a_{\mathcal{G}_t}$, the force of inertia, \mathcal{F}^{g_i} , is computed as

$$\mathcal{F}^{g_i} = mg - ma_{\mathcal{G}} \quad (\text{S.3})$$

where m is the body mass. The moment around the projected CoM, \mathcal{C}_m , is

$$\mathcal{M}_{\mathcal{C}}^{g_i} = \overrightarrow{\mathcal{C}_m \mathcal{G}} \times mg - \overrightarrow{\mathcal{C}_m \mathcal{G}} \times ma_{\mathcal{G}} - \dot{\mathcal{H}}_{\mathcal{G}} \quad (\text{S.4})$$

where $\overrightarrow{\mathcal{C}_m \mathcal{G}}$ is the vector joining the projected CoM, \mathcal{C}_m with the actual CoM, \mathcal{G} and $\dot{\mathcal{H}}_{\mathcal{G}}$ is the rate of change of angular momentum at the CoM. For $\dot{\mathcal{H}}_{\mathcal{G}}$, we equally distribute the total m to point masses at the vertices of the body mesh proportional to the volume of the body part they are part of. The per-vertex mass and acceleration is

$$m_{v_i} = \frac{\mathcal{V}^{P_{v_i}}}{\sum_{i=1}^{N_U} \mathcal{V}^{P_{v_i}}} m, \quad a_{v_i} = \frac{v_{i_{t+1}} - 2v_{i_t} + v_{i_{t-1}}}{\Delta t^2} \quad (\text{S.5})$$

And $\dot{\mathcal{H}}_G$ is

$$\dot{\mathcal{H}}_G = \sum_{i=1}^{N_U} \vec{v}_i \dot{\mathcal{G}} \times m_{v_i} a_{v_i} \quad (\text{S.6})$$

Finally, the ZMP is computed in closed-form as

$$\mathcal{Z} = \mathcal{C}_m - \frac{n \times \mathcal{M}_{\mathcal{C}_m}^{g_i}}{\mathcal{F}^{g_i} \cdot n} \quad (\text{S.7})$$

For CoP computation, we follow Tripathi *et al.* and uniformly sample the body mesh into $N_p = 20000$ uniformly sampled surface points. We, then, use their heuristic pressure field to compute per-point, p_i , pressure as

$$\rho_i = \begin{cases} 1 - \alpha h(p_i) & \text{if } h(p_i) < 0, \\ e^{-\gamma h(p_i)} & \text{if } h(p_i) \geq 0, \end{cases} \quad (\text{S.8})$$

where $\alpha = 100$ and $\gamma = 10$ are scalar hyperparameters set empirically. The CoP is computed as,

$$\mathcal{C}_p = \frac{\sum_{i=1}^{N_p} \rho_i p_i}{\sum_{i=1}^{N_p} \rho_i}. \quad (\text{S.9})$$

With the ZMP and the CoP, known the dynamic stability loss is defined as,

$$\mathcal{L}_{\text{dyn}} = \rho(\|\mathcal{C}_p - \mathcal{Z}\|_2) \quad (\text{S.10})$$

where ρ is the Geman-McClure penalty function [25].

B Effect of body shape on motion

In Fig. S.1 (left), we assess the diversity of HUMOS generated motions across 100 β parameters obtained by interpolating between a short male and a tall male body. We report the maximum right knee joint angle ($|\theta|$) for the same walk sequence shown in the Sup. Mat. (SM) video (05:34). The graph illustrates that taller people bend their knees less for the same walking motion, indicating body parameters affect movement. Similarly, in Fig. S.1 (center), we plot the right hand joint velocity across six different identities in the same walk sequence [SM video (05:34)]. The joint velocities differ across subjects in corresponding frames implying diversity induced by body shape variation. In Fig. S.1 (right), we qualitatively show the same frame of the jumping jack sequence [SM video (05:47)] where the different arm positions indicate motion diversity.

C Additional Qualitative Results

We include additional comparisons with baselines in Fig. S.2. For video results, we recommend watching the **supplementary video**.

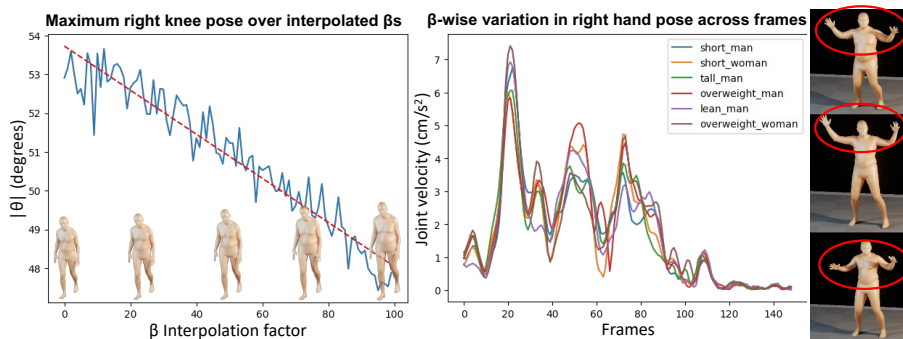


Fig. S.1: Effect of body shape across (left) interpolated β parameters, (center) 150 frames for 6 different identities and (right) different identities for the same jumping jack frame **Q Zoom in**.

D Additional Ablations

In Tab. S.1, we conduct additional ablations to analyze the effect of latent embedding losses, \mathcal{L}_E and \mathcal{L}_{KL} . We take the HUMOS model and successively remove the two loss terms individually. On ablating \mathcal{L}_E during training, we observe a small improvement in ground penetration and float. However, the skate and dynamic stability metrics worsen. While the effect of \mathcal{L}_E is minimal in terms of metrics, we empirically note faster and stable convergence when using it during training. \mathcal{L}_{KL} also results in a slight improvement in physics metrics at the cost of dynamic stability. A significant advantage, however, of using \mathcal{L}_{KL} is that it adds structure to the shape-agnostic latent space, making realistic motion generation easier.

Table S.1: Ablation study for latent embedding losses, \mathcal{L}_E and \mathcal{L}_{KL}

Method	Penetrate (cm) ↓	Float (cm) ↓	Skate (%) ↓	Dyn. Stability (%) ↑	BoS Dist (cm) ↓
HUMOS	1.23	1.04	7.37	71.9	14.62
HUMOS - \mathcal{L}_E	1.20	0.98	9.3	71.0	15.01
HUMOS - \mathcal{L}_{KL}	1.14	0.93	6.96	71.05	15.21

E AMASS Shape Statistics

For training and evaluation, we use the AMASS dataset [49]. AMASS is a comprehensive collection of human motion data, unifying various optical marker-based motion capture datasets. This dataset stands out due to its extensive volume, containing over 50 hours of motion data from 480 unique subjects, encompassing more than 11,000 distinct motions. Among the 480 unique subjects,

we have 274 male and 206 female subjects. To understand the diversity of body shapes included in AMASS, in Fig. S.3, we plot the mean and standard deviation of each principal component for the AMASS beta parameters. Following prior work, we use the first 10 shape principal components to represent body shape.

F Additional Implementation Details.

Data processing. AMASS captures diverse human motions performed by real participants. Therefore, motions in AMASS start at arbitrary locations and facing directions. AMASS also includes motions where the person is supported by objects such as chairs, stairs or raised platforms. Only the human is captured in such sequences, and given the lack of a supporting object, these motions are physically implausible. These sequence, along with arbitrary start locations and facing directions, add unnecessary ambiguity and make the raw AMASS data unsuitable for neural network training. To prevent this, we process the raw AMASS data by removing all sequences where the lowest vertex in at least 5 frames is higher than 0.25m from the ground. Next, as described in the main text, we canonicalize all sequences to start at the origin with the same facing direction. To augment our training data, we mirror the pose parameters and global root translation from left-to-right and vice-versa, effectively doubling the training data. Figure S.4 shows the effect of each step in our data processing pipeline.

Motion representation. The SMPL body model parameterizes the human body into body pose, shape and global root translation. The SMPL body pose is represented as parent-relative rotations in the axis-angle format. For our motion representation, we follow NeMF [35] and convert the parent-relative joint rotations to global root-relative rotations in 6d format [86]. This helps with convergence and produces better performance than using the SMPL parameters directly. We also experiment with using deltas in joint rotations and global translation in our motion representation. We empirically observe worse performance in this setting due to the propagation of errors in the integration step when recovering the joint rotations and translation from the predicted deltas.

F.1 Perceptual Study

We show the layout for our perceptual study in Fig. S.5. We randomly sample sequence generations from our methods and baselines and every video is rated by 25 participants on Amazon Mechanical Turk. We ensure quality in ratings by adding two ground-truth videos and two catch-trial videos per worker with extreme ground penetrations or floating sequences. Additionally, every participant is shown 5 *warming-up* sequences at the start of their annotation task which we discard. This allows the participant to get a sense of the task before they can reliably rate the generated motions. We report average ratings across all participants who qualify the quality checks.

To test statistical significance, we performed one-way ANOVA tests, yielding a significant p-value of $1.5 \times e^{-10}$. Tukey's HSD statistical test indicates that our method has statistically significant differences with TEMOS-Rokoko (mean diff = 0.386, $p < 0.001$) and TEMOS-Rokoko-G (mean diff = 0.447, $p < 0.001$).



Fig.S.2: Additional qualitative comparison of shape-conditioned motion generation. Each row represents generations across different methods for a unique body shape and gender. The difference in quality between methods is particularly evident in their interaction with the ground. **Q Zoom in.**

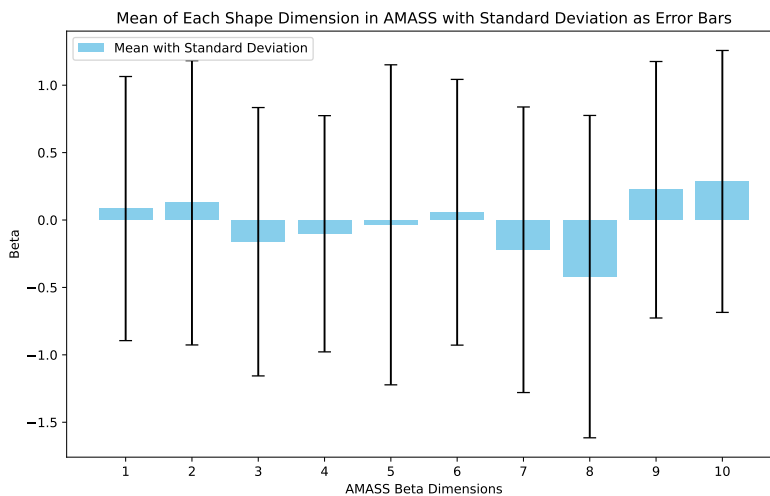


Fig. S.3: Mean and standard deviation of the first 10 betas parameters in AMASS. This represents the diversity in body shapes.

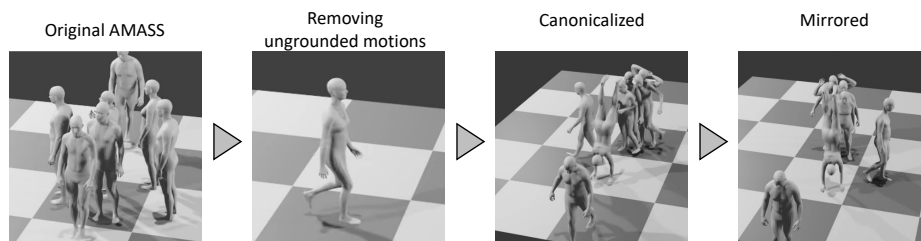
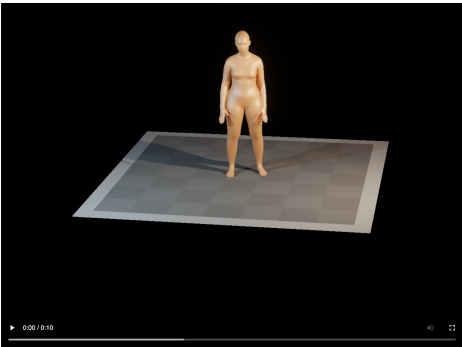


Fig. S.4: We process the raw data from AMASS by 1) removing unsupported physically implausible motions *e.g.* walking up the stairs 2) canonicalizing all motions to start facing the same direction at origin and 3) mirroring the pose and root translation to augment data

How realistic is the motion given this body shape?

In this task, you will see an animation of a character performing various motions.
Your task is to rate how realistic the motion is based on the character's body shape and gender.
For example, a lean person moves differently than an obese person.
Please pay attention to the motion of the body and the contact with the ground. The motion should look plausible and the body should be realistically moving, without penetrating the floor, floating, or sliding on the ground.
The hands are not animated, please ignore them in your rating.
You must watch every video until the end to be able to rate it.
Once you have rated all videos, the 'SUBMIT' button will be activated and you can submit the HIT.



How realistic is the motion given this body shape?

1 Completely unrealistic 2 Somewhat unrealistic 3 Neither realistic nor unrealistic 4 Somewhat realistic 5 Completely realistic

Next Video
0:00/0:10

3 videos left

Fig. S.5: Layout of the perceptual study.