# HACK: Learning a Parametric Head and Neck Model for High-fidelity Animation

LONGWEN ZHANG*, ShanghaiTech University and Deemos Technology Co., Ltd., China
ZIJUN ZHAO*, ShanghaiTech University and Deemos Technology Co., Ltd., China
XINZHOU CONG*, ShanghaiTech University and Deemos Technology Co., Ltd., China
QIXUAN ZHANG, ShanghaiTech University and Deemos Technology Co., Ltd., China
SHUQI GU, ShanghaiTech University, China
YUCHONG GAO, ShanghaiTech University, China
RUI ZHENG, ShanghaiTech University, China
WEI YANG, Huazhong University of Science and Technology, China
LAN XU[†], ShanghaiTech University, China
JINGYI YU[†], ShanghaiTech University, China

Fig. 1. we present HACK (Head-And-neCK), a novel parametric model for constructing the cervical region of digital humans. By combining rich physically-based appearance modeling and inner anatomical structures, HACK achieves more accurate and expressive results than existing head and neck models.

Significant advancements have been made in developing parametric models for digital humans, with various approaches concentrating on parts such as the human body, hand, or face. Nevertheless, connectors such as the neck have been overlooked in these models, with rich anatomical priors often unutilized. In this paper, we introduce HACK (Head-And-neCK), a novel parametric model for constructing the head and cervical region of digital humans. Our model seeks to disentangle the full spectrum of neck and larynx motions, facial expressions, and appearance variations, providing personalized and anatomically consistent controls, particularly for the neck regions. To build our HACK model, we acquire a comprehensive multi-modal dataset of the head and neck under various facial expressions. We employ a 3D ultrasound imaging scheme to extract the inner biomechanical structures, namely the precise 3D rotation information of the seven vertebrae of the cervical spine. We then adopt a multi-view photometric approach to capture the geometry and physically-based textures of diverse subjects, who exhibit a diverse range of static expressions as well as sequential head-and-neck movements. Using the multi-modal dataset, we train the parametric HACK model by separating the 3D head and neck depiction into various shape, pose, expression, and larynx blendshapes from the neutral expression and the rest skeletal pose. We adopt an anatomically-consistent skeletal design for the cervical region, and the expression is linked to facial action units for artist-friendly controls. We also propose to optimize the mapping from the identical shape space to the PCA spaces of personalized blendshapes to augment the pose and expression blendshapes, providing personalized properties within

*Equal contributions.
†Corresponding author.

Authors' addresses: Longwen Zhang, ShanghaiTech University and Deemos Technology Co., Ltd., Shanghai, China, zhanglw2@shanghaitech.edu.cn; Zijun Zhao, ShanghaiTech University and Deemos Technology Co., Ltd., Shanghai, China, zhaozj2022@shanghaitech.edu.cn; Xinzhou Cong, ShanghaiTech University and Deemos Technology Co., Ltd., Shanghai, China, congxzh2022@shanghaitech.edu.cn; Qixuan Zhang, ShanghaiTech University and Deemos Technology Co., Ltd., Shanghai, China, zhangqx1@shanghaitech.edu.cn; Shuqi Gu, ShanghaiTech University, Shanghai, China, gushq@shanghaitech.edu.cn; Yuchong Gao, ShanghaiTech University, Shanghai, China, gaoych@shanghaitech.edu.cn; Rui Zheng, ShanghaiTech University, Shanghai, China, zhengrui@shanghaitech.edu.cn; Wei Yang, Huazhong University of Science and Technology, Wuhan, China, weiyangcs@hust.edu.cn; Lan Xu, ShanghaiTech University, Shanghai, China, xulan1@shanghaitech.edu.cn; Jingyi Yu, ShanghaiTech University, Shanghai, China, yujingyi@shanghaitech.edu.cn.

the framework of the generic model. Furthermore, we use larynx blendshapes to accurately control the larynx deformation and force the larynx slicing motions along the vertical direction in the UV-space for precise modeling of the larynx beneath the neck skin. HACK addresses the head and neck as a unified entity, offering more accurate and expressive controls, with a new level of realism, particularly for the neck regions. This approach has significant benefits for numerous applications, including geometric fitting and animation, and enables inter-correlation analysis between head and neck for fine-grained motion synthesis and transfer.

CCS Concepts: • **Computing methodologies** → **Mesh models**.

Additional Key Words and Phrases: Head and neck modeling, Anatomical model, Facial expressions, Neck animation, Parametric learning

## 1 INTRODUCTION

Realism in human modeling goes beyond surface-level characteristics. Anatomical structures, physical attributes, and their induced motions capture unique and nuanced personalities, encompassing their inner thoughts, emotions, and experiences. In constructing digital humans, tremendous efforts have been focused on modeling the skin, skeleton, and muscle components of faces [Kähler et al. 2001; Qiu et al. 2022], hands [Li et al. 2021; Romero et al. 2017], bodies [Loper et al. 2015] that exhibit more obvious features. In contrast, body parts that serve as connectors have been largely overlooked. A significant example is the human neck which connects the jaw, head, and shoulder, allowing the head and face to move in a more natural way. The neck is often visible in a variety of poses and angles, and it greatly impacts the overall appearance of a digital human character. For example, the prominent jugular vein on the neck of Michelangelo's monumental David sculpture vividly illustrates the adrenaline-fueled situation of David, before facing the giant Goliath.

Convincing neck movements further add subtlety and realism to facial performances as well as enable nonverbal communication and social interaction. For instance, we humans turn to an unconscious habit of air swallowing when overly anxious. In fact, the idiosyncrasies of the head and neck are the defining characteristics of various human characters, real or virtual, physical or digital. However, so far very few attempts have been made to model human necks [Liu et al. 2021; Luo et al. 2013], let alone a generic parametric model for the face and neck as a whole. The challenges are multifold. Anatomically the neck is made up of bones, muscles, and other tissues that are difficult to model. Therefore, to allow for a wide range of movement including flexion, extension, rotation, and lateral bending, it is critical to accurately model the joints in accord with the head. Further, the neck is a flexible and deformable structure under facial motions. Yet by far very limited datasets are available that simultaneously capture the facial and neck geometry under a rich variety of motions.

The majority of existing parametric models build upon 3DMM [Blanz and Vetter 1999; Booth et al. 2017; Brunton et al. 2014] to model the human head [Qiu et al. 2022; Wang et al. 2022b; Yang et al. 2020], as large as the whole front face and as small as individual facial components such as eyeballs, teeth, lips, etc [Bérard et al. 2016; Garrido et al. 2016; Wu et al. 2016]. A few attempts [Li et al. 2020a; Liu et al. 2021] aim to extend the parametric model to the neck region. They unanimously employ a highly simplified model, e.g., using a single

empirical joint, to represent the kinematics between head and neck, without considering more sophisticated anatomical priors of inner structures such as the cervical spine. Physiologically, the animation of human head and neck is an intricate orchestration, from rich facial expressions coordinated with neck movements to nuanced deformation of the neck shape caused by sliding larynx beneath the skin. Without accurately incorporating these anatomical priors, brute-force extensions of the parametric model can often lead to unconvincing and sometimes biomechanically absurd results as shown in Fig. 15.

Early 3DMM-based models [Li et al. 2020a; Yang et al. 2020] also tend to retain the statistical properties of the underlying 3D scan dataset based on an inherent low-rank approximation and therefore are insufficient to represent local and high-frequency surface details. More recent data-driven approaches [Giebenhain et al. 2022; Hong et al. 2022; Wang et al. 2022a; Yenamandra et al. 2021; Zhuang et al. 2022] adopt neural rendering techniques to provide more personalized and realistic results across various identities. A drawback though is that the latest neural models cannot readily support the existing CG production pipeline to conduct either rendering or controllable editing. The most recent trend [Cao et al. 2022; Gao et al. 2022b; Li et al. 2020b,a] is to explicitly employ more personalized characteristics into the generic parametric models. However, they rely on tedious subject-specific training or network inference to obtain personalized facial assets, sacrificing the compact and efficient controls of 3DMM-based models.

In this paper, we present HACK (Head-And-neCK), a novel parametric model for constructing the cervical region of digital humans. By combining the outer physically-based appearance and inner anatomical structures (i.e., the cervical spine that is composed of seven vertebrae), HACK tackles the full spectrum of neck and larynx motions, offering more personalized and anatomically-consistent controls with a new level of realism (see Fig. 1). As a parametric model analogous to previous blendshape-based techniques [Li et al. 2017; Loper et al. 2015], HACK is differentiable, computationally efficient, and compatible with existing CG engines.

The first step in building HACK is data collection for effective model training. We first acquire a comprehensive multi-modal dataset that covers both internal anatomical structures and external appearances of the head and neck under facial expressions. To extract the cervical spline as biomechanical priors, we use the portable 3D ultrasound imaging (US) system [Chen et al. 2020] to obtain the sonography scans of the neck regions of individuals. Such a solution is radiation-free and cost-effective and the process has obtained IRB approval. We then ask experienced radiologists to label 3D landmarks on the seven vertebrae of the cervical spine and subsequently extract their anatomically-consistent rotation information with regard to the skull and the external neck geometry. To correlate the inner structures of the neck with its outer appearance, we conduct physically-based scanning using a photometric scanning solution [Debevec et al. 2000; Debevec 2012; Zhang et al. 2022]. Specifically, we employ the multi-view photometric capture system to capture both the geometry and physically-based textures of head and neck regions. We further conduct topology-consistent reconstructions to capture the physically-based head-and-neck attributes on subjects who perform a diverse set of static expressions as well

as sequential motions. Large-scale geometry variations can be effectively recovered in terms of 3D surfaces via 3D/4D scans. Besides, we directly obtain the high-resolution 2D RGB images and normal maps to model small-scale and nuanced variations (e.g., movements of the larynx).

To learn the parametric HACK model, we follow a similar discipline as in the human face and body modeling [Li et al. 2017; Loper et al. 2015] by separating the depiction of 3D head and neck into various shape, pose, expression, and larynx blendshapes from the neutral expression and the rest skeletal pose. Specifically, to learn pose-dependent blendshapes, we adopt an anatomy-consistent skeleton of the neck that consists of 8 joints corresponding to the 7 cervical vertebrae and the head skull, respectively. We also tie the expression blendshapes to the action units of the Facial Action Coding System (FACS) [Li et al. 2010, 2020b; Prince et al. 2015] to provide a compact and artist-friendly expression control. Instead of using the same generic blendshapes across identities similar to previous parametric models [Li et al. 2017; Loper et al. 2015], we further tailor schemes to learn personalized pose and expression blendshapes, by optimizing the generic mapping from the identical shape space to personalized blendshapes. We show such a strategy significantly enhances HACK's capability of modeling personalized properties while maintaining the generalization across identities as a generic model. Inspired by the recent work [Liu et al. 2021], we further adopt the larynx blendshapes to control the larynx deformation on top of the larynx-removed neck under the rest pose, and subsequently force the larynx slicing motions along the vertical direction in the UV-space. Such a strategy provides the disentanglement to anatomically mimic two kinds of muscles related to the larynx: one moves the vocal folds and hence changes the larynx's size while the other causes vertical slicing of the larynx beneath the neck skin. Once trained, our HACK model is differentiable and compatible with standard CG software, and hence readily benefits a variety of applications like geometric fitting, animation, and visual inference. Most importantly, our HACK model provides more accurate and expressive controls for spine-driven neck poses and larynx motions. It also enables fine-grained analysis of the inter-correlation of head and neck, including: (1) Faithful motion synthesis by leaning the temporal mapping from the head pose to the skeletal pose of the cervical spine, from facial expression to the larynx slicing. (2) Biology-consistently transferring the head/neck motions from a human to another mammal i.e., the giraffe, since through evolution most mammals share the same skeletal structure of cervical spines. To summarize, our main contributions include:

- We present HACK, a generic parametric model that jointly considers human identity, facial expression, anatomy-inspired neck, and larynx motions, as well as physically-based appearance.
- We propose to jointly utilize the comprehensive inner biomechanical priors and external appearances during parameter learning, which provides personalized and anatomically-consistent controls for the neck regions.
- We make available our trained HACK model and showcase various applications to demonstrate its effectiveness, ranging

from model fitting and inference, to motion synthesis and transfer.

## 2 RELATED WORKS

*Parametric Head and Neck Models.* Parametric human modeling is characterized by a fixed number of parameters and a specific functional form, which makes modeling of the human body [Anguelov et al. 2005; Joo et al. 2018; Loper et al. 2015], hands [Romero et al. 2017], heads [Li et al. 2017; Paysan et al. 2009; Ploumpis et al. 2019], and faces [Abrevaya et al. 2018; Booth et al. 2017; Brunton et al. 2014; Dai et al. 2020; Huber et al. 2016; Paysan et al. 2009; Smith et al. 2020; Wang et al. 2022b] relatively simple to work with and interpret. Modeling and analysis of human faces are particularly important in the field as it plays a key role in many applications such as facial animation, virtual try-on, and face editing. Since Blanz and Vetter [1999] propose the first key ideas of general face representation, i.e., linear combinations of faces produce morphologically realistic faces, and separating facial shape and color to disentangle illumination, camera parameters and etc, many advances have been emerging [Kim et al. 2018; Li et al. 2017; Paysan et al. 2009]. We refer the readers to Egger et al. [2020] for a comprehensive survey on 3DMM. Though achieving high quality on the face region, most existing approaches pay little attention to the neck, which results in unrealistic neck models and animations. Some exceptions [Choi et al. 2022; Danecek et al. 2022; Li et al. 2017; Loper et al. 2015; Osman et al. 2020, 2022] incorporate the neck modeling, which is still relatively simple and exhibit unrealistic deformations. Especially, Liu et al. [2021] and Li et al. [2013] take the larynx into account and map 2D UV texture to the vertex displacement on 3D meshes, which produce good results.

A recent trend in human body modeling is to enforce anatomical constraints, several exemplary works include the modeling of the full human body [Xu et al. 2020], head [Duan et al. 2015, 2013; Qiu et al. 2022], and hand models [Li et al. 2021]. However, anatomical constraints of the neck region are rarely adopted in generic head and neck modeling [Lee et al. 2009; Lee and Terzopoulos 2006; Luo et al. 2013]. The reason is twofold: 1. anatomical constraints are good for the physical simulation of a specific subject but are difficult to generalize; 2. they lack details of external neck skin deformation. In this paper, we adequately consider anatomical constraints and construct a head and neck model with finer neck deformation.

Generating a realistic and convincing person-specific head-and-neck model requires tedious steps of manual intervention, including the construction of physically-based appearance attributes, the capture of fine facial movements and their migration [Laine et al. 2017; Li et al. 2020a; Thies et al. 2016]. These tasks generally require a lot of manual adjustments by professional artists, making the whole process too expensive to scale to a larger audience. Recently many learning-based approaches enable generic models that are highly robust for different lighting and motion [Bao et al. 2021; Feng et al. 2021; Li et al. 2017; Raj et al. 2021] or performer-specific high-quality geometry and appearance [Gafni et al. 2021; Gecer et al. 2019; Lattas et al. 2022; Lombardi et al. 2019, 2021] by processing datasets of different volumes. The former is able to learn person-specific blendshapes and pore-level dynamic physical materials [Li et al. 2020b],
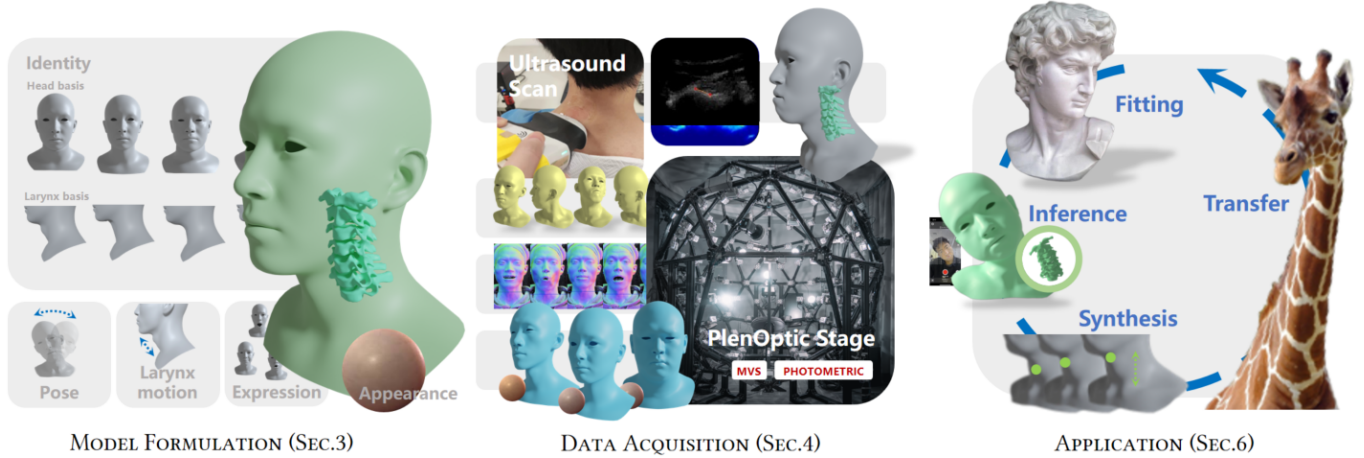
Fig. 2. Overview of HACK. We first introduce the model formulation of HACK in Sec. 3, and the comprehensive multi-modal dataset that covers both internal anatomical structures and external appearances under facial expressions in Sec. 4. After building HACK, various applications are shown in Sec. 6 to demonstrate its effectiveness.

but generic blendshapes tend to lead to the loss of fine-grained expressions and dynamic physically-based textures. The latter can achieve very fine geometry and appearance for specific performers, and can even be re-lighted for animation and cross-performer driving. However, they mostly lack explicit control support and are difficult to deploy into traditional production processes.

*Facial Animation.* The goal of 3D facial animation is to make the character's facial expressions and movements as realistic and believable as possible as if the character is a real human or creature. Typically, a 3D mesh is attached to a skeleton (or bones) and moves the bones to deform the mesh [Wang et al. 2017]. Instead, muscle-based animation uses physics-based models to simulate the movement and deformation of muscles in a 3D mesh. One example is the Facial Action Coding System (FACS), which annotates and codes facial expressions based on the movement of specific facial muscles [Ekman and Friesen 2002; Jackson et al. 2017; Mattar and Gao 2017; Mattar et al. 2015]. Another example is the Facewarehouse model [Deng et al. 2019; Yang et al. 2020], which represents the shape and appearance of a face using a set of blendshapes and a physics-based muscle simulation model. Additionally, a face model can be driven by audio signals or video footage of a person's face. The Audio-Driven Facial Animation (ADFA) framework [Liu et al. 2019] uses deep learning to map between audio features and facial expression parameters, and the Audio-Driven Facial Animation System (ADFAS) [Fan et al. 2019] combines rule-based and data-driven techniques to synthesize realistic facial animations from audio signals. 2D Video-Driven Facial Animation (VDFA) systems combine optical flow and deep learning to synthesize facial animations from video input [Guo et al. 2018], while 3D VDFA systems use morphable model fitting and facial expression synthesis to generate realistic facial animations [Wang et al. 2020b].

*Geometry and Appearance.* Capturing fine facial movements is an essential setup of head and neck modeling, and several advanced acquisition systems have already been proposed. Early approaches

adopt single acquisition devices, such as laser scanners [Blanz and Vetter 2003; Levoy et al. 2001; Phillips et al. 2008] and structured light scanners [Geng 2011]. Structured light systems capture at a single viewpoint and may not be sufficient to reconstruct a 3D geometry of the human face accurately. Multi-view stereo (MVS) can overcome this limitation by reconstructing a 3D geometry from 2D images of multiple views [Goesele et al. 2006; Seitz et al. 2006; Wrobel 2001], where several successful works produce good results in recovering the human body [Dou et al. 2016; Zhang et al. 2022], face [Chen et al. 2019; Klaudiny and Hilton 2012; Lombardi et al. 2018, 2021; Smith et al. 2020], hand [Romero et al. 2017], hair [Nam et al. 2019], etc.

Photometric stereo is another prevalent technique for facial data acquisition [Jain et al. 2017; Ma et al. 2006]. Instead of locating 3D points from pixel correspondences across multi-view images, photometric stereo systems capture images from single or sparse viewpoints under various lighting conditions and use the intensity variations to estimate the surface normals [Villarini et al. 2017; Wang et al. 2020a; Zafeiriou et al. 2013]. Photometric stereo can be used to reconstruct the surface detail of a face, including wrinkles, texture, and color. The most well-known facial capture system using photometric stereo is the LightStage [Debevec et al. 2000], which has been widely used and modified in a large variety of fields in recent years [Debevec 2012; Dutta 2010; Wenger et al. 2005; Weyrich et al. 2006]. The FaStage proposed in [Zhang et al. 2022] extended the LightStage by combining multi-view reconstruction and photometric reconstruction to recover the dynamic geometry and physically-based texture of the performer in a motion sequence.

*Anatomical Data Acquisition.* Exploiting anatomically prior information for human modeling has been a widely researched area, with great success being achieved through the use of Magnetic Resonance Imaging (MRI) and Computer Tomography (CT) scans. These techniques are widely used to reconstruct 3D models for evaluating bone and joint morphology and preoperative planning [Stephen et al. 2021; Touati et al. 2021]. CT scans provide a clear bone-soft

tissue contrast but expose the subject to high doses of ionizing radiation. On the other hand, MRI is radiation-free, but it is more often used to scan soft tissue [Iacono et al. 2015; Misaki et al. 2015] rather than bone structures [Samim 2021].

Ultrasound technology then is an alternative technique for capturing human bones in various applications [Nelson and Pretorius 1998]. Wang et al. [2018] used ultrasound range finding to capture detailed 3D models of the neck. van Eerd et al. [2014] scanned cervical spine specimens using an ultrasound probe that moves on a fixed track. In other studies, ultrasound waves have been used to track and analyze neck movements in various tasks such as head turning and tilting [Zhang et al. 2018, 2019]. Recently, Chen et al. [2020] proposed a portable mobile 3D ultrasound scanning system, which we adopt to capture and reconstruct the interior structure of the cervical spine for anatomically prior information.

## 3 MODEL FORMULATION

Here we introduce a novel parametric model, HACK, for jointly constructing the head and neck region of digital humans. As shown in Fig. 2, our HACK model combines the rich observations from both physically-based appearance and inner anatomical structures. It achieves full-spectrum modeling of neck and larynx motions, providing personalized and anatomically-consistent controls. Analogous to the human head and body model [Li et al. 2017; Loper et al. 2015], our HACK models the depiction of 3D head and neck into various shape, pose, expression and larynx blendshapes from the neutral expression and the rest skeletal pose. Specifically, the general formulation of HACK is defined as follows:

$$\text{HACK}(\boldsymbol{\beta}, \boldsymbol{\psi}, \boldsymbol{\theta}, \eta, \tau, \boldsymbol{\alpha}) = \{\text{G}(\boldsymbol{\beta}, \boldsymbol{\psi}, \boldsymbol{\theta}, \eta, \tau), \text{A}(\boldsymbol{\alpha})\}, \quad (1)$$

where G denotes the head-and-neck geometry, and A produces the physically-based appearance. $\boldsymbol{\beta}$, $\boldsymbol{\theta}$, $\boldsymbol{\psi}$ and $\boldsymbol{\alpha}$ are parameters to control the shape, pose, expression and appearance, respectively. Besides, we use $\eta$ and $\tau$ to control the larynx size and larynx slicing beneath the neck skin, so as to separate larynx motions from facial expressions for more subtle and realistic modeling (e.g., for handling swallowing). Specifically, the geometry model is obtained through a skinning process as follows:

$$\text{G}(\boldsymbol{\beta}, \boldsymbol{\psi}, \boldsymbol{\theta}, \eta, \tau) = \text{LBS}\big(T(\boldsymbol{\beta}, \boldsymbol{\psi}, \boldsymbol{\theta}, \eta, \tau), J(\boldsymbol{\beta}), \boldsymbol{\theta}, \mathcal{W}\big), \quad (2)$$

where $\text{LBS}(\cdot)$ denotes the Linear Blend Skinning (LBS) function; $T(\cdot)$ is the person-specific mesh with corrective deformations ascribed to identity, pose, expression, and larynx parameters; $J(\boldsymbol{\beta})$ represents the joint location; $\mathcal{W}$ is the learned skinning weight of $\text{LBS}(\cdot)$. In stark contrast, we adopt an anatomy-consistent skeleton that consists of 8 joints corresponding to the 7 cervical vertebrae and the head skull. Hence, the joint position regressor $J(\cdot)$ infers the 8 person-specific biomechanical joint positions given identity parameter $\beta$ for accurate skinning (Sec. 3.1).

Then, the personalized template under rest pose is a linear combination of the universal template and various blendshapes:

$$
\begin{aligned}
T(\boldsymbol{\beta}, \boldsymbol{\psi}, \boldsymbol{\theta}, \eta, \tau) = & \bar{T} + \text{B}_\text{S}(\boldsymbol{\beta}; \mathcal{S}) + \text{B}_\text{E}(\boldsymbol{\psi}; \mathcal{E}_{\boldsymbol{\beta}}) + \text{B}_\text{P}(\boldsymbol{\theta}; \mathcal{P}_{\boldsymbol{\beta}}) \\
& + L(\boldsymbol{\beta}, \eta, \tau; \mathcal{L}).
\end{aligned}
\quad (3)
$$

Note that the universal template $\bar{T} \in \mathbb{R}^{3N}$ is under the rest pose with mean shape, no expression, and larynx-removed. $\text{B}_\text{S}$ is the
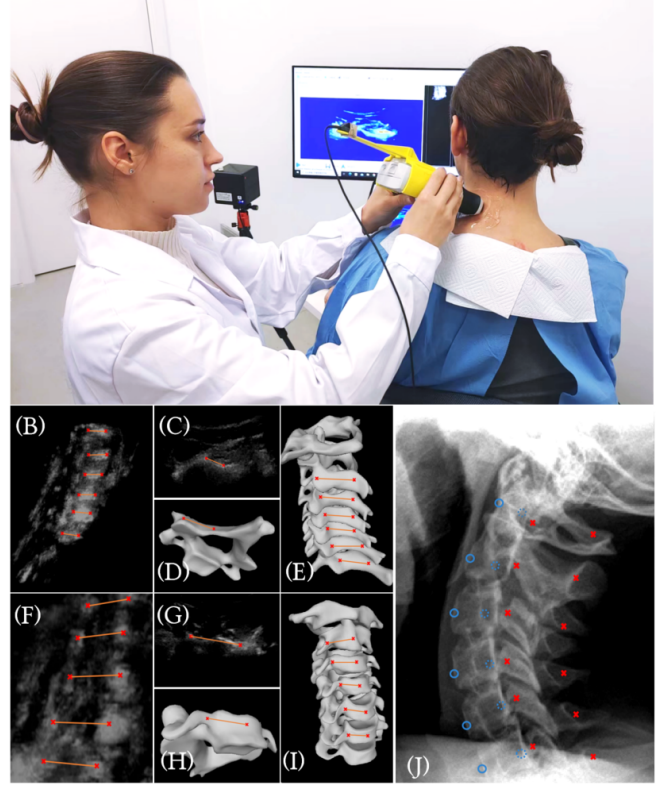


Fig. 3. Ultrasound data acquisition pipeline and results. (A) A performer was scanned using the portable real-time 3D ultrasound imaging system. (B) Reconstructed image from rear-left side scan and feature marker. (C) 2D slice of rear-left side scan with feature marker. (D) Single vertebra mesh corresponding to rear-left side US scan. (E) Full cervical spine mesh corresponding to rear-left side US scan. (F) Reconstructed image from front-right side scan and feature marker. (G) 2D slice of front-right side scan with feature marker. (H) Single vertebra mesh corresponding to front-right side US scan. (I) Full cervical spine mesh observed from the front-right side. (J) Feature marks on the sagittal human neck X-ray image, blue markers and red markers correspond to the front-right scan and rear-left scan, respectively.

multiplication of orthonormal PCA of shape blendshapes $\mathcal{S}$ and $\boldsymbol{\beta}$. $\text{B}_\text{E}$ is the multiplication of expression blenshapes $\mathcal{E}_{\boldsymbol{\beta}}$ and $\boldsymbol{\psi}$. $\text{B}_\text{P}$ is the multiplication of pose blendshapes $\mathcal{P}_{\boldsymbol{\beta}}$ and pose rotation matrix. Please refer to previous work [Li et al. 2017; Loper et al. 2015] for detailed formulation. Yet, such unified blendshapes turn to retain the statistics of the whole dataset, thus losing personalized details. To this end, we propose to model the PCA spaces of the various person-specific $\mathcal{E}_{\boldsymbol{\beta}}$ and $\mathcal{P}_{\boldsymbol{\beta}}$, respectively, and subsequently learn an efficient mapping from the identity parameter $\boldsymbol{\beta}$ to their corresponding PCA spaces, so as to obtain person-specific facial motion traits (Sec. 3.2).

We further adopt the larynx deformation $L(\boldsymbol{\beta}, \eta, \tau; \mathcal{L})$ following previous work [Liu et al. 2021], which will be explained in Sec. 3.3. Specifically, we first predict the larynx geometry on top of the larynx-removed neck by controlling its size according to $\eta$. Then, we formulate the larynx slicing motions along the vertical direction in the UV-space in terms of $\tau$.

## 3.1 Anatomy-aware Skeleton in HACK

An over-simplified skeleton, e.g., using a single empirical joint, to represent the kinematics between head and neck is insufficient and often leads to unconvincing animations. To this end, in our HACK model, we design an anatomy-inspired skeletal structure with 8 joints, which correspond to 7 cervical vertebrae (C1-C7) and the head skull (the apex of C1), respectively. These joints, denoted as c7-t1, c6-c7, c5-c6, c4-c5, c3-c4, c2-c3, c1-c2, and o-c1, are approximately located behind the intervertebral discs, and provide accurate kinematic control for nuanced neck motions.

To accurately determine the set of person-specific joint positions $J \in \mathbb{R}^{3K}$, we resort to a portable 3D ultrasound imaging (US) system to obtain sonography scans with 3D information of the seven vertebrae, labeled by experienced radiologists (see Sec.4.1). Accompanied by our anatomy-consistent skeleton design, we then optimize the regressor $J(\cdot)$ to predict joint positions from the rich ultrasound data in the rest pose. Compared with the existing one-joint solutions, our HACK provides more accurate control over head and neck poses by accurately locating rotation centers and fully modeling the rotation of each cervical vertebra.

## 3.2 More Personalized Expression/Pose Blendshapes

Recall that we follow the similar discipline of blendshapes as in human face and body modeling [Li et al. 2017; Loper et al. 2015]. Specifically, we adopt the orthonormal PCA of shape blendshapes, to correlate the shape parameter $\beta$ with rich identity information. Instead of using the same generic expression and pose blendshapes that remain for all individuals, we further tailor schemes to employ more personalized properties into these blendshapes. We propose to learn the extra mappings $\mathcal{M}_E$ and $\mathcal{M}_P$ from $\beta$ to expression and pose blendshapes basis $\mathcal{E}_\beta$, $\mathcal{P}_\beta$ for novel identities, i.e., $\mathcal{M}_E(\beta) \mapsto \mathcal{E}_\beta$ and $\mathcal{M}_P(\beta) \mapsto \mathcal{P}_\beta$.

For expression deformations, we tie the expression blendshapes to the action units of FACS [Prince et al. 2015] for more artist-friendly expression controls, following previous work [Li et al. 2020b]. Specifically, for each captured subject with shape parameter $\beta$, we obtain the corresponding expression blendshapes from the static captured scans with FACS expressions: $\mathcal{E}_\beta = [E_1^\beta, \ldots, E_{|\psi|}^\beta]$. Thus, all the individuals share the same latent structure for the expression parameters $\psi$. To efficiently learn the mapping $\mathcal{M}_E$, we first formulate the PCA space of expression blendshapes from the set $\{\mathcal{E}_\beta\}$ across various identities. We subsequently train the mapping network $\mathcal{M}_E$ as a shallow MLP to predict the PCA weights, hence generating the personalized expression blendshapes.

For pose-dependent deformations, similar to previous work [Loper et al. 2015], we employ the pose blendshapes $\mathcal{P}_\beta$ with the rotation matrix interpreted from the skeletal parameter $\theta$. Here we adopt an anatomy-aware skeletal design, where the $\theta$ denotes the concatenated rotation vector of the joints for 7 cervical vertebrae and the head skull. Again, we obtain the person-specific $\mathcal{P}_\beta$ for the performer with identity $\beta$ from the captured dynamic sequential scans, and subsequently optimize the PCA space of pose blendshapes from $\{\mathcal{P}_\beta\}$. Analogous to $\mathcal{M}_E$, we adopt the same mapping network with shallow MLP to predict the PCA weights from the identity

parameter. Such a strategy significantly improves the ability of our parametric model for generating more personalized controls.

## 3.3 Larynx Modeling

Accurate geometric modeling of the larynx part is vital for realistic head and neck presentation, especially for actions with subtle larynx movement (e.g., swallowing and talking). Anatomically, the movement of the larynx is controlled by two groups of muscles, one of which moves the vocal folds and changes the size of the larynx, and the other moves the position of the larynx in the neck vertically. Hence, we model the larynx as vertex displacements added to the larynx-removed rest-pose mesh with shape change capability and constrained movement along the vertical direction of the neck, which is also similar to [Liu et al. 2021]. $L(\beta, \eta, \tau; \mathcal{L}) : \mathbb{R}^{2+|\beta|} \mapsto \mathbb{R}^{3N}$ is the larynx shape function that maps identity $\beta$, larynx size $\eta$ and position $\tau$ to vertex displacements to simulate the sliding effect beneath neck skin.

$$L(\beta, \eta, \tau; \mathcal{L}) = \eta \cdot \sum_{i=1}^{|\beta|} \beta_i L_i(\tau), \qquad (4)$$

where $\mathcal{L} = [L_1(\tau), \ldots, L_{|\beta|}(\tau)] \in \mathbb{R}^{3N \times |\beta|}$ is the larynx blendshape basis that is able to move vertically according to $\tau$. In actual implementation, we model the larynx geometric displacement in the UV space for convenience. Specifically, $L_i(\tau) \in \mathbb{R}^{3 \times H \times W}$ is a larynx blendshape basis represented as a 2D image through texture atlasing and storing the vertex displacement into corresponding uv coordinate. Then the larynx function can be reformulated as: $L(\beta, \eta, \tau; \mathcal{L}) = \{v(u,v)\}$, $v(u,v) = \eta \cdot \sum_{i=1}^{|\beta|} L_i(u, v+\tau)$, where $v(u,v)$ denotes the displacement of vertex v ascribe to larynx, which maps to the value at $u,v$ in L. This formulation guarantees that the larynx moves vertically along the neck and hence increases robustness.

## 4 DATA ACQUISITION AND PROCESSING

In order to build HACK, a comprehensive and multi-modal dataset that covers both internal anatomical structures and the external appearance of the head and neck is essential for effective training. To acquire this dataset, we first introduce an ultrasound pipeline to construct biomechanical priors from sonography scans of the cervical spines (Sec. 4.1). Then, we utilize a multi-view photometric capture system to capture both the geometry and appearance of the head and neck across diverse facial expressions and sequential motions, where geometry variations are effectively recovered (Sec. 4.2).

## 4.1 Neck Ultrasound Scanning

To accurately extract the cervical spine as biomechanical priors, we employ a portable real-time 3D ultrasound imaging system [Chen et al. 2020] to acquire sonography scans of the interior structure of the neck region, whose process includes scanning, annotating, and alignment. It is worth noting that our scanning pipeline is radiation-free and cost-effective, and the process has been approved by the IRB, which is demonstrated in Fig. 3.

*Ultrasound scan process.* During the ultrasound data collection process, the subject sits with their back against the chair to minimize
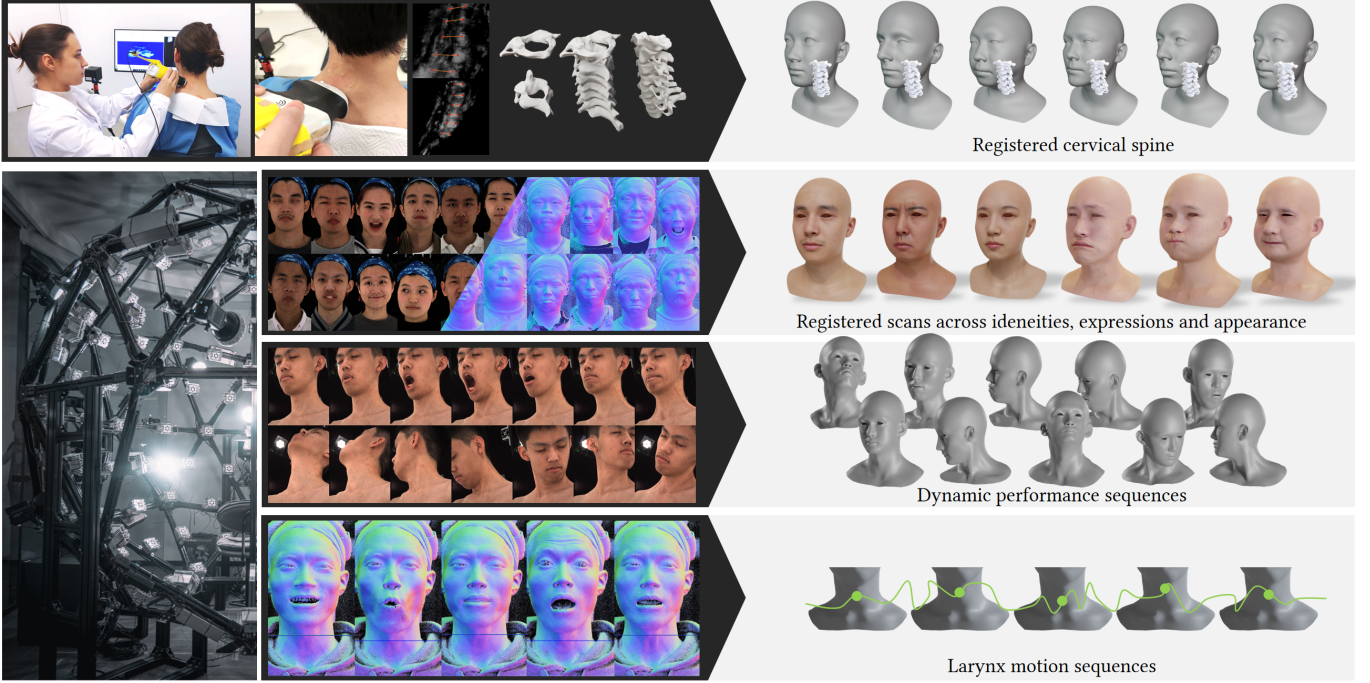
Fig. 4. Data processing pipeline. Utilizing both the ultrasound imaging system and multi-view photometric capture system, we further process these multi-modal data for model learning. The processed data includes registered cervical spine joints, registered neutral meshes, solved personalized expressions, dynamic performance sequences, larynx motion sequences, and physically-based textures.

body movement during scanning. An ultrasound probe, equipped with an electromagnetic (EM) tracking sensor, is pressed against the subject's neck and moved from C1 to C7 along the cervical curve, as illustrated in Fig. 3(A), or capturing continuous ultrasound images of the cervical spine. Before starting, an ultrasound coupling agent is applied to enable effective transmission of the ultrasound waves. We scan the neck regions of each subject in the following order: rear-left, rear-right, front-left, and front-right. Fig. 3(C) and Fig. 3(G) present ultrasound images of a neck in the neutral pose, scanned in the rear-left and front-right directions, respectively.

*Landmarks annotation.* To accurately locate the joints, we first reconstruct the 3D ultrasound volume of the neck and invite medical expertise to mark a set of pre-defined vertebra feature points. Specifically, the expertise annotates different points in ultrasound images captured from different neck regions. For instance, the points of interest of the rear-left side scans are the peak of the lamina and the bottom of the vertebra arch, while those of the front-right side scans are the end and center points of the vertebra body. Fig. 3(B, C), (F, G) shows the annotated feature points for the rear-left, and front-right side scans respectively. And Fig. 3(D, H, E, I) exhibit the corresponding feature points on the template cervical spine mesh.

*Neck template alignment.* The final stage is to deform a standard cervical spine template to match each identity using annotated feature points and then put the deformed cervical spine into the model in the rest-pose. We first resize the template cervical spine mesh to fit with the real-world size ultrasound scan, i.e., we scale

the vertebra in mesh to match the length in scans. The vertebra is rigid, and its motion and deformation could be approximated as the translation and rotation of the whole vertebra. We need to estimate the rotation and position of the vertebra for accurate alignment [Gao et al. 2022a] , and we define the position of a vertebra as the center of the feature point pair on it. Then, the rotation is calculated through the line that connects the feature point pair. Applying the translation and rotation to the scaled template mesh, we can align the template with ultrasound data accurately, as exhibited in Fig. 3(E, I, J) that the aligned template is consistent with the ultrasound scan. Further, we need to put the cervical spine template into the neck and head mesh in the rest pose. For this purpose, we ask the experienced radiologists to additionally annotate several markers on the outer skin surface on the ultrasound images, and subsequentially align the rest pose mesh with the neck template.

*Anatomical-prior head and neck skeleton.* With the aligned cervical spine template in the rest pose mesh, we can finally define the joints of our head and neck skeleton model. We design the HACK to have $K = 8$ joints corresponding to the bottom points of 7 vertebrae (C1-C7) and the apex of C1. This setup yields 8 bone transformations denoted as c7-t1, c6-c7, c5-c6, c4-c5, c3-c4, c2-c3, c1-c2, and o-c1, that will serve as the foundation for skinning as in Eqn. 2. By using sonography scans of the neck regions, we build biomechanical priors for HACK, which play an important role in further parametric learning and realistic animation.

## 4.2 Multi-view Photometric Scanning

To further correlate the inner structures of the neck with its outer appearance, and to achieve a diverse range of blendshapes for identity, expression, and pose, we conduct physically-based scanning using a photometric scanning solution. Specifically, we use a multi-view photometric capture system to acquire both the geometry and physically-based textures of the head and neck regions, ranging from a diverse set of static expressions to dynamic performance sequences. This allows us to effectively recover large-scale geometry variations and obtain detailed deformation of the head and neck, which is essential for building HACK and achieving realistic and nuanced movements.

*Static head and neck capture.* To learn the shape and expression space, we utilize the multi-view photometric capture system to capture the static head and neck of various identities in rest pose with a neutral expression, including 274 females and 350 males, ranging from age 15 to 65, with different races, skin tones, and face shapes.

Besides neutral faces in rest pose, we further capture scans of subjects with various expressions (28 predefined expressions following the FACS standard [Prince et al. 2015]) for creating personalized expression blendshapes for each identity. We conduct multi-view stereo reconstruction and register a template head and neck mesh with the reconstructed geometry following the registration pipeline as described in Li et al. [2020b], under the supervision of densely painted markers on faces. Then, we invite several professional artists to remove the larynx from the registered template mesh for explicitly estimating the larynx shape and the subsequential blendshapes space construction.Additionally, we use the captured photometric data to recover the high-resolution physically-based textures pertaining to reflectance, including the diffuse, specularity, and normal maps, which are further used to build up the appearance space.

*Dynamic performance sequence capture.* To analyze the geometric deformation of head and neck-related poses, we capture the performing data under continuously changing poses of various identities. In each captured clip, we require the performer to perform the following actions that will affect the head and neck muscles sequentially:

- rest pose, neutral expression;
- rest pose, swallow once;
- rest pose, open and close the mouth once;
- rest pose, stretch the mouth and contract the platysma muscle once;
- flex, extend and rotate the neck to extreme positions, neutral expression.
- cervical extension pose, swallow once;
- cervical extension pose, open and close the mouth once;
- cervical extension pose, stretch mouth and contract the platysma muscle once;
- side-bend cervical on both sides, neutral expression;
- neck-shifting (twisting neck), neutral expression;

The above actions are deliberately designed for comprehensive disclosure of personalized head and neck motion space together with

larynx motion. Similarly, we register the template mesh with recovered dynamic sequences and use the registered sequences for learning pose blendshapes and larynx deformation.

*Larynx motion during speaking.* Besides actions like mouth stretching and swallowing, another type of action that affects larynx motion greatly is speaking. To further analyze the speaking-related larynx motion and nuanced variations, we introduce an efficient way to capture larynx and speaking-caused mouth motion simultaneously to establish their connections. Instead of relying on complicated 3D reconstructions, we analyze the larynx and mouth movements on RGB images and normal maps estimated via photometric stereo. Specifically, we efficiently track the larynx in normal maps while extracting mouth movements using off-the-shelf expression identification techniques (please see Sec. 6.1 for details).

## 4.3 Data Finalization

With ultrasound scans, aligned head and neck geometries, and physically-based textures, we can finalize our data. As illustrated in Fig. 4, our comprehensive multi-modal dataset captures the full spectrum of head and neck and consists of 624 identities with 16078 mesh registrations. From ultrasound scans and static geometry scans of the same subject, we obtain the joint positions of the head and neck skeleton at rest-pose $\mathcal{J}_U^r$ and skeleton-aligned head and neck meshes $\mathcal{N}_U^r$ from 30 identities, where the symbol $\mathcal{N}$ means meshes recovered under neutral expression, the superscript $r$ denotes rest-pose, while subscript $U$ means the ultrasound data. To cover more varieties of identities, expressions, and poses, we aggregate samples from ICT-FaceKit with our static captures at the rest pose, and obtain $\mathcal{N}^r = \{\mathcal{N}_H^r, \mathcal{N}_I^r\} \in \mathbb{R}^{3N \times (N_H^r + N_I^r)}$, with larynx removed, where subscript H means HACK data and I means ICT-FaceKit data, and specifically $N_H^r = 624$, $N_I^r = 600$. Notice that $\mathcal{N}_U^r$ is a subset of $\mathcal{N}_H^r$. We then subtract $\mathcal{N}_H^r$ from the original meshes without removing the larynx to obtain the pure larynx geometries $\Gamma_H^r$. Further, we have also captured subjects with various expressions in static scans, we denote this data as $\mathcal{T}_H^r$, where symbol $\mathcal{T}$ means meshes making expressions, $r$ denotes rest pose, and we have 208 such subjects. Notice we do not remove the larynx geometry from $\mathcal{T}$. For the captured dynamic performance sequences, we have the mesh sequence $\mathcal{D}_H = \{\mathcal{T}_H^p(i)\}$, $i$ denotes the $i$-th subject, that both making expressions and changing head and neck poses, where we have twelve identities in total. At last, we use all recovered textures, including diffuse, specularity, and normal maps, together with textures from the online dataset as $\mathcal{X}$ for our appearance learning, which consists of 360 identities. We have also captured 2D images and normal maps $\zeta_H^c$, $\zeta_H^n$ of the speaking sequences with a total of 6000 frames for larynx motion applications. In summary, the full collected dataset for our HACK model learning is:

$$\text{DATA} = \{\mathcal{J}_U^r, \mathcal{N}_U^r, \mathcal{N}^r, \Gamma_H^r, \mathcal{T}_H^r, \mathcal{D}_H, \mathcal{X}, \zeta_H^c, \zeta_H^n\}, \tag{5}$$

with skeleton $\mathcal{J}_U^r$ well matching neutral meshes in rest pose $\mathcal{N}_U^r$, combined neutral meshes $\mathcal{N}^r = \{\mathcal{N}_H^r, \mathcal{N}_I^r\}$, larynx geometry $\Gamma_H^r$ is calculated from scanned neutral meshes $\mathcal{N}_H^r$ and dynamic sequence $\mathcal{D}_H = \{\mathcal{T}_H^p(i)\}$. Our comprehensive multi-modal dataset covers both internal anatomical structures and external appearances, which

builds an important biomechanical prior for HACK, allowing for anatomically-consistent and realistic modeling and animation.

## 5 LEARNING HACK

Recall that our HACK model provides anatomy-consistent and expressive disentanglement of the head and neck regions through a set of parameters, i.e., $\{\boldsymbol{\beta}, \boldsymbol{\psi}, \boldsymbol{\theta}, \eta, \tau, \boldsymbol{\alpha}\}$ to represent the shape, expression, skeletal pose, larynx motions, and appearance, respectively. Similarly to previous models [Li et al. 2017; Loper et al. 2015], we utilize blendshapes to transfer the parameters into geometric deformations. We hence formulate the mean shape $\bar{\mathrm{T}}$, shape blendshapes $\mathcal{S}$, expression blendshapes $\mathcal{E}_{\boldsymbol{\beta}}$, pose blendshapes $\mathcal{P}_{\boldsymbol{\beta}}$ and larynx function $L(\cdot)$, as in Eqn. 3. We further introduce learning the mapping $\mathcal{M}_E$ and $\mathcal{M}_P$ from the identity parameter $\boldsymbol{\beta}$ to expression and pose blendshapes, respectively, to augment the capability for personalized modeling. Besides, for the LBS skinning process, we construct the joint position regressor $J(\cdot)$ and skinning weights $\mathcal{W}$. In a nutshell, we set out to learn the model parameters including $\{\bar{\mathrm{T}}, \mathcal{S}, J, \mathcal{E}_{\boldsymbol{\beta}}, \mathcal{P}_{\boldsymbol{\beta}}, L, \mathcal{W}\}$ in a two-stage optimization framework based on various data modalities. We first learn the shape blendshapes $\mathcal{S}$, joint regressor $J$, expression and larynx blendshapes $\mathcal{E}_{\boldsymbol{\beta}}$, $\mathcal{L}_{\boldsymbol{\beta}}$ from the captured data under the rest pose (Sec. 5.1). Then, using the dynamic and sequential data, we further estimate the pose blendshapes $\mathcal{P}_{\boldsymbol{\beta}}$ and skinning weight $\mathcal{W}$ with carefully designed regularizations (Sec. 5.2).

### 5.1 Static Geometry Modeling

*Identity Learning.* The identity learning involves computation of the mean head and neck mesh $\bar{\mathrm{T}}$ and building shape blendshapes $\mathcal{S}$, where we use the static neutral meshes in rest pose $\mathcal{N}^r$ as mentioned in Sec. 4.3. The mean template mesh $\bar{\mathrm{T}}$ is calculated as the mean of $\mathcal{N}^r$. We then apply the PCA on the displacement between $\mathcal{N}^r$ and $\bar{\mathrm{T}}$ to find the shape blendshapes $\mathcal{S} = [\mathrm{S}_1, \ldots, \mathrm{S}_{|\boldsymbol{\beta}|}] \in \mathbb{R}^{3N \times |\boldsymbol{\beta}|}$ by only preserving the first $|\boldsymbol{\beta}|$ principal components. Notice that for $\mathcal{N}^r$, we have removed the larynx from the neutral meshes, hence the shape blendshapes $\mathcal{S}$ does not contain the larynx geometry.

*Larynx blendshapes.* Similar to the shape blendshapes, we learn the larynx blendshapes from $\Gamma_{\mathrm{H}}^r$ using the PCA decomposition. Instead of conducting PCA on 3D geometry directly, we model the larynx geometric displacement w.r.t. neutral mesh in the UV space for convenience. Specifically, we un-warp $\Gamma_{\mathrm{H}}^r$ into a 2D image through texture atlasing and storing the vertex displacement into corresponding uv coordinate. We conduct PCA on the obtained images and obtain larynx blendshapes basis $\mathrm{L}_i(\tau = 0) \in \mathbb{R}^{3 \times H \times W}$, where $\tau = 0$ means the larynx blendshapes basis is computed at the rest pose.

*Head and neck skeleton regression.* Here we adopt the anatomically-consistent skeletal structure from Sec. 3.1. Similar to preivious models [Li et al. 2017; Loper et al. 2015], we learn a mapping function $J(\boldsymbol{\beta}) : \mathbb{R}^{|\boldsymbol{\beta}|} \mapsto \mathbb{R}^{3K}$ which takes the identity parameter $\boldsymbol{\beta}$ as input and regresses the joint positions. We use the aligned joints and neutral mesh data, i.e., $\mathcal{J}_{\mathrm{U}}^r$ and $\mathcal{N}_{\mathrm{U}}^r$, to learn $J$, by minimizing the

Euclidean distance loss:

$$E_{\mathrm{joint}} = \sum_i \|J(\boldsymbol{\beta}_{\mathrm{U}}(i)) - \mathcal{J}_{\mathrm{U}}^r(i)\|_2^2, \tag{6}$$

where $\boldsymbol{\beta}_{\mathrm{U}}(i)$ is the identity parameter of the $i$-th subject estimated from $\mathcal{N}_{\mathrm{U}}^r$ using the shape blendshapes $\mathcal{S}$. With the learned $J$, we are able to predict joints $\mathcal{J}^r$ for all neutral meshes $\mathcal{N}^r$.

*Person-specific expression blendshapes.* From the static scans of subjects with 28 FACS expressions $\mathcal{T}_{\mathrm{H}}^r$, we further construct the person-specific expression blendshapes for the captured subject with identity parameter $\boldsymbol{\beta}$, denoted as $\mathcal{E}_{H,\boldsymbol{\beta}} = [\mathrm{E}_{H,1}^{\boldsymbol{\beta}}, \ldots, \mathrm{E}_{H,|\boldsymbol{\psi}|}^{\boldsymbol{\beta}}] \in \mathbb{R}^{3N \times |\boldsymbol{\psi}|}$, using the technique proposed by Li et al. [2010]. Then we have a set of person-specific expression blendshapes, denoted as $\{\mathcal{E}_{H,\boldsymbol{\beta}}\} = \{\mathcal{E}_{H,\boldsymbol{\beta}(H,1)}, \ldots, \mathcal{E}_{H,\boldsymbol{\beta}(H,n)}\}$, for subjects in $\mathcal{T}_{\mathrm{H}}^r$, where $\boldsymbol{\beta}(\mathrm{H}, i)$ indicates the $i$-th subject in $\mathcal{T}_{\mathrm{H}}^r$. Different from previous face methods, such as FLAME [Li et al. 2017] and ICT-FaceKit [Li et al. 2020a], that model a general expression blendshapes independent from identity information, we set out to model personalized expressions by learning a mapping network $\mathcal{M}_E$ which maps one's identity parameter $\boldsymbol{\beta}$ to its personalized expression blendshapes. Specifically, we would like to train $\mathcal{M}_E$ on our set of person-specific expression blendshapes $\{\mathcal{E}_{H,\boldsymbol{\beta}}\}$, with the following loss:

$$E_{\mathrm{exp}} = \sum_i \left\| \mathcal{M}_E(\boldsymbol{\beta}(H, i)) - \mathcal{E}_{H,\boldsymbol{\beta}(H,i)} \right\|_2. \tag{7}$$

Notice that the expression blendshapes, $\mathcal{E}_{H,\boldsymbol{\beta}}$, exist in a high dimensional space that may be difficult for the mapping network to learn. To address this challenge, we first apply PCA to the set of person-specific expression blendshapes, $\{\mathcal{E}_{H,\boldsymbol{\beta}}\}$, and the mapping network, $\mathcal{M}_E$, is then trained to predict the PCA weights. This allows for more efficient learning and improves the robustness while preserving personalized expressions as the predicted blendshapes rely on $\boldsymbol{\beta}$. In the implementation, $\mathcal{M}_E$ consists of three linear layers with 64 neurons and ReLU activation, and predicts the weights of the first 50 principal components. The training is conducted using Pytorch Adam optimizer with a learning rate of 0.0001, and converges in 2 hours on a single Nvidia Titan GPU.

### 5.2 Dynamic Deformation Learning

We have obtained the shape blendshapes $\mathcal{S}$, person-specific expression blendshapes $\mathcal{E}_{\boldsymbol{\beta}}$, and larynx blendshapes $\mathcal{L}$ for all identities. To further model the deformations when changing the head and neck pose, we rely on the dynamic sequences $\mathcal{D}_{\mathrm{H}}$ for learning the pose blendshapes, larynx function $L(\boldsymbol{\beta}, \eta, \tau; \mathcal{L})$ and skinning weights $\mathcal{W}$. Notice that $\mathcal{W}$ is shared across all identities, and we hire professional animation artists to create initial $\tilde{\mathcal{W}}$ as a prior skinning weight.

To this end, we jointly learn $\mathcal{P}_{\boldsymbol{\beta}}$, $L(\boldsymbol{\beta}, \eta, \tau; \mathcal{L})$ and $\mathcal{W}$ via minimizing the reconstruction loss on $\mathcal{D}_{\mathrm{H}}$ as:

$$E_{\mathrm{rec}} = \sum_{i=1}^{|\mathcal{D}_{\mathrm{H}}|} \sum_{j=1}^{|\mathcal{T}_{\mathrm{H}}^p(i)|} \left\| \mathrm{G}(\boldsymbol{\beta}(i), \boldsymbol{\psi}(i, j), \boldsymbol{\theta}(i, j), \eta(i), \tau(i, j)) - \mathcal{T}_{\mathrm{H}}^p(i, j) \right\|_2^2, \tag{8}$$

where $\mathrm{G}(\cdot)$ is the geometry modeling as described in Sec. 3. $i$ is the index of the subject and $j$ indicates the frame in a sequence of a
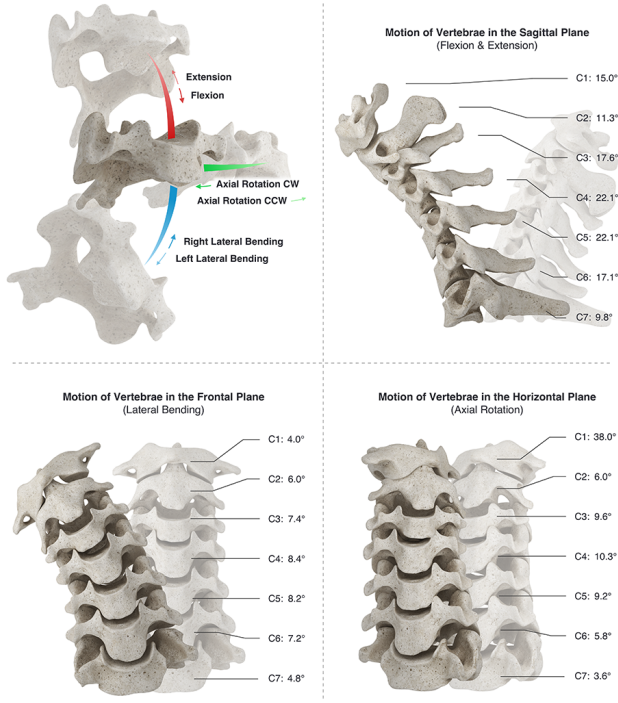
Fig. 5. Rotation limit of the cervical spine. We first illustrate the degree of freedom of rotation of cervical vertebrae. Then we show the average rotation limit in angles of each cervical vertebra under different cervical motions, including flexion and extension (upper right), lateral bending (bottom left), and axial rotation (bottom right). We design the rotation limit upon a comprehensive anatomical survey [Savlovskis 2022].

certain identity. Notice that $\mathcal{W}$ is updated in the skinning step as in Eqn. 2. Moreover, we also fine-tune the person-specific expression blendshapes $\mathcal{E}_{\boldsymbol{\beta}}$ during the training process for better modeling the neck motion in this stage.

Besides supervision from captured meshes, we additionally design several regularization terms, including joint rotation limits, pose regularization, and collision penalization, based on anatomy and physical priors to improve the robustness of training.

*Joint rotation limits.* The moving range of the human cervical vertebrae is limited due to its unique anatomy structure, as illustrated in Fig. 5. To simulate the head and neck motion correctly and obtain realistic model deformations, we limit the rotation angles that can be applied to each joint according to Fig. 5. Specifically, we convert the pose parameter $\boldsymbol{\theta}$ to Euler angles and enforce penalization when any angle of a certain joint exceeds the corresponding limits. The penalization is formulated as the following loss:

$$E_{\text{rot}} = \sum_{i,j} \sum_{1 \leq k \leq K} \left\| \max\{|\theta^k(i,j)| - b^k, 0\} \right\|_1, \tag{9}$$

where $i, j$ denote the identity and frame in dynamic sequence respectively. $\theta^k$ is the Euler angle of $k$-th joint of given pose $\boldsymbol{\theta}$, $b_k$ is the corresponding rotation limits exhibited as in Fig. 5.

*Adjacent joints consistency.* Further, the joints do not rotate independently, and the rotations of adjacent cervical vertebrae exhibit certain similarities. Hence, we minimize the rotation differences of Euler angles between the adjacent joints, as:

$$E_{\text{sim}} = \sum_{1 \leq i < K} \|\theta_i - \theta_{i+1}\|_2^2, \tag{10}$$

where $\theta_i \in \boldsymbol{\theta}$ is the $i$th rotation vector of given pose $\boldsymbol{\theta}$.

*Collision penalization.* We find that the neck modeling is prone to noise and inaccuracies due to heavy occlusion, particularly when the head is tilted. To address this, we enforce collision penalization to ensure that the cervical spine is not intersected with the neck skin, and we re-use the collision term $E_{col}$ proposed in [Hasson et al. 2019].

*Temporal smoothness.* As the per-frame learning scheme leads to noises, we introduce temporal consistency terms for parameters related to motion, including $\boldsymbol{\psi}, \boldsymbol{\theta}, \eta, \tau$. For every scalar element $v$ in $\{\boldsymbol{\psi}, \boldsymbol{\theta}, \eta, \tau\}$, we aggregate the estimated $v$ for a captured sequence $s$ and form the vectors $\{v(t)\}_s$ for each $v$, where $t$ indicates the frame. We enforce the Lipschitz continuity and minimize the second derivative of $\{v\}_s$ as follows:

$$E_{\text{tem}} = \sum_s \lambda_v \sum_{v(t) \in \{v(t)\}_s} \left( \lambda_1 \max\{|\frac{d[v(t)]}{dt}| - \epsilon, 0\}^2 + \lambda_2 (\frac{d^2[v(t)]}{dt^2})^2 \right), \tag{11}$$

where $\lambda_v, \lambda_1, \lambda_2$ are the weighting factors, $\epsilon$ is threshold that to tolerate small deviations. For $\boldsymbol{\psi}, \lambda_v = 1, \lambda_1 = 1, \lambda_2 = 0.01, \epsilon = 0.1$; for $\boldsymbol{\theta}, \lambda_v = 1, \lambda_1 = 1, \lambda_2 = 5000, \epsilon = 0.15$; for $\boldsymbol{\eta}, \lambda_v = 1, \lambda_1 = 1, \lambda_2 = 10, \epsilon = 0.005$; for $\boldsymbol{\tau}, \lambda_v = 1, \lambda_1 = 1, \lambda_2 = 1, \epsilon = 0.001$.

*Regularization.* Similar to the temporal smoothness of parameters, we want to ensure the generated geometries are also smooth in temporal dimensions. We apply the Laplacian smoothness penalty $E_{\text{smo}}$ for all pose blendshapes. Also, $E_{\text{ski}} = \|\mathcal{W} - \tilde{\mathcal{W}}\|_2^2$ is introduced to regularize skinning weight to not deviate from artist-created initial weights too much. The final optimization loss is the combination as follows:

$$E = E_{\text{rec}} + E_{\text{rot}} + E_{\text{sim}} + E_{\text{col}} + E_{\text{tem}} + E_{\text{smo}} + E_{\text{ski}}, \tag{12}$$

where the weighting factors for each penalty are 1e5, 1e6, 5e3, 5e5, 1e6, 5e-2, 1, respectively. During training, we estimate the person-specific pose blendshapes $\mathcal{P}_{\boldsymbol{\beta}}$ with pose parameters $\boldsymbol{\theta}$, expression parameters $\boldsymbol{\psi}$ and larynx parameters $\eta, \tau$, jointly. During the learning process, HACK can be considered as a differentiable layer, and the optimization process is carried out using Adam optimizer in PyTorch. Each dynamic performance sequence has an average length of 850 frames, and the overall loss converges after 45000 epochs with the learning rate set at 0.001.

After every person-specific pose blendshapes converge for each identity, we have a set of person-specific pose blendshapes, denoted as $\{\mathcal{P}_{\boldsymbol{\beta}}\} = \{\mathcal{P}_{\boldsymbol{\beta}(H,1)}, \ldots, \mathcal{P}_{\boldsymbol{\beta}(H,n)}\}$, where $\boldsymbol{\beta}(H, j)$ indicates the $j$-th subject in $\mathcal{D}_H$. Still, instead of modeling a general pose blendshapes as previous methods [Li et al. 2017; Loper et al. 2015], we further model personalized pose blendshapes by learning the mapping network $\mathcal{M}_P$ that maps identity code to its personalized

pose blendshapes. The learning process of $\mathcal{M}_P$ is the same as expression mapping network $\mathcal{M}_E$ as in Eqn. 7. Similarly, to reduce the dimensionality of this even larger space $\mathcal{M}_P \in \mathbb{R}^{3N \times 72}$, PCA is applied on $\{\mathcal{P}_\beta\}$ similar to learning $\mathcal{M}_E$. Due to the complexity of capturing dynamic head and neck performance sequences, we captured twelve identities that have person-specific pose blendshapes for $\mathcal{M}_P$ to learn. However, modeling personalized pose blendshapes instead of the general ones is still proved effective that better restores the identity-dependent characteristics, which will be evaluated in Sec. 7.1. As more dynamic head and neck performance sequences are captured, modeling personalized motion features, including expression and pose blendshapes, will demonstrate its strong capability and continue to benefit the community of parametric models.

### 5.3 Appearance Learning

For authentic appearance generation, we create a parametric appearance model from our physically-based texture dataset, including diffuse albedo, specular intensity, and normal map. Using the multi-view photometric system, we capture the rich physically-based appearances for the whole head-and-neck regions of various subjects, and subsequently unwarp the appearance into unified texture maps from the neutral scans at the rest pose. We follow previous work [Qian et al. 2020; Qiu et al. 2022] to unify all the physically-based texture data, and subsequently perform principal component analysis (PCA) using singular value decomposition directly to obtain the principal components in the UV space. To this end, we obtain the appearance model $A(\boldsymbol{\alpha})$ to generate realistic textures from a random appearance parameter $\boldsymbol{\alpha}$. Since our textures have uniform texture UV mapping as our template skin mesh, we could directly apply generated physically-based textures with various shapes and produce a photo-realistic appearance.

## 6 APPLICATIONS

As a parametric model that tackles the head and neck as a whole with expressive controls, HACK is differentiable and compatible with standard CG software, making it suitable for a range of applications like geometric fitting, animation, and inference. HACK also enables fine-grained analysis of the inter-correlation of head and neck. In the following, we demonstrate various applications like motion synthesis or transfer towards such unique characteristics.

### 6.1 Motion Synthesis from Head to Neck

The head poses and facial expressions are highly correlated with the skeletal pose of the cervical spine and neck deformation. Take the speaking action as a typical example, where the movement of the larynx is coordinated with the movements of the mouth, which could be faithfully captured by the expression parameters, as it controls the vocal cords. Hence, we use the transformer architecture from FaceFormer [Fan et al. 2022] to regress the larynx movement sequence $\{\tau\}_m$ from $\Psi$, as shown in Fig. 6.

We train the transformer using our captured larynx motion sequences $\zeta_{\mathrm{H}}^c$ and $\zeta_{\mathrm{H}}^n$, where $\zeta_{\mathrm{H}}^c$ and $\zeta_{\mathrm{H}}^n$ are terms of RGB images and normal maps, respectively. We utilize the facial tracker NPFA [Zhang et al. 2022] to predict the expression parameters from $\zeta_{\mathrm{H}}^c$. We also adopt a novel convolutional method to track the movement of the
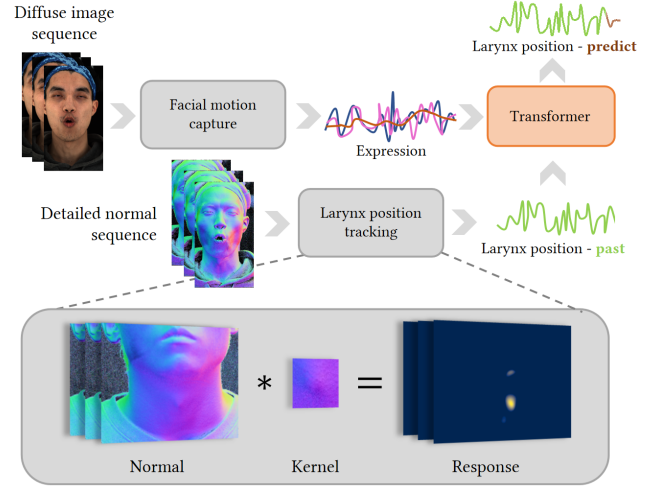


Fig. 6. The pipeline of generating larynx motion sequence using transformer. After extracting expressions and larynx positions from captured sequences, we train the transformer to predict larynx motion in an auto-regressive manner.

larynx in the vertical direction directly from image inputs. Specifically, we define a kernel k of size $70 \times 70 \times 3$ denoting the generic shape (represented as normal) of the larynx. We extract the larynx position by finding the maximum response on the dot product of k and $\zeta_{\mathrm{H}}^n$ in a convolutional manner, which can be formulated as follows:

$$\bar{\tau}(i, j) \leftarrow \arg\max_{x} \left[ \mathrm{k} * \zeta_{\mathrm{H}}^n(i, j) \right](x) - \tau_0(i), \qquad (13)$$

where $*$ denotes the convolution operator; $x$ denotes the spatial position in normal maps; $\tau_0(i)$ is the initial position of the larynx in the rest pose of the $i$-th identity; $j$ denotes the frame index in the sequence. We follow the framework of the FaceFormer. In particular, we fed the expression $\Psi$ into the attention layers of our transformer, and use the past sequence $\bar{\tau}(i, 1 \ldots k - 1)$ to predict the next signal $\bar{\tau}(i, k)$. Once trained, the adopted transformer is able to predict the larynx position sequence $\{\tau\}$ from unseen expressions $\{\Psi\}$ for novel speaking. Then, we can use $\{\tau\}$ and $\{\Psi\}$ to synthesize the coordinated larynx motion and speaking action.

Similarly, we can capture the correlations between the head pose and the skeletal pose of the cervical spine to effectively drive a HACK model with existing facial motion capture techniques, like DECA [Feng et al. 2021] and Apple ARKit [Apple 2023]. Specifically, we calculate the facial orientation using the global rotation of the head skull (the o-c1 joint) on our dynamic performance sequences $\mathcal{D}_{\mathrm{H}}$, and set up a small MLP that consists of 2 layers with 512 neurons to learn the mapping from facial orientations to pose parameters $\boldsymbol{\theta}$ in our definition. Then, with the correct mapping, we can manipulate the HACK model to display a range of personalized expressions, poses, and appearances with accurate neck deformations.

Fig. 7. Samples using HACK, of various identities, head and neck poses, expressions, and appearance. Our model demonstrates its strong ability for realistic modeling, animation and rendering.

## 6.2 Cross Species Motion Transfer

Motion transfer or retargeting is the process of transferring expressions, poses, and other motions from one character to another. In this section, we demonstrate another unique ability of HACK: the ability to transfer expressions and poses from humans to another mammal model, in this case, a giraffe. Existing generic models cannot perform this type of transfer as they do not model neck motion. However, with HACK's neck modeling ability, the transferred neck motion achieves a high level of realism due to the neck skeleton's structural similarity between the characters.

We start by using a giraffe model with defined neck skeleton joints, which has been created by artists. We require the giraffe model to have the same topology as the HACK template model and in rest pose. Next, we apply the estimated expression and pose parameters, which were fitted using human models, to the giraffe model. Because HACK's head and neck skeleton is anatomically consistent and mammals share the same neck skeleton structure, the transferred poses on the giraffe are highly realistic, as demonstrated in Fig. 18.

## 7 RESULTS

### 7.1 Model Evaluation

*Qualitative results.* HACK offers personalized and anatomically-consistent controls for the neck regions by leveraging comprehensive inner biomechanical priors. As demonstrated in Fig. 7, by controlling the parameters of HACK, HACK can generate diverse results and even form a delicate artwork. Being a parametric model, HACK has the ability to fit on existing meshes under various neck poses. Enhanced by modeling personalized characteristics such as expression and pose blendshapes, HACK is powerful in restoring their unique traits and nuanced personalities, as illustrated in Fig. 8. HACK combines inner anatomical structures and physically-based appearance to model the full spectrum of neck motions and detailed facial expressions. It supports personalized and anatomically-consistent controls using existing facial trackers, while providing realistic animation

Table 1. Quantitative comparison on neutral registration on the FaceScape dataset [Yang et al. 2020], VOCASET [Cudeiro et al. 2019], ICT-3DRFE [Stratou et al. 2011] and Multiface [Wuu et al. 2022].

| Dataset | FLAME | ICT-FaceKit | HACK |
|---|---|---|---|
| FaceScape | $2.401 \pm 0.915$ | $2.194 \pm 1.644$ | $1.929 \pm 0.709$ |
| VOCASET | $1.406 \pm 0.280$ | $0.958 \pm 0.148$ | $0.913 \pm 0.133$ |
| ICT-3DRFE | $0.397 \pm 0.045$ | $0.366 \pm 0.041$ | $0.376 \pm 0.048$ |
| Multiface | $0.943 \pm 0.082$ | $0.901 \pm 0.120$ | $0.842 \pm 0.127$ |
| avg. | $1.298 \pm 0.972$ | $1.132 \pm 1.173$ | $1.035 \pm 0.750$ |

and rendering results, as shown in Fig. 9. As shown in Fig. 10, HACK demonstrates strong generalization ability by successfully registering various released datasets and generating novel appearances and poses with realistic rendering.

*Quantitative evaluation.* Fig. 11 shows the compactness of learned spaces in HACK. These red curves describe the explained-variance ratio in the training data with respect to the number of components. In shape space, the curve in Fig. 11(a) implies that with the first 50 principal components, the shape space is able to cover 95.5% of the entire space. Meanwhile, 200 principal components are sufficient to express 99.5% of the entire space. In larynx shape space, Fig. 11(b) shows that the first 10 principal components achieve over 95.2% of the larynx shape space, while 30 components are able to cover 99.0% of the larynx shape space. In expression blendshapes space, Fig. 11(c) shows that the first 30 principal components achieve over 86.4% of the expression blendshapes space, while 50 components are able to cover 93.7% of the whole space. In pose blendshapes space, Fig. 11(d) shows that 6 components are able to cover 87.9% of the space.

To evaluate the generalization of the shape space, larynx shape space, and expression blendshapes space, we randomly divided 90% of the data for training and the remaining 10% for testing.
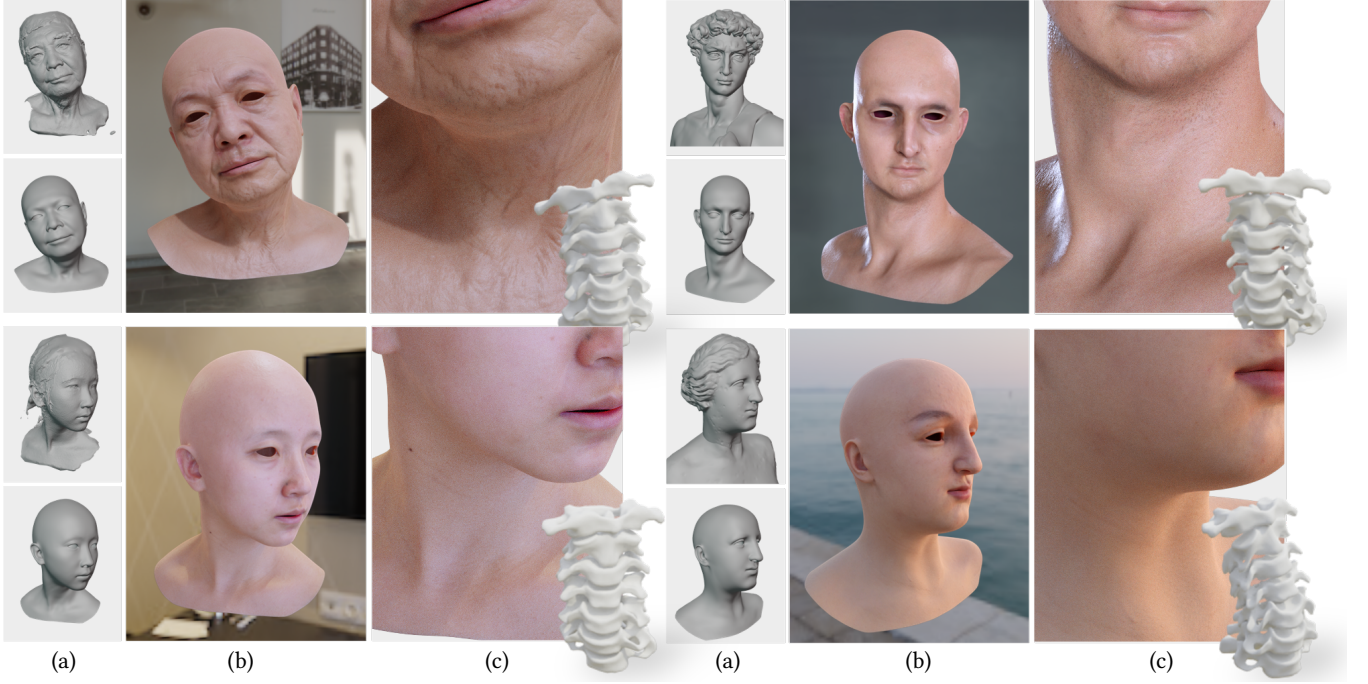
Fig. 8. Registration results of HACK on our testing dataset and masterpiece statues. As discussed in Sec. 6, we can fit HACK on different data, ranging from our captured unseen scans to masterpiece statues. In the left part, we first demonstrate the registration of an old man and a young lady from our captured data, then, in the right part, we show the registration of *David*, a masterpiece of Renaissance sculpture, and *Venus de Milo*, an ancient Greek sculpture that was created during the Hellenistic period, both with an elegant neck. For each case, we show (a) the target and HACK's registration; (b) the registration with appearance sampled from HACK's appearance space; (c) the zoom-in view to demonstrate the details and corresponding neck pose. HACK successfully models across identities, poses, and expressions, with details restored and can be realistically rendered with appearance applied.



Fig. 9. Novel pose and inference results of previous HACK registrations. For each case, we provide (a) the same identities with synthetic novel poses; (b) driven results using facial performance capture, as discussed in Sec. 6. Note that for *David*, our driving result provides the realistic bowstring effect where the platysma is contracting, thanks to our joint modeling of the head and neck. As shown in the figure, HACK provides realistic driven and rendering results.

Fig. 10. Registration results on various released dataset, including samples from FaceScape [Yang et al. 2020] (upper left), 3D Scan Store [Store 2022] (upper right), MultiFace [Wuu et al. 2022] (lower left), and VOCA [Cudeiro et al. 2019] (lower right). For each sample, we show (a) the original mesh, (b) HACK registration result, (c) the realistic rendering result using our physically-based appearance, and (d) with novel expressions and poses.
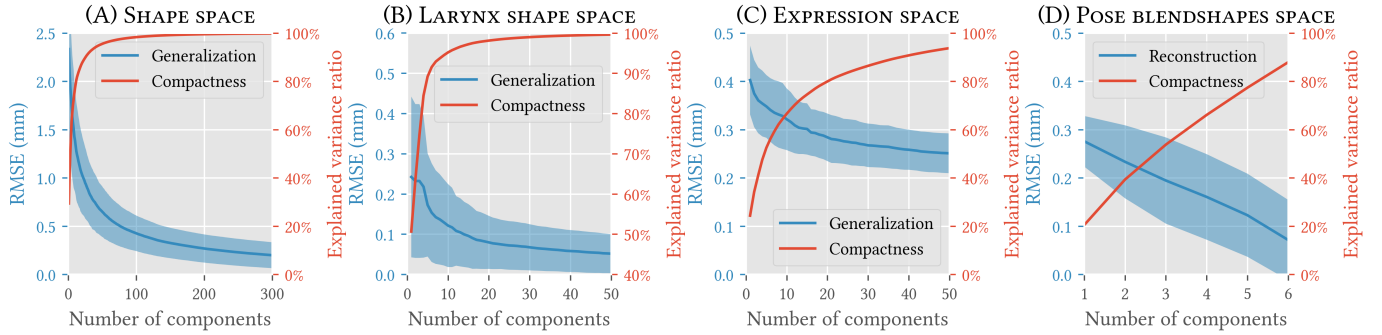


Fig. 11. Quantitative evaluation of compactness and generalization.

Fig. 11(a)(b)(c) plot the evaluation result on the generalization ability of the HACK shape, larynx shape, and expression blendshapes spaces, respectively. The generalization error is shown as the mean and standard deviation of RMSE with respect to the number of components. The shape space's generalization error decreases monotonically on the testing shape as the number of components increases. The error curve reaches below 0.65 mm and 0.26 mm with 50 and 200 components, respectively. Similarly, the larynx shape space and expression blendshapes space also exhibit decreasing error curves as the number of components increases. Fig.11(d) shows the reconstruction error of the pose blendshapes space, which decreases as the number of components increases.

*Ablation.* Modeling the space of expression blendshapes provides fine-grained geometry for different identities under different expressions. In contrast, using generic expression blendshapes ignores the diversity of facial muscles across identities and under different expressions. We denote the original model and the variant one using generic expression blendshapes as **personalized** $\mathcal{E}$ and **generic** $\mathcal{E}$,

respectively. We utilize both models to fit unseen expressions. As illustrated in Fig. 12, modeling person-specific expression blendshapes results in better reconstruction of facial geometry under different expressions, including single and complex expressions.

Similarly, modeling the space of pose blendshapes provides fine-grained geometry for different identities under different poses. In contrast, using generic pose blendshapes ignores person-specific attributes in animation. We denote the original model and the variant one using generic pose blendshapes as **personalized** $\mathcal{P}$ and **generic** $\mathcal{P}$, respectively. We utilize both models to fit unseen poses. As illustrated in Fig. 13, modeling person-specific pose blendshapes results in better reconstruction of neck geometry under different poses. One observation is that body fat can obscure the underlying muscle tissue, making the muscles less visible and less defined, and vice versa. Therefore, modeling pose blendshapes space results in better fitting accuracy and higher expressiveness.
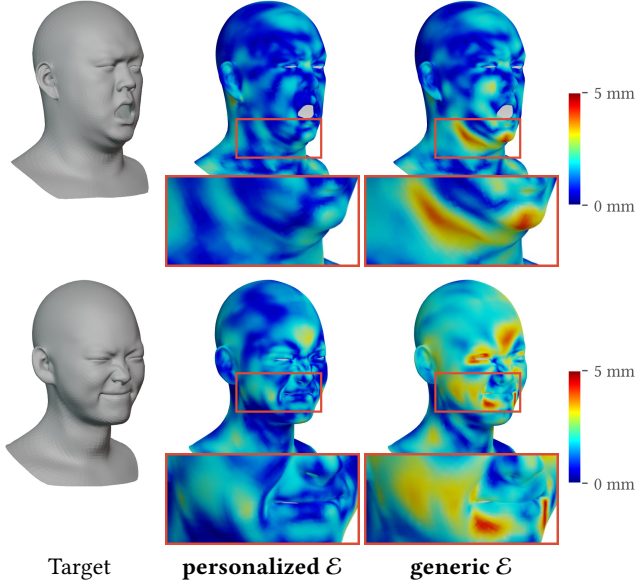
Fig. 12. Qualitative evaluation of whether to use person-specific expressions. The upper performer opens his mouth but with slightly different jaw movements than normal people due to his high body fat rate. The lower performer makes a "fully compressed" face that is extremely difficult to fit expression using blendshapes. With modeling personalized expressions, our model achieves better reconstruction accuracy on rich expressions, especially around the mouth and jaw.
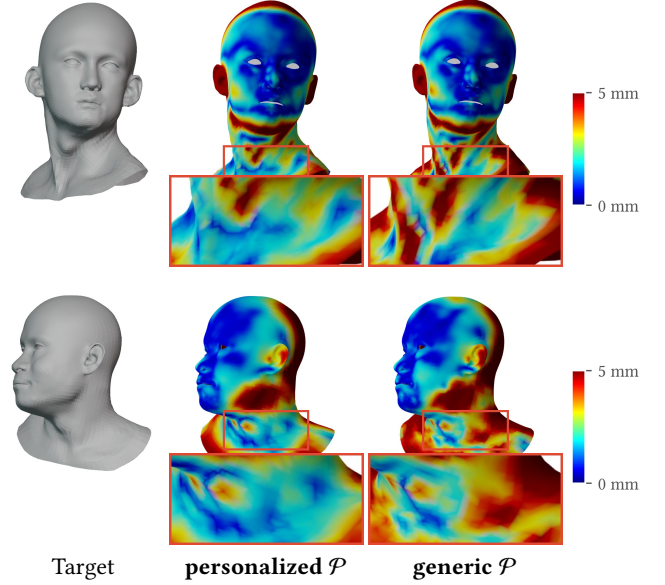


Fig. 13. Qualitative evaluation of whether to use person-specific pose blendshapes. Using generic pose blendshapes ignores personalized attributes of head and neck movements. For example, performers with extreme body fat rates will result in the Sternocleidomastoid muscle being either stronger prominent bulged (upper) or less defined (lower). With the modeling of personalized pose blendshapes, our model has lower reconstruction error on posed scans.
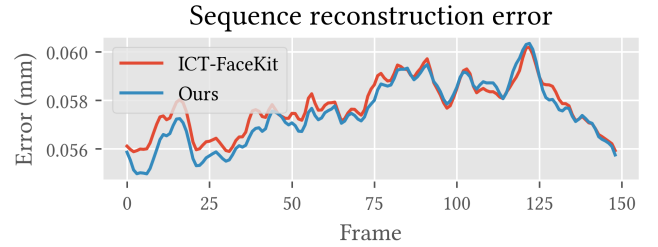
## 7.2 Comparison

*Quantitative results.* In the task of reconstructing faces in a rest pose, we compared HACK with two state-of-the-art parametric models, FLAME [Li et al. 2017] and ICT-FaceKit [Li et al. 2020a], to demonstrate its generalization ability. We experimented with various identities, including 20 scans from FaceScape [Yang et al. 2020], 12 scans from VOCASET [Cudeiro et al. 2019], 22 scans from ICT-3DRFE [Stratou et al. 2011], and 10 scans from Multiface [Wuu et al. 2022]. To fit identity shape, we used the official *flame-fitting* repository [Li et al. 2017], with the weights of "scan-to-mesh", "landmark" and "shape regularization" setting to 2, 0.01 and 0.00005, respectively. To ensure a fair comparison, we used the first 100 shape components in each model and removed outlier point clouds in the raw scans during optimization. Table 1 reports the quantitative results of the mean reconstruction error, as the mean and the standard deviation of the average scan-to-mesh distance in millimeters. In Fig. 14, we compared the use of generic expressions from ICT-FaceKit to personalized expressions of HACK by reconstructing talking sequences from VOCASET.

*Qualitative results.* We further compare the qualitative reconstruction results between FLAME, ICT-FaceKit and HACK. In Fig. 15, we reconstruct neutral scans from FaceScape [Yang et al. 2020], VO-CASET [Cudeiro et al. 2019], ICT-3DRFE [Stratou et al. 2011], and Multiface [Wuu et al. 2022]. We provide qualitative registration results and the corresponding mesh-to-scan distance. Note that a



Fig. 14. Quantitative comparison with ICT-FaceKit on sequence from VO-CASET, which uses generic expression blendshapes instead of personalized ones.

regularization term is applied during the fitting process to minimize the fitting objective and avoid artifacts. In Fig. 16, we compare the fitting of posed scans from our captured testing data for both FLAME and HACK, and provide corresponding error maps.

These comparison results illustrate that our HACK model provides accurate disentanglement for human identity, facial expression, skeletal pose for neck region, and larynx motions. Note that HACK is not intended to outperform previous models in every individual component such as shape or expression. Rather, it provides a comprehensive and anatomically-consistent model for the neck regions. More significantly, we show that by further considering nuanced neck and larynx motions in human parametric modeling, HACK is
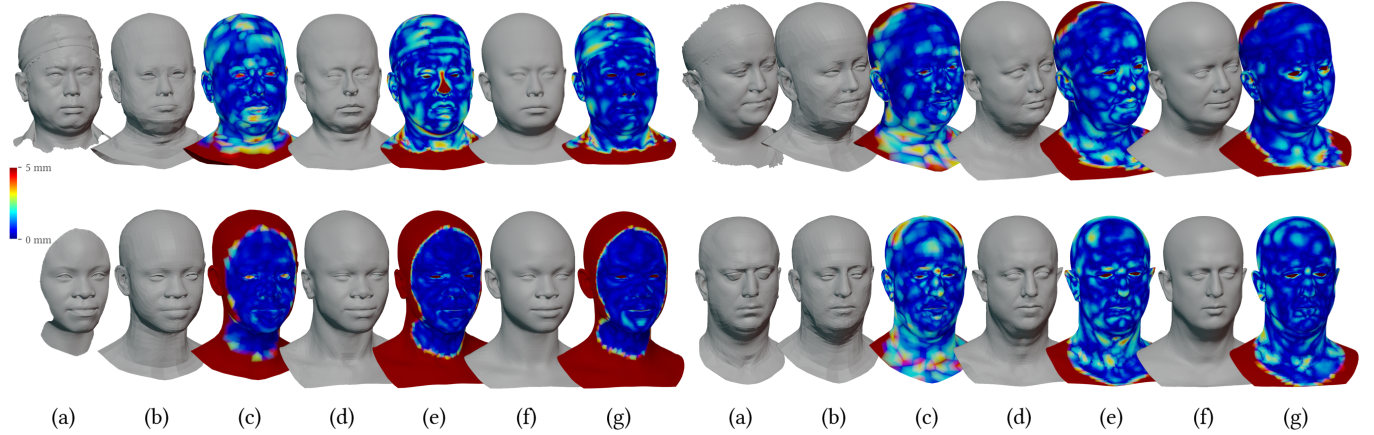
Fig. 15. Qualitative comparison on registration of samples from FaceScape dataset [Yang et al. 2020] (upper left), VOCASET [Cudeiro et al. 2019] (upper right), ICT-3DRFE [Stratou et al. 2011] (lower left) and Multiface [Wuu et al. 2022] (lower right). In the figure, we show (a) the original scan, (b) FLAME registration result, (c) mesh-to-scan distance of FLAME, (d) ICT-FaceKit registration result, (e) mesh-to-scan distance of ICT-FaceKit, (f) HACK registration result, (g) mesh-to-scan distance of HACK.
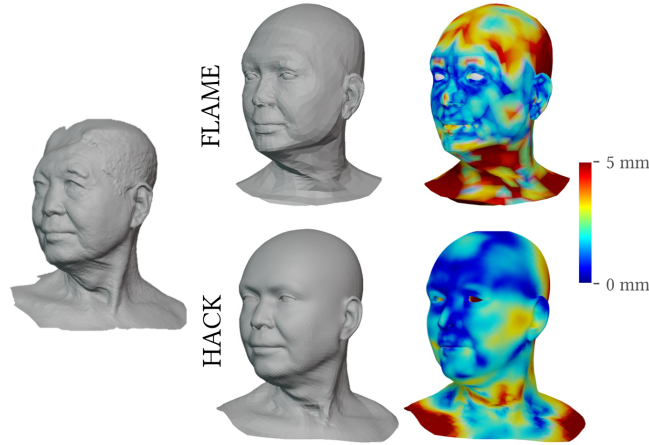


Fig. 16. Qualitative comparison with FLAME on registration of scan from our captured testing data with a head pose.



Fig. 17. The results of generating larynx motion from expressions when speaking. The performer is speaking the sentence "She always tells me to smile and put on a happy face". The figures show the larynx position in different frames (upper) and the corresponding generated position sequence (lower).

able to provide personalized and realistic controls for a wide range of applications.

## 7.3 Application result

As a parametric model, we can fit the parameters of HACK to a human head and neck mesh, and further synthesize new expressions, poses, and larynx motions with high fidelity. Existing head modeling techniques lack the ability to accurately estimate neck geometry. To demonstrate this, we fit the parameters of HACK to the 3D model of Michelangelo's highly regarded sculpture, *David*, which relies on the neck for conveying movement and energy. Using HACK, we are able to animate *David* with more realism by exploiting HACK's neck modeling ability. Given a target mesh of *David*, we fit HACK parameters on it while keeping his unique traits of head and neck using personalized pose blendshapes. As illustrated in the
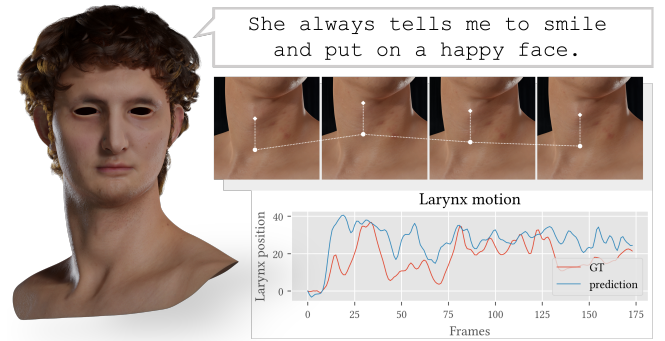
upper right of Fig. 8, HACK faithfully reconstructs his elegant neck with anatomically-consistent cervical spines. We further animate the HACK *David* model. As shown in the upper right of Fig. 9. Notice that with HACK's personalized characteristic modeling ability, *David* preserves personalized traits during animation by contracting his platysma and making a defined bowstring, which demonstrates his strength and strong emotion under different poses, as shown in the upper right of Fig. 9(b).

Larynx motion is an important but also complicated movement, and HACK adopts a fully disentangled design to fully control the appearance and motion of the larynx. To demonstrate the fine-grained control over larynx motion enabled by HACK, we simulate the swallowing action on the *David* model. Specifically, we extract the corresponding larynx slicing parameter $\tau$ for the swallowing action in the dynamic performance sequence. We align sequences
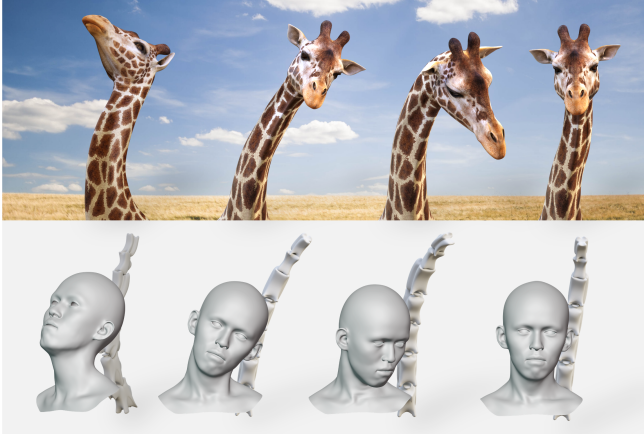
Fig. 18. The results of motion transfer from a human to a giraffe. From left to right, we demonstrate different poses including cervical extension, side-bend, flexion, and twisting. The upper figure shows the rendering result with transferred motion, and the lower figure shows the original human pose. As discussed in Sec. 6.2, HACK can transfer the motions from one character to another character even if they belong to different species. With HACK's modeling of the cervical spine, the new character's head and neck movements are transferred naturally and accurately.

$\tau$ from different clips temporally via peaks and valleys, and then average all sequences to obtain the controlling sequence for the swallowing action. We then apply the obtained $\tau$ sequence on the *David* HACK model to simulate the swallowing action. Please refer to the accompanying video for a demonstration of the realistic swallowing simulation.

*Larynx motion generation.* As discussed in Sec. 6.1, we can utilize a transformer to generate larynx motion from the expression sequence. We use a state-of-the-art facial expression extractor, ARKit [Apple 2023], to obtain the expression sequence $\Psi$. The transformer is then applied to $\{\Psi\}$, generating the larynx motion sequence $\{\tau\}$. In Fig. 17, we show a generated sequence using our testing data, which faithfully restores the larynx motion using given expressions.

*Motion transfer.* HACK can be used to transfer the neck animation from a human to a giraffe, as shown in Fig. 18. By transferring the pose parameters between our models, we are able to reproduce realistic neck movements and pose changes in the giraffe, while keeping the skin acting naturally and accurately. Thus, we can generate photo-realistic animations that accurately simulate how a giraffe's neck would move and deform under different conditions. In comparison to conventional surface models that lack an inner anatomically-prior cervical spine, HACK is able to achieve more realistic and believable neck animations, as it is able to accurately reproduce the way that the skin would behave and interact with the underlying bones. Please see the animation sequence in the accompanying video for a demonstration of the capabilities of our model.

## 7.4 Limitation and Discussion

HACK provides a compelling parametric disentanglement of the human head and neck, complete with rich inner anatomical structures. However, it still has some limitations. Here, we provide a detailed analysis and discuss potential future extensions.

First, similar to other parametric models [Ploumpis et al. 2020], our HACK model could be enhanced by integrating more individual facial components such as eyeballs or teeth [Bérard et al. 2016; Wu et al. 2016] for more comprehensive modeling. However, modeling human hair remains challenging and requires more complicated geometry representation such as strands [Winberg et al. 2022]. Further combining recent advances in neural rendering (NR) [Lombardi et al. 2021; Luo et al. 2021; Wang et al. 2022c] with HACK could provide photo-realistic hair rendering, but this may sacrifice compatibility with existing CG production pipelines. Furthermore, building a hybrid model that combines both parametric and NR modeling for various components of human characters remains an unsolved research direction with enormous potential. Nevertheless, we believe that our publicly available HACK model will serve as a solid foundation for future explorations Besides, HACK aims to strike a delicate balance between general and person-specific characteristics within the framework of parametric models. As a result, we achieve more personalized pose and expression results than previous general models such as FLAME [Li et al. 2017]. However, it remains extremely challenging to encode more nuanced personalized properties, such as dynamic textures with wrinkles or pore-level details, into the current parametric model while maintaining its convenient controls. Recent NR advances hold promise for synthesizing such personalized details, but the compatibility with CG engines remains an unsolved bottleneck issue.

## 8 CONCLUSION

We have presented HACK, a novel parametric model for constructing the head and cervical region of digital humans. It significantly models neck and larynx motions, achieving more realistic and anatomically consistent controls for the neck regions that are compatible with CG engines. Our comprehensive capturing combines 3D ultrasound imaging with multi-view photometric capture system, so as to extract the inner biomechanical structures for the vertebrae of the cervical spine, as well as the external geometry and physically-based textures. HACK separates the depiction of 3D head and neck into various blendshapes on top of the neutral one. Our anatomically-consistent skeletal pose encodes the rich cervical priors, while our expression blenshapes tied to the facial action units enable artist-friendly controls. The adopted larynx blendshapes with both larynx deformation and slicing in the UV-space further provide accurate and nuanced modeling for the larynx beneath the neck skin. We also tackle schemes to augment the pose and expression blendshapes by optimizing an efficient mapping from the identical shape space to the PCA spaces of personalized blendshapes, which significantly improves the personalized characteristics. We demonstrate the capabilities of HACK model for more realistic and nuanced controls, especially for the neck regions, and showcase the applications using the inter-correlation between head and neck for motion synthesis and transfer. It makes a solid step for covering

more spectrum of parametric digital humans, and hence facilitating numerous potential applications for entertainment, gaming, and immersive experience in VR/AR.

## ACKNOWLEDGMENTS

## REFERENCES

Victoria Fernández Abrevaya, Stefanie Wuhrer, and Edmond Boyer. 2018. Multilinear autoencoder for 3D face model learning. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 1–9.

Dragomir Anguelov, Praveen Srinivasan, Daphne Koller, Sebastian Thrun, Jim Rodgers, and James Davis. 2005. SCAPE: Shape Completion and Animation of People. In *ACM SIGGRAPH 2005 Papers* (Los Angeles, California) *(SIGGRAPH '05)*. Association for Computing Machinery, New York, NY, USA, 408–416. https://doi.org/10.1145/1186822.1073207

Apple. 2023. ARKit. https://developer.apple.com/arkit/.

Linchao Bao, Xiangkai Lin, Yajing Chen, Haoxian Zhang, Sheng Wang, Xuefei Zhe, Di Kang, Haozhi Huang, Xinwei Jiang, Jue Wang, Dong Yu, and Zhengyou Zhang. 2021. High-Fidelity 3D Digital Human Head Creation from RGB-D Selfies. *ACM Trans. Graph.* 41, 1, Article 3 (nov 2021), 21 pages. https://doi.org/10.1145/3472954

Pascal Bérard, Derek Bradley, Markus Gross, and Thabo Beeler. 2016. Lightweight Eye Capture Using a Parametric Model. 35, 4, Article 117 (jul 2016), 12 pages. https://doi.org/10.1145/2897824.2925962

Volker Blanz and Thomas Vetter. 1999. A morphable model for the synthesis of 3D faces. In *Proceedings of the 26th annual conference on Computer graphics and interactive techniques*. 187–194.

V Blanz and T Vetter. 2003. Face recognition based on fitting a 3D morphable model. *Transactions on Pattern Analysis and Machine Intelligence* 25, 9 (2003), 1063–1074.

James Booth, Anastasios Roussos, Allan Ponniah, David J. Dunaway, and Stefanos Zafeiriou. 2017. Large Scale 3D Morphable Models. *International Journal of Computer Vision* 126 (2017), 233 – 254.

Alan Brunton, Timo Bolkart, and Stefanie Wuhrer. 2014. Multilinear wavelets: A statistical shape space for human faces. In *European Conference on Computer Vision (ECCV)*. 297–312.

Chen Cao, Tomas Simon, Jin Kyu Kim, Gabe Schwartz, Michael Zollhoefer, Shun-Suke Saito, Stephen Lombardi, Shih-En Wei, Danielle Belko, Shoou-I Yu, Yaser Sheikh, and Jason Saragih. 2022. Authentic Volumetric Avatars from a Phone Scan. 41, 4, Article 163 (jul 2022), 19 pages. https://doi.org/10.1145/3528223.3530143

Anpei Chen, Z. Chen, Guli Zhang, Ziheng Zhang, Kenny Mitchell, and Jingyi Yu. 2019. Photo-Realistic Facial Details Synthesis From Single Image. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)* (2019), 9428–9438.

Hongbo Chen, Rui Zheng, Edmond Lou, and Lawrence H Le. 2020. Compact and Wireless Freehand 3D Ultrasound Real-time Spine Imaging System: A pilot study. In *2020 42nd Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC)*. IEEE, 2105–2108.

Byungkuk Choi, Haekwang Eom, Benjamin Mouscadet, Stephen Cullingford, Kurt Ma, Stefanie Gassel, Suzi Kim, Andrew Moffat, Millicent Maier, Marco Revelant, Joe Letteri, and Karan Singh. 2022. Animatomy: An Animator-Centric, Anatomically Inspired System for 3D Facial Modeling, Animation and Transfer. In *SIGGRAPH Asia 2022 Conference Papers* (Daegu, Republic of Korea) *(SA '22)*. Association for Computing Machinery, New York, NY, USA, Article 16, 9 pages. https://doi.org/10.1145/3550469.3555398

Daniel Cudeiro, Timo Bolkart, Cassidy Laidlaw, Anurag Ranjan, and Michael Black. 2019. Capture, Learning, and Synthesis of 3D Speaking Styles. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*. 10101–10111.

Hang Dai, Nick Pears, William Smith, and Christian Duncan. 2020. Statistical Modeling of Craniofacial Shape and Texture. *International Journal of Computer Vision* 128 (02 2020). https://doi.org/10.1007/s11263-019-01260-7

Radek Danecek, Michael J. Black, and Timo Bolkart. 2022. EMOCA: Emotion Driven Monocular Face Capture and Animation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*. 20311–20322.

Paul Debevec, Tim Hawkins, Chris Tchou, Haarm-Pieter Duiker, Westley Sarokin, and Mark Sagar. 2000. Acquiring the Reflectance Field of a Human Face *(SIGGRAPH '00)*. ACM Press/Addison-Wesley Publishing Co., USA, 145–156. https://doi.org/10.1145/344779.344855

Paul E. Debevec. 2012. The Light Stages and Their Applications to Photoreal Digital Actors. In *International Conference on Computer Graphics and Interactive Techniques*.

Yu Deng, Jiaolong Yang, Sicheng Xu, Dong Chen, Yunde Jia, and Xin Tong. 2019. Accurate 3d face reconstruction with weakly-supervised learning: From single image to image set. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 0–0.

Mingsong Dou, S. Khamis, Yu.G. Degtyarev, Philip L. Davidson, S. Fanello, Adarsh Kowdle, Sergio Orts, Christoph Rhemann, David Kim, Jonathan Taylor, Pushmeet Kohli, Vladimir Tankovich, and Shahram Izadi. 2016. Fusion4D: real-time performance capture of challenging scenes. *ACM Trans. Graph.* 35 (2016), 114:1–114:13.

Fuqing Duan, Donghua Huang, Yun Tian, Ke Lu, Zhongke Wu, and Mingquan Zhou. 2015. 3D face reconstruction from skull by regression modeling in shape parameter spaces. *Neurocomputing* 151 (2015), 674–682. https://doi.org/10.1016/j.neucom.2014.04.089

Fuqing Duan, Sen Yang, Donghua Huang, Yongli Hu, zk wu, and Mingquan Zhou. 2013. Craniofacial reconstruction based on multi-linear subspace analysis. *Multimedia Tools and Applications* 73 (11 2013). https://doi.org/10.1007/s11042-012-1351-2

Abhishek Dutta. 2010. *Face Shape and Reflectance Acquisition using a Multispectral Light Stage*. Ph. D. Dissertation. University of York.

Bernhard Egger, William AP Smith, Ayush Tewari, Stefanie Wuhrer, Michael Zollhoefer, Thabo Beeler, Florian Bernard, Timo Bolkart, Adam Kortylewski, Sami Romdhani, et al. 2020. 3d morphable face models—past, present, and future. *ACM Transactions on Graphics (TOG)* 39, 5 (2020), 1–38.

Paul Ekman and Wallace V Friesen. 2002. *Facial action coding system*. Vol. 2. Consulting Psychologists Press.

Yingruo Fan, Zhaojiang Lin, Jun Saito, Wenping Wang, and Taku Komura. 2022. FaceFormer: Speech-Driven 3D Facial Animation with Transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

Z Fan, X Liu, Y Guo, Y Yan, and B Chen. 2019. Audio-Driven Facial Animation System. *ACM Transactions on Graphics* 38, 4 (2019), 1–11.

Yao Feng, Haiwen Feng, Michael J. Black, and Timo Bolkart. 2021. Learning an Animatable Detailed 3D Face Model from In-The-Wild Images. *ACM Transactions on Graphics, (Proc. SIGGRAPH)* 40, 8. https://doi.org/10.1145/3450626.3459936

Guy Gafni, Justus Thies, Michael Zollhofer, and Matthias Nießner. 2021. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8649–8658.

Xuan Gao, Chenglai Zhong, Jun Xiang, Yang Hong, Yudong Guo, and Juyong Zhang. 2022b. Reconstructing Personalized Semantic Facial NeRF Models from Monocular Video. *ACM Trans. Graph.* 41, 6, Article 200 (nov 2022), 12 pages.

Yuchong Gao, Hongye Zeng, Jianhao Zhao, Mingbo Zhang, and Rui Zheng. 2022a. 3D Ultrasound Parametric Modeling Imaging for Spine Deformity – A Preliminary Study. In *2022 IEEE International Ultrasonics Symposium (IUS)*. 1–4. https://doi.org/10.1109/IUS54386.2022.9958051

Pablo Garrido, Michael Zollhöfer, Chenglei Wu, Derek Bradley, Patrick Pérez, Thabo Beeler, and Christian Theobalt. 2016. Corrective 3D Reconstruction of Lips from Monocular Video. *ACM Trans. Graph.* 35, 6, Article 219 (dec 2016), 11 pages. https://doi.org/10.1145/2980179.2982419

Baris Gecer, Stylianos Ploumpis, Irene Kotsia, and Stefanos Zafeiriou. 2019. GANFIT: Generative Adversarial Network Fitting for High Fidelity 3D Face Reconstruction. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 1155–1164. https://doi.org/10.1109/CVPR.2019.00125

Jason Geng. 2011. Structured-light 3D surface imaging: a tutorial. *Advances in Optics and Photonics* 3, 2 (2011), 128–160.

Simon Giebenhain, Tobias Kirschstein, Markos Georgopoulos, Martin Rünz, Lourdes Agapito, and Matthias Nießner. 2022. Learning Neural Parametric Head Models. *arXiv preprint arXiv:2212.02761* (2022).

Michael Goesele, Brian Curless, and Steven M. Seitz. 2006. Multi-View Stereo Revisited. *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)* 2 (2006), 2402–2409.

Yulan Guo, Xiaoming Li, Jie Cheng, and Chi-Keung Yang. 2018. Video-Driven Facial Animation. *ACM Transactions on Graphics (TOG)* 37, 4 (2018), 1–12.

Yana Hasson, Gül Varol, Dimitrios Tzionas, Igor Kalevatykh, Michael J. Black, Ivan Laptev, and Cordelia Schmid. 2019. Learning Joint Reconstruction of Hands and Manipulated Objects. In *IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*. 11807–11816. https://doi.org/10.1109/CVPR.2019.01208

Yang Hong, Bo Peng, Haiyao Xiao, Ligang Liu, and Juyong Zhang. 2022. HeadNeRF: A Real-time NeRF-based Parametric Head Model. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

Patrik Huber, Guosheng Hu, Rafael Tena, Pouria Mortazavian, P Koppen, William J Christmas, Matthias Ratsch, and Josef Kittler. 2016. A multiresolution 3d morphable face model and fitting framework. In *Proceedings of the 11th international joint conference on computer vision, imaging and computer graphics theory and applications*.

Maria Ida Iacono, Esra Neufeld, Esther Akinnagbe, Kelsey Bower, Johanna Wolf, Ioannis Vogiatzis Oikonomidis, Deepika Sharma, Bryn Lloyd, Bertram J. Wilm, Michael Wyss, Klaas P. Pruessmann, Andras Jakab, Nikos Makris, Ethan D. Cohen, Niels Kuster, Wolfgang Kainz, and Leonardo M. Angelone. 2015. MIDA: A Multimodal Imaging-Based Detailed Anatomical Model of the Human Head and Neck. *PLOS ONE* 10, 4 (04 2015), 1–35. https://doi.org/10.1371/journal.pone.0124126

Lisa Jackson, Paulo Esteves, Ben Sussmilch, Yuhong Zhang, and Costas Iliopoulos. 2017. Automatic facial expression recognition based on FACS. *IEEE Transactions on Affective Computing* 8, 2 (2017), 201–213.

AK Jain, SZ Li, and Y Chen. 2017. The 3D Dynamic Facial Expression dataset. In *Proc. Int'l Conf. on Computer Vision.*

Hanbyul Joo, Tomas Simon, and Yaser Sheikh. 2018. Total Capture: A 3D Deformation Model for Tracking Faces, Hands, and Bodies. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).*

Kolja Kähler, Jörg Haber, and Hans-Peter Seidel. 2001. Geometry-Based Muscle Modeling for Facial Animation. In *Proceedings of Graphics Interface 2001* (Ottawa, Ontario, Canada) *(GI '01)*. Canadian Information Processing Society, CAN, 37–46.

Hyeongwoo Kim, Pablo Garrido, Ayush Tewari, Weipeng Xu, Justus Thies, Matthias Niessner, Patrick Pérez, Christian Richardt, Michael Zollhöfer, and Christian Theobalt. 2018. Deep video portraits. *ACM Transactions on Graphics (TOG)* 37, 4 (2018), 1–14.

Martin Klaudiny and Adrian Hilton. 2012. High-Detail 3D Capture and Non-sequential Alignment of Facial Performance. In *2012 Second International Conference on 3D Imaging, Modeling, Processing, Visualization & Transmission.* 17–24. https://doi.org/10.1109/3DIMPVT.2012.67

Samuli Laine, Tero Karras, Timo Aila, Antti Herva, Shunsuke Saito, Ronald Yu, Hao Li, and Jaakko Lehtinen. 2017. Production-Level Facial Performance Capture Using Deep Convolutional Neural Networks. In *Proceedings of the ACM SIGGRAPH / Eurographics Symposium on Computer Animation* (Los Angeles, California) *(SCA '17)*. Association for Computing Machinery, New York, NY, USA, Article 10, 10 pages. https://doi.org/10.1145/3099564.3099581

Alexandros Lattas, Stylianos Moschoglou, Stylianos Ploumpis, Baris Gecer, Abhijeet Ghosh, and Stefanos Zafeiriou. 2022. AvatarMe++: Facial Shape and BRDF Inference With Photorealistic Rendering-Aware GANs. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44, 12 (2022), 9269–9284. https://doi.org/10.1109/TPAMI.2021.3125598

Sung-Hee Lee, Eftychios Sifakis, and Demetri Terzopoulos. 2009. Comprehensive Biomechanical Modeling and Simulation of the Upper Body. 28, 4, Article 99 (sep 2009), 17 pages. https://doi.org/10.1145/1559755.1559756

Sung-Hee Lee and Demetri Terzopoulos. 2006. Heads up! Biomechanical Modeling and Neuromuscular Control of the Neck. *ACM Trans. Graph.* 25, 3 (jul 2006), 1188–1198. https://doi.org/10.1145/1141911.1142013

Marc Levoy, Kari Pulli, Brian Curless, Szymon Rusinkiewicz, David Koller, Lucas Pereira, Matt Ginzton, Sean Anderson, James Davis, Jeremy Ginsberg, Jonathan Shade, Duane Fulk, and Cyberware Inc. 2001. The Digital Michelangelo Project: 3D Scanning of Large Statues. *Proceedings of the ACM SIGGRAPH Conference on Computer Graphics* 1 (03 2001).

Duo Li, Shinjiro Sueda, Debanga R. Neog, and Dinesh K. Pai. 2013. Thin Skin Elastodynamics. 32, 4, Article 49 (jul 2013), 10 pages. https://doi.org/10.1145/2461912.2462008

Hao Li, Thibaut Weise, and Mark Pauly. 2010. Example-based facial rigging. *Acm transactions on graphics (tog)* 29, 4 (2010), 1–6.

Jiaman Li, Zhengfei Kuang, Yajie Zhao, Mingming He, Karl Bladin, and Hao Li. 2020b. Dynamic facial asset and rig generation from a single scan. *ACM Trans. Graph.* 39, 6 (2020), 215–1.

Ruilong Li, Karl Bladin, Yajie Zhao, Chinmay Chinara, Owen Ingraham, Pengda Xiang, Xinglei Ren, Pratusha Prasad, Bipin Kishore, Jun Xing, et al. 2020a. Learning formation of physically-based face attributes. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition.* 3410–3419.

Tianye Li, Timo Bolkart, Michael. J. Black, Hao Li, and Javier Romero. 2017. Learning a model of facial shape and expression from 4D scans. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)* 36, 6 (2017), 194:1–194:17. https://doi.org/10.1145/3130800.3130813

Yuwei Li, Minye Wu, Yuyao Zhang, Lan Xu, and Jingyi Yu. 2021. PIANO: A Parametric Hand Bone Model from Magnetic Resonance Imaging. In *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21.* 816–822. https://doi.org/10.24963/ijcai.2021/113

X Liu, Y Liu, Y Chen, Z Fan, Y Yan, and B Chen. 2019. Audio-Driven Facial Animation Framework. *ACM Transactions on Graphics* 38, 4 (2019), 1–12.

Yilong Liu, Chengwei Zheng, Feng Xu, Xin Tong, and Baining Guo. 2021. Data-Driven 3D Neck Modeling and Animation. *IEEE Transactions on Visualization and Computer Graphics* 27, 7 (2021), 3226–3237. https://doi.org/10.1109/TVCG.2020.2967036

Stephen Lombardi, Jason Saragih, Tomas Simon, and Yaser Sheikh. 2018. Deep Appearance Models for Face Rendering. 37, 4, Article 68 (jul 2018), 13 pages. https://doi.org/10.1145/3197517.3201401

Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. 2019. Neural volumes. *ACM Transactions on Graphics* 38, 4 (jul 2019), 1–14. https://doi.org/10.1145/3306346.3323020

Stephen Lombardi, Tomas Simon, Gabriel Schwartz, Michael Zollhoefer, Yaser Sheikh, and Jason Saragih. 2021. Mixture of volumetric primitives for efficient neural rendering. *ACM Transactions on Graphics (TOG)* 40, 4 (2021), 1–13.

M Loper, M Maher, V Choutas, G Pons-Moll, and M Black. 2015. SMPL: A Skinned Multi-Person Linear Model. *ACM Transactions on Graphics* (2015).

Haimin Luo, Anpei Chen, Qixuan Zhang, Bai Pang, Minye Wu, Lan Xu, and Jingyi Yu. 2021. Convolutional neural opacity radiance fields. In *2021 IEEE International Conference on Computational Photography (ICCP).* IEEE, 1–12.

Zhiping Luo, Nicolas Pronost, and Arjan Egges. 2013. Physically-based Human Neck Simulation.

LPP Ma, PJ Phillips, and PJ Flynn. 2006. Facial reflectance map estimation using photometric stereo. In *Proc. Int'l Conf. on Computer Vision.*

Muammar Mattar and Xinbo Gao. 2017. FACS-based facial expression recognition using convolutional neural networks. In *2017 15th IEEE International Conference on Automatic Face and Gesture Recognition (FG).* IEEE, 86–93.

Muammar Mattar, Xun Yang, and Xinbo Gao. 2015. FACS-based facial expression analysis using PCA and LDA. In *2015 13th IEEE International Conference on Automatic Face and Gesture Recognition (FG).* IEEE, 1–8.

Masaya Misaki, Jonathan Savitz, Vadim Zotev, Raquel Phillips, Han Yuan, Kymberly D. Young, Wayne C. Drevets, and Jerzy Bodurka. 2015. Contrast enhancement by combining T1- and T2-weighted structural brain MR images. *Magnetic Resonance in Medicine* 74, 6 (2015), 1609–1620. https://doi.org/10.1002/mrm.25560 arXiv:https://onlinelibrary.wiley.com/doi/pdf/10.1002/mrm.25560

Giljoo Nam, Chenglei Wu, Min H. Kim, and Yaser Sheikh. 2019. Strand-Accurate Multi-View Hair Capture. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2019), 155–164.

Thomas R. Nelson and Dolores H. Pretorius. 1998. Three-dimensional ultrasound imaging. *Ultrasound in Medicine & Biology* 24, 9 (1998), 1243–1270. https://doi.org/10.1016/S0301-5629(98)00043-X

Ahmed A A Osman, Timo Bolkart, and Michael J. Black. 2020. STAR: A Sparse Trained Articulated Human Body Regressor. In *European Conference on Computer Vision (ECCV).* 598–613. https://star.is.tue.mpg.de

Ahmed A. A. Osman, Timo Bolkart, Dimitrios Tzionas, and Michael J. Black. 2022. SUPR: A Sparse Unified Part-Based Human Representation. In *European Conference on Computer Vision (ECCV).* Springer International Publishing.

Pascal Paysan, Reinhard Knothe, Brian Amberg, Sami Romdhani, and Thomas Vetter. 2009. A 3D Face Model for Pose and Illumination Invariant Face Recognition. In *2009 Sixth IEEE International Conference on Advanced Video and Signal Based Surveillance.* 296–301. https://doi.org/10.1109/AVSS.2009.58

PJ Phillips, PJ Flynn, KW Bowyer, and H Schott. 2008. The FRGC 2.0 dataset. In *Proc. Int'l Conf. on Computer Vision.*

Stylianos Ploumpis, Evangelos Ververas, Eimear O'Sullivan, Stylianos Moschoglou, Haoyang Wang, Nick Pears, William AP Smith, Baris Gecer, and Stefanos Zafeiriou. 2020. Towards a complete 3D morphable model of the human head. *IEEE transactions on pattern analysis and machine intelligence* 43, 11 (2020), 4142–4160.

Stylianos Ploumpis, Haoyang Wang, Nick Pears, William A. P. Smith, and Stefanos Zafeiriou. 2019. Combining 3D Morphable Models: A Large Scale Face-And-Head Model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR).*

Emily B. Prince, Katherine B. Martin, and Daniel S. Messinger. 2015. Facial Action Coding System.

Neng Qian, Jiayi Wang, Franziska Mueller, Florian Bernard, Vladislav Golyanik, and Christian Theobalt. 2020. Html: A parametric hand texture model for 3d hand reconstruction and personalization. In *European Conference on Computer Vision.* Springer, 54–71.

Zesong Qiu, Yuwei Li, Dongming He, Qixuan Zhang, Longwen Zhang, Yinghao Zhang, Jingya Wang, Lan Xu, Xudong Wang, Yuyao Zhang, and Jingyi Yu. 2022. SCULPTOR: Skeleton-Consistent Face Creation Using a Learned Parametric Generator. *ACM Trans. Graph.* 41, 6, Article 213 (nov 2022), 17 pages. https://doi.org/10.1145/3550454.3555462

Amit Raj, Michael Zollhofer, Tomas Simon, Jason Saragih, Shunsuke Saito, James Hays, and Stephen Lombardi. 2021. Pixel-aligned volumetric avatars. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 11733–11742.

Javier Romero, Dimitrios Tzionas, and Michael J. Black. 2017. Embodied Hands: Modeling and Capturing Hands and Bodies Together. *ACM Transactions on Graphics, (Proc. SIGGRAPH Asia)* 36, 6 (Nov. 2017).

Mohammad Samim. 2021. 3D MRI Models of the Musculoskeletal System. *Seminars in musculoskeletal radiology* 25 3 (2021), 388–396.

Janis Savlovskis. 2022. Anatomy Standard. https://www.anatomystandard.com/biomechanics/spine/rom-of-vertebrae.html.

Steven M. Seitz, Brian Curless, James Diebel, Daniel Scharstein, and Richard Szeliski. 2006. A Comparison and Evaluation of Multi-View Stereo Reconstruction Algorithms. *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)* 1 (2006), 519–528.

William A. P. Smith, Alassane Seck, Hannah Dee, Bernard Tiddeman, Joshua Tenenbaum, and Bernhard Egger. 2020. A Morphable Face Albedo Model. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).* 5011–5020.

Joanna M. Stephen, James DF Calder, Andy Williams, and Hadi El Daou. 2021. Comparative accuracy of lower limb bone geometry determined using MRI, CT, and direct bone 3D models. *Journal of Orthopaedic Research* 39, 9 (2021), 1870–1876. https://doi.org/10.1002/jor.24923

arXiv:https://onlinelibrary.wiley.com/doi/pdf/10.1002/jor.24923

3D Scan Store. 2022. https://www.3dscanstore.com/.

Giota Stratou, Abhijeet Ghosh, Paul Debevec, and Louis-Philippe Morency. 2011. Effect of illumination on automatic expression recognition: a novel 3D relightable facial database. In *2011 IEEE International Conference on Automatic Face & Gesture Recognition (FG)*. IEEE, 611–618.

Justus Thies, Michael Zollhöfer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. 2016. Face2Face: Real-Time Face Capture and Reenactment of RGB Videos. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2387–2395. https://doi.org/10.1109/CVPR.2016.262

Redha Touati, William Trung Le, and Samuel Kadoury. 2021. A feature invariant generative adversarial network for head and neck MRI/CT image synthesis. *Physics in Medicine & Biology* 66, 9 (apr 2021), 095001. https://doi.org/10.1088/1361-6560/abf1bb

Maarten van Eerd, Jacob Patijn, Judith M Sieben, Mischa Sommer, Jan Van Zundert, Maarten van Kleef, and Arno Lataster. 2014. Ultrasonography of the cervical spine: an in vitro anatomical validation model. *Anesthesiology* 120, 1 (2014), 86–96.

Barbara Villarini, Athanasios Gkelias, and Vasilios Argyriou. 2017. Photometric Stereo for 3D Face Reconstruction Using Non Linear Illumination Models. In *Multimodal Pattern Recognition of Social Signals in Human-Computer-Interaction*, Friedhelm Schwenker and Stefan Scherer (Eds.). Springer International Publishing, Cham, 140–152.

Daoye Wang, Prashanth Chandran, Gaspard Zoss, Derek Bradley, and Paulo Gotardo. 2022a. MoRF: Morphable Radiance Fields for Multiview Neural Head Modeling. In *ACM SIGGRAPH 2022 Conference Proceedings* (Vancouver, BC, Canada) *(SIGGRAPH '22)*. Association for Computing Machinery, New York, NY, USA, Article 55, 9 pages. https://doi.org/10.1145/3528233.3530753

L Wang, X Chen, Y Liu, Y Fang, and Y Chen. 2018. Ultrasonic range finding for custom-fit neck brace design. In *Proc. Int'l Conf. on Computer Vision*.

Lizhen Wang, Zhiyuan Chen, Tao Yu, Chenguang Ma, Liang Li, and Yebin Liu. 2022b. FaceVerse: a Fine-grained and Detail-controllable 3D Face Morphable Model from a Hybrid Dataset. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 20333–20342.

Lijuan Wang, Hongxun Zhang, Qionghai Fu, Xing Li, and Hanqing Lu. 2017. Realistic facial animation using blend shapes and skeletal deformation. *IEEE Transactions on Visualization and Computer Graphics* 23, 12 (2017), 2640–2651.

Xueying Wang, Yudong Guo, Bailin Deng, and Juyong Zhang. 2020a. Lightweight Photometric Stereo for Facial Details Recovery. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

Xing Wang, Yida Liu, Zhiwei Wang, Xudong Wang, Yajie Hu, Yiliang Chen, and Hubert P. H. Shum. 2020b. VDFA-S: A Video-Driven Facial Animation System. *ACM Transactions on Graphics (TOG)* 39, 4 (2020), 1–12.

Ziyan Wang, Giljoo Nam, Tuur Stuyck, Stephen Lombardi, Michael Zollhöfer, Jessica Hodgins, and Christoph Lassner. 2022c. HVH: Learning a Hybrid Neural Volumetric Representation for Dynamic Hair Performance Capture. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6143–6154.

Andreas Wenger, Andrew Gardner, Chris Tchou, Jonas Unger, Tim Hawkins, and Paul Debevec. 2005. Performance Relighting and Reflectance Transformation with Time-Multiplexed Illumination. *ACM Trans. Graph.* 24, 3 (jul 2005), 756–764. https://doi.org/10.1145/1073204.1073258

Tim Weyrich, Wojciech Matusik, Hanspeter Pfister, Bernd Bickel, Craig Donner, Chien Tu, Janet McAndless, Jinho Lee, Addy Ngan, Henrik Wann Jensen, and Markus Gross. 2006. Analysis of Human Faces Using a Measurement-Based Skin Reflectance Model. *ACM Trans. Graph.* 25, 3 (jul 2006), 1013–1024. https://doi.org/10.1145/1141911.1141987

Sebastian Winberg, Gaspard Zoss, Prashanth Chandran, Paulo Gotardo, and Derek Bradley. 2022. Facial Hair Tracking for High Fidelity Performance Capture. *ACM Trans. Graph.* 41, 4, Article 165 (jul 2022), 12 pages.

Bernhard P. Wrobel. 2001. Multiple View Geometry in Computer Vision. *Künstliche Intell.* 15 (2001), 41.

Chenglei Wu, Derek Bradley, Pablo Garrido, Michael Zollhöfer, Christian Theobalt, Markus Gross, and Thabo Beeler. 2016. Model-Based Teeth Reconstruction. *ACM Trans. Graph.* 35, 6, Article 220 (dec 2016), 13 pages. https://doi.org/10.1145/2980179.2980233

Cheng-hsin Wuu, Ningyuan Zheng, Scott Ardisson, Rohan Bali, Danielle Belko, Eric Brockmeyer, Lucas Evans, Timothy Godisart, Hyowon Ha, Alexander Hypes, Taylor Koska, Steven Krenn, Stephen Lombardi, Xiaomin Luo, Kevyn McPhail, Laura Millerschoen, Michal Perdoch, Mark Pitts, Alexander Richard, Jason Saragih, Junko Saragih, Takaaki Shiratori, Tomas Simon, Matt Stewart, Autumn Trimble, Xinshuo Weng, David Whitewolf, Chenglei Wu, Shoou-I Yu, and Yaser Sheikh. 2022. Multiface: A Dataset for Neural Face Rendering. In *arXiv*. https://doi.org/10.48550/ARXIV.2207.11243

Hongyi Xu, Eduard Gabriel Bazavan, Andrei Zanfir, William T. Freeman, Rahul Sukthankar, and Cristian Sminchisescu. 2020. GHUM & GHUML: Generative 3D Human Shape and Articulated Pose Models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

Haotian Yang, Hao Zhu, Yanru Wang, Mingkai Huang, Qiu Shen, Ruigang Yang, and Xun Cao. 2020. Facescape: a large-scale high quality 3d face dataset and detailed riggable 3d face prediction. In *Proceedings of the ieee/cvf conference on computer vision and pattern recognition*. 601–610.

Tarun Yenamandra, Ayush Tewari, Florian Bernard, Hans-Peter Seidel, Mohamed Elgharib, Daniel Cremers, and Christian Theobalt. 2021. i3DMM: Deep Implicit 3D Morphable Model of Human Heads. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12803–12813.

Stefanos Zafeiriou, Gary A. Atkinson, Mark F. Hansen, William A. P. Smith, Vasileios Argyriou, Maria Petrou, Melvyn L. Smith, and Lyndon N. Smith. 2013. Face Recognition and Verification Using Photometric Stereo: The Photoface Database and a Comprehensive Evaluation. *IEEE Transactions on Information Forensics and Security* 8, 1 (2013), 121–135. https://doi.org/10.1109/TIFS.2012.2224109

J Zhang, X Chen, Y Liu, Y Fang, and Y Chen. 2018. Ultrasonic imaging for neck movement measurement and analysis. In *Proc. Int'l Conf. on Computer Vision*.

J Zhang, X Chen, Y Liu, Y Fang, and Y Chen. 2019. Ultrasonic imaging for neck movement analysis. In *Proc. Int'l Conf. on Computer Vision*.

Longwen Zhang, Chuxiao Zeng, Qixuan Zhang, Hongyang Lin, Ruixiang Cao, Wei Yang, Lan Xu, and Jingyi Yu. 2022. Video-Driven Neural Physically-Based Facial Asset for Production. *ACM Transactions on Graphics (TOG)* 41 (2022), 1 – 16.

Yiyu Zhuang, Hao Zhu, Xusen Sun, and Xun Cao. 2022. Mofanerf: Morphable facial neural radiance field. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part III*. Springer, 268–285.