



Malek-Podjaski, M. and Deligianni, F. (2022) Towards Explainable, Privacy-Preserved Human-Motion Affect Recognition. In: 2021 IEEE Symposium Series on Computational Intelligence (SSCI), Orlando, FL, USA, 04-07 Dec 2021, ISBN 9781728190488

(doi: [10.1109/SSCI50451.2021.9660129](https://doi.org/10.1109/SSCI50451.2021.9660129))

This is the Author Accepted Manuscript.

© 2021 IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.

There may be differences between this version and the published version. You are advised to consult the publisher's version if you wish to cite from it.

<http://eprints.gla.ac.uk/253811/>

Deposited on: 4 October 2021

Towards Explainable, Privacy-Preserved Human-Motion Affect Recognition¹

Matthew Malek-Podjaski
School of Computing Science
University of Glasgow
Glasgow, United Kingdom
2323841m@student.gla.ac.uk

Fani Deligianni
School of Computing Science
University of Glasgow
Glasgow, United Kingdom
fani.deligianni@glasgow.ac.uk

Abstract—77Human motion characteristics are used to monitor the progression of neurological diseases and mood disorders. Since perceptions of emotions are also interleaved with body posture and movements, emotion recognition from human gait can be used to quantitatively monitor mood changes. Many existing solutions often use shallow machine learning models with raw positional data or manually extracted features to achieve this. However, gait is composed of many highly expressive characteristics that can be used to identify human subjects, and most solutions fail to address this, disregarding the subject’s privacy. This work introduces a novel deep neural network architecture to disentangle human emotions and biometrics. In particular, we propose a cross-subject transfer learning technique for training a multi-encoder autoencoder deep neural network to learn disentangled latent representations of human motion features. By disentangling subject biometrics from the gait data, we show that the subject’s privacy is preserved while the affect recognition performance outperforms traditional methods. Furthermore, we exploit Guided Grad-CAM to provide global explanations of the model’s decision across gait cycles. We evaluate the effectiveness of our method to existing methods at recognizing emotions using both 3D temporal joint signals and manually extracted features. We also show that this data can easily be exploited to expose a subject’s identity. Our method shows up to 7% improvement and highlights the joints with the most significant influence across the average gait cycle.

Index Terms—human motion analysis, privacy, disentanglement, affect, deep learning

I. INTRODUCTION

Strong neuroscientific evidence [1] shows that there is an interaction between brain networks involved in gait and emotion. It does not come as a surprise that certain gait characteristics are related to mood disorders, such as depression and anxiety. Gait analysis has revealed several indices that correlate to emotional well-being. For example, increased gait speed, step length and arm swing have been related to positive emotions. In contrast, a low gait initiation reaction time and flexion of posture have been linked to negative feelings [2], [3].

However, developing robust algorithms for human motion and gait analysis requires monitoring patients continuously at their home. This poses tremendous ethical and privacy challenges that have to be addressed in order for the technology to be successfully accepted and adapted. Privacy-preserving

deep learning technologies for computer vision applications is a relatively new research area that is usually based on disentangling biometric features representation from other attributes [4], [5]. Some old-fashion techniques of face obfuscation such as pixelation, blurring and masking to protect privacy offer limited protection while they compromise the ability to track human motion precisely. For example, inference in deep learning models is compromised if there is an uncertainty on whether the view is anterior or posterior with relation to the human’s head. Furthermore, facial features also encode relevant information, such as the subject’s mood.

To disentangle human biometrics and affects, we exploit a Multi-encoder Autoencoder structure that achieves cross-subject transfer learning as inspired by recent works by Gu et al. [6], who highlighted their efficiency for disentanglement learning in abnormal gait recognition based on joints movements of the lower limbs. This concept is also similar to motion retargeting, which allows transferring motions from one subject to another by disentangling motion, skeleton, and camera view data [7]. To our knowledge, this is the first study that disentangles affects and biometrics and demonstrates that this approach preserves subjects privacy while it improves classification performance.

Furthermore, the ability to provide insight into which factors are affecting the output and by how much is crucial for developing explainable machine learning models and a necessary condition in evaluating and accepting these models in healthcare applications [8]. We choose a gradient-based explainability method, which is inherently model-specific and local, to obtain intuitive visualisation of the deep neural network. Subsequently, we aggregate information across testing samples in each class which allows us to inspect global explanations of the model.

Overall, the contribution of our work is threefold: a) we introduce a novel method to disentangle affects and biometrics from gait data, which significantly improves the performance of the affect recognition model, b) we demonstrate that our architecture of disentangling affect and biometrics enhances subject’s privacy considerably and c) we exploit a gradient-based neural network visualisation method to examine global explanations of the model across classes. This reveals the joints that contribute significantly to the prediction and identify

¹ Authors acknowledge funding from EPSRC EP/R045178/1

² Codebase will be released upon acceptance.

which phases of the gait cycle are most informative.

II. RELATED WORK

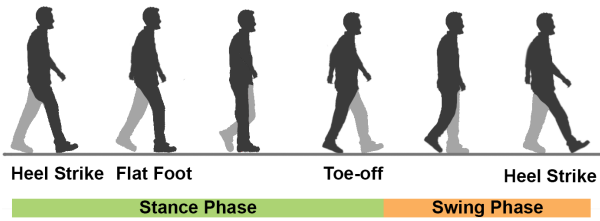


Fig. 1. An illustration of a single gait cycle, highlighting key events of the right foot and the two phases (stance and swing) that compose the gait cycle.

A. Gait Cycle

Human gait involves repetitive sets of motions, iterating with each step. When analysing such motions we confine the descriptions to a single gait cycle, defined as the period between two consecutive heel strikes of the same heel [9]. A single gait cycle is composed of two main phases, the stance phase, and the swing phase. With the stance phase being composed of the single limb stance, during which only one foot is touching the ground, and a double support where both feet are touching the ground [10]. The stance phase begins as the heel strikes the ground and ends with the toe-off, when the foot lifts off the ground. Subsequently, the gait enters into the swing phase, the phase with the foot off the ground, which is broken down into the initial swing, mid swing and the terminal swing. The latter happens as the heel strikes the ground for the second time. A visual representation of the gait cycle along with the key gait events is shown in Figure 1.

B. Emotion within Gait

Embodiment theory as described by Winkelman et al. [1], "is the idea that higher-level processing is grounded in the organism's sensory and motor experiences", suggesting that the perception of emotions and cognitions is interleaved with body postures and movements. For example, when people recall joyful memories, they partially reproduce that state. Embodiment theory suggests that such re-enactments are crucial for identifying and interpreting the differences between affects. Most importantly, they distinguish that such re-enactments are not always conscious. This concept builds the foundation that human motions can be used to identify internal emotional states.

C. Gait Emotion Characteristics

Extensive research has been done into manually extracted features from human motion in order to understand how humans perceive emotion from body language. Studies such as those by Gross et al. [11], show that when a subject's specific target emotions are stimulated during their gait, those same emotions are subconsciously re-enacted through the subject's motions, and could be identified reliably by human observers at a ratio greater than chance. The same study also identified an exhaustive list of quantitative gait metrics and their importance

to each emotion. In particular, metrics such as: velocity, stride length, hip/shoulder/elbow flexion, and pelvis/trunk rotation were shown to have statistically significant differences between emotions. Similar studies [12]–[14] found stride lengths and gait speed to be a significant characteristic, as well as increased arm swing, thigh elevation angles and cadence being strongly correlated with happiness and anger.

D. Gait Subject Biometrics

Aside from neurological conditions, gait is also highly expressive regarding subjects' identity. In fact, subject's gait styles usually dominate over any subtle neurological effect. There is significant research on identifying subjects and their gender based on their gait characteristics [15]. We highlight some of these aspects in our work and show that we can easily distinguish between subjects' identity from raw gait sensor data.

E. Manual Feature Extraction

Much of human motion analysis literature focuses on analysing manually extracted features with statistical analysis or use standard machine learning models such as Support Vector Machines (SVM) for classification. Since extracted features can be strong indicators of emotion [11], feature engineering has been used to find an optimal set of gait characteristics that relate with emotions. Studies have identified emotions by computing image moments from gait images [16], extracting temporal, frequency, and temporal-frequency domain features from wearable accelerometer data [17], and using similarity indices to compare sets of features extracted based on dimensionality reduction techniques such as Principle Component Analysis (PCA) [18]. On the other hand, there are also studies that use pose-based gait data to extract angular or velocity features. Srivastava et al. [19] extracted angular features such as head tilt, shoulder, right elbow and right knee, and reported high accuracy with shallow machine learning models.

F. Deep Learning Models

Deep learning models have become increasingly popular due to their ability to automatically extract optimum features that have the potential to capture better the complexity of the gait data than predefined manually extracted features. Commonly applied deep learning models within this domain are Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), such as the Long Short-Term Memory (LSTM) models, which have been used for gait recognition and gender identification [20]–[22], as well as for the detection of abnormalities in a subject's gait [23].

G. Privacy-Preserved Human Motion Analysis

When exploiting human biological data, there is always a privacy infringement concern even when the intended processing of the data does not aim to extract subject-specific information. Deep learning models can memorise training data and this information leakage can be exploited to reveal subject

biometrics [24], [25]. For example, Song and Shmatikov [26] introduces the concept of overlearning where a model trained with the intent of being a gender binary classifier it also unintentionally learns to recognize subjects based on identity and race, although race is not explicitly defined. Therefore, it has become evident that AI approaches should be designed with privacy and ethical concerns taking into consideration [27]. This problem is particularly profound in healthcare applications. Therefore, it is important to develop algorithms that can dynamically filter streams of data to eliminate unnecessary information that can compromise privacy. Our approach, is using a modular architecture, with a part of it using identity information to guide the learning process to disentangle affects and biometrics. In this way, a part of the network is focused on affect recognition while identify information is inherently filtered out in a measurable manner.

Recent work has demonstrated that it is possible to separate the latent representation of emotions from the identity representation of a face in a stack of images [28]–[31]. Typically this method is used to synthesise new data by sampling the latent space representations. Here we suggest exploiting this deep learning approach to disentangle subtle gait features related to affects from the noisy subject-specific features to generate a latent representations of emotions, reducing the complexity of the dataset and enhancing the subjects’ privacy in the process. Disentanglement in representation learning refers to the ability to break down data features in key categories and represent them in their own latent spaces. This method is particularly successful in computer vision applications.

Typically to achieve this sort of subject-data disentanglement, a cross-subject training approach is required, where latent representations of the subject data and the target disentanglement data are learned and transferred between each other for different subjects. This type of training has been shown to generalise well and it has been applied in a range of problems, from disentangling identity from human faces [32] to disentangling pose, expression, and illumination from faces [33], allowing the model to synthesise a face from different orientations with different expressions and lighting conditions. For human motion applications, this methodology was used to disentangle human pose from images [34], to estimate a 3D human mesh from 2D images [35] by learning to disentangle the skeleton from the image, and even for a predictive model of locomotion [36]. The major strength of the technique is that it reduces complexity to a low dimension latent representation.

H. Explainability in Deep Learning Models

Deep learning models do not provide direct explanations, and therefore their predictions are difficult to understand and trust, which hinders their adaptation into safety-critical applications. For these reasons, the effort to explain deep learning models has been intensified in the last few years. Explainability methods are categorised into local and global, depending on whether they provide explanations for individual samples or the entire method/group of samples, respectively.

Furthermore, methods can be model agnostic or model specific depending on whether they can apply to any model or are tailored to the specific model under investigation [8]. Gradient-based methods are both local and model specific, since they exploit the neural network activation’s mechanisms to directly infer the attributions of the network with relation to a decision. We suggest summing this information across testing samples to obtain global explanations for each class.

III. METHODS

A. Cross-subject transfer/Multi-encoder Autoencoder Network

To perform a disentanglement between subject biometrics and affects we use a separate encoder for each of the features, which transfers the low dimensional feature representations between each other to generate new targets that are a cross over of both features. Figure 2(a) shows the implemented architecture, which consists of a subject-specific encoder and an affect-specific encoder. A decoder is trained to use the concatenated input of the subject and affect encodings to reconstruct a gait cycle and transfer the emotional state of one person onto another. In this way, the encoded features are combined during training to generate two different sets of output, a reconstruction batch that matches the original input batch, as well as a cross-subject generated batch that contains new data generated by the autoencoder during training as a result of the cross-subject transfer. Note that for the best results, training using this method requires the ground-truth labels for the cross-subject transfer data to measure the accuracy of the generated data after the latent space transfer.

B. Autoencoder Formulation

Let S and A denote the sets of subject biometrics and affect features that form any given gait cycle. Then let $\{x_{i,j}\} \in \mathbb{R}^{T \times 3J}$ be a motion, where x is a gait cycle described as a subject identity ($i \in S$), and an affect ($j \in A$). T is the temporal duration of the motion, and J is the number of joints each specified by (x, y, z) coordinates.

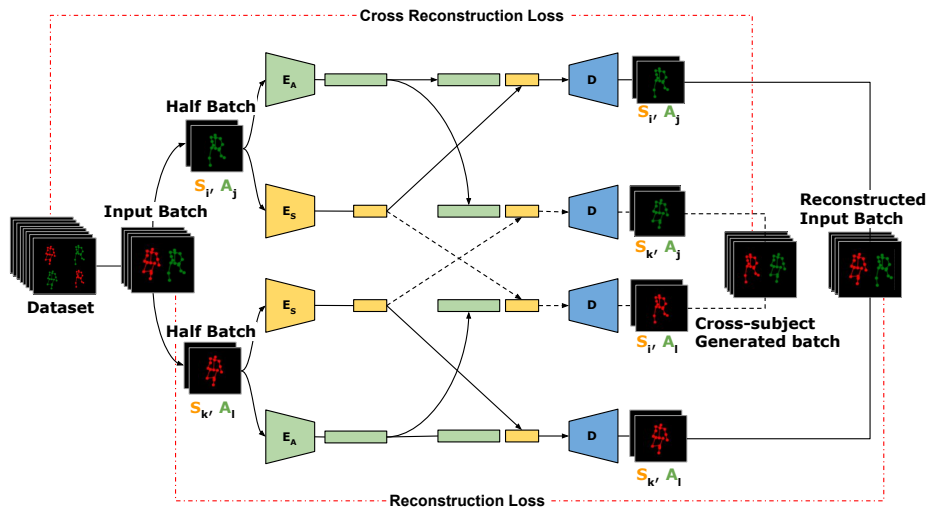
Our aim is to create an architecture that takes a gait cycle and decomposes it into two factors. Towards this end two encoders E_S and E_A are trained simultaneously to decompose the gait cycle into two separate latent representation of subject biometrics and affects, respectively. After this decomposition a gait cycle can be defined as $x_{i,j} = D(s_i, a_j)$, where s_i are subject biometric features encoded by $E_S \mapsto s_i$, and a_j are affect features encoded by $E_A \mapsto a_j$. D is a decoder trained to generate the reconstruction of a decomposed gait cycle, $D : (s_i, a_j) \mapsto \hat{x}_{i,j}$.

Given two data samples, $x_{i,j}, x_{k,l} \in \mathbb{R}^{T \times 3J}$, where $i, k \in S$ and $j, l \in A$, we want to optimise our network such that:

$$\forall_{i,k \in S | j,l \in A} \|D(E_S(x_{k,l}), E_A(x_{i,j})) - x_{k,j}\| \approx 0 \quad (1)$$

C. Loss functions

We use three different types of loss functions to train the autoencoder: a Reconstruction loss, a Cross Reconstruction loss, and a Triplet loss. With a combination of all of them, we can



(a) Model Training Procedure

Network	Layer	Parameters	In/Out
E_S	Conv-1D	$k = 8, s = 2$	45/45
	IC + LeakyReLU	dropout = 0.05	-
	Conv-1D	$k = 8, s = 2$	45/25
	IC + LeakyReLU	dropout = 0.05	-
	Conv-1D	$k = 8, s = 2$	25/15
	IC + LeakyReLU	dropout = 0.05	-
	Conv-1D	$k = 8, s = 2$	15/10
	IC + LeakyReLU	dropout = 0.05	-
	Conv-1D	$k = 8, s = 2$	10/4
E_A	Conv-1D	$k = 8, s = 2$	45/96
	IC + LeakyReLU	dropout = 0.05	-
	Conv-1D	$k = 8, s = 2$	96/128
	IC + LeakyReLU	dropout = 0.05	-
	Conv-1D	$k = 8, s = 2$	128/128
	IC + LeakyReLU	dropout = 0.05	-
	Conv-1D	$k = 8, s = 2$	128/128
	IC + LeakyReLU	dropout = 0.05	-
	Conv-1D	$k = 8, s = 2$	128/64
D	UP + Conv1D	scale = 2, $k = 7, s = 1$	68/128
	IC + LeakyReLU	dropout = 0.05	-
	UP + Conv1D	scale = 2, $k = 7, s = 1$	128/128
	IC + LeakyReLU	dropout = 0.05	-
	UP + Conv1D	scale = 2, $k = 7, s = 1$	128/128
	IC + LeakyReLU	dropout = 0.05	-
	UP + Conv1D	scale = 2, $k = 7, s = 1$	128/96
	IC + LeakyReLU	dropout = 0.05	-
	UP + Conv1D	scale = 2, $k = 7, s = 1$	96/45
C_A	Linear + ReLU	-	16/32
	Dropout + Linear	dropout = 0.5	32/4

(b) Model Architecture

Fig. 2. (a) Illustration of cross-subject training procedure. Given two motions $x_{i,j}$ and $x_{k,l}$, the subject and affect features are extracted, their latent spaces are transferred across each other. Then the transferred and original feature combinations are reconstructed, resulting in four different motions, $\hat{x}_{i,j}$, $\hat{x}_{k,j}$, $\hat{x}_{i,l}$, $\hat{x}_{k,l}$. The reconstructions of the original subject/affect pairings $\hat{x}_{i,j}$ and $\hat{x}_{k,l}$ are compared with the original input $x_{i,j}$ and $x_{k,l}$ to calculate the reconstruction loss, and $\hat{x}_{k,j}$ and $\hat{x}_{i,l}$ are compared to the ground truth motions $x_{k,j}$ and $x_{i,l}$ that exist in the full dataset. (b) Architecture summary of Autoencoder model.

ensure that the reconstructed output is a faithful reconstruction of the original motion while emphasising clustering in the low dimensional encodings.

D. Reconstruction loss

Firstly, we expect the reconstruction produced by the autoencoder to be faithful to the original input. As such for a sample data point from the training set \mathcal{X} , we aim to minimise the difference between the mapped reconstruction of the data sample, and itself using:

$$\mathcal{L}_{\text{rec}} = \mathbb{E}_{\mathbf{x}_{i,j} \sim \mathcal{X}} \left[\|D(E_A(\mathbf{x}_{i,j}), E_S(\mathbf{x}_{i,j})) - \mathbf{x}_{i,j}\|^2 \right] \quad (2)$$

where the full loss for any given batch is calculated using the mean squared error (MSE).

E. Cross-reconstruction loss

To really strengthen the disentanglement between the different features we also use a cross-reconstruction loss during training. Where given two data samples from the data set \mathcal{X} , we transfer the disentangled affect between them to produce two new data samples, computing the reconstruction loss for each using the following:

$$\begin{aligned} \mathcal{L}_{\text{cross}} = & \mathbb{E}_{\mathbf{x}_{i,j}, \mathbf{x}_{k,l} \sim \mathcal{X} \times \mathcal{X}} \left[\|D(E_A(\mathbf{x}_{i,j}), E_S(\mathbf{x}_{k,l})) - \mathbf{x}_{i,l}\|^2 \right] \\ & + \mathbb{E}_{\mathbf{x}_{i,j}, \mathbf{x}_{k,l} \sim \mathcal{X} \times \mathcal{X}} \left[\|D(E_A(\mathbf{x}_{k,l}), E_S(\mathbf{x}_{i,j})) - \mathbf{x}_{k,j}\|^2 \right]. \end{aligned} \quad (3)$$

This is possible because the dataset \mathcal{X} contains the ground truth for each affect/subject label combination. Similarly to the reconstruction loss, the loss for a full batch is calculated using MSE.

F. Triplet loss

Finally, we use a triplet loss function to encourage the autoencoder to cluster the low dimensional encodings tightly. While the reconstruction loss disentangles the subject and affects biometrics, there are no explicit requirements for separating different features in each subject/affect encodings. A triplet loss enforces the separation between the different feature classes in each latent space while also encouraging the same class's features to cluster closer together.

The triplet loss function is defined with regards to three different data points, an anchor, a positive and a negative. The anchor is the current data point that we are looking at, the positive is another data point of the same feature to whom we want to maximise the similarity, and the negative is a sample of a different feature, for which we want to maximise the difference. Resulting in the following two loss functions for subject and affect:

$$\begin{aligned} \mathcal{L}_{\text{trip}}^S = & \mathbb{E}_{\mathbf{x}_{i,l}, \mathbf{x}_{i,j}, \mathbf{x}_{k,l} \sim \mathcal{X}} \left[\|E_S(\mathbf{x}_{i,l}) - E_S(\mathbf{x}_{i,j})\| - \right. \\ & \left. \|E_S(\mathbf{x}_{i,l}) - E_S(\mathbf{x}_{k,l})\| + \alpha \right] \\ \mathcal{L}_{\text{trip}}^A = & \mathbb{E}_{\mathbf{x}_{i,l}, \mathbf{x}_{j,l}, \mathbf{x}_{i,k} \sim \mathcal{X}} \left[\|E_A(\mathbf{x}_{i,l}) - E_A(\mathbf{x}_{j,l})\| - \right. \\ & \left. \|E_A(\mathbf{x}_{i,l}) - E_A(\mathbf{x}_{i,k})\| + \alpha \right]. \end{aligned} \quad (4)$$

Where $\mathbf{x}_{i,l}$ is the anchor, $\mathbf{x}_{i,j}$ and $\mathbf{x}_{j,l}$ are the subject and affect positives, and $\mathbf{x}_{k,l}$ and $\mathbf{x}_{i,k}$ are the subject and affect negatives. α is called a margin that controls the distance between clusters. In other words, it makes the positive data points cluster tighter and separate themselves from the negative data points. Finally, we achieve a significantly faster convergence rate by combining the batch normalisation technique with

dropout into a single Independent-Component layer as Chen et al. has suggested to whiten the inputs of a neural network [37]. The integration of which can be seen in Figure/Table 2(b) which details the architecture of our model.

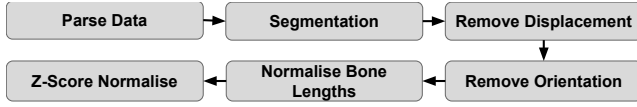


Fig. 3. Overview of the preprocessing pipeline.

G. Guided Grad-Cam for Model Explainability

The ability to explain the model decisions and link them to the gait cycle is a key factor in enhancing confidence in the model and adopting it in healthcare applications. Here we propose to exploit a local interpretability method to examine global properties of the model. To evaluate how our model attributes each of the joint signals to determine the final classification we make use of Guided Grad-Cam. Grad-CAM is a convolutional network visualisation method that exploits the fact that convolutional layers retain spatial information, and layers deeper within a network aim to capture high-frequency details. As such, the final convolutional layer of a network is expected to be a good compromise between high-level semantic and detailed spatial information. Grad-CAM computes the gradients for the output of a class label with respect to the feature map of a convolutional layer. These gradients are then global average pooled and passed through a ReLU function to reduce the impact of unwanted activations from other classes. Resulting in a low dimension heatmap of the general areas of activation of the same shape as the final convolutional layer. However, our main architecture is an autoencoder and thus the shape of the last convolutional layer is much smaller than the original input.

On the other hand, guided backpropagation visualises individual pixels detected by neurons in the network by backpropagating through ReLU layers and allows for very fine-grain high-resolution visualisation of the activation levels at each timeframe. However, these visualisations are not class-dependent, therefore to combine the best aspects of both approaches, Guided Grad-CAM performs an element-wise multiplication of the high-resolution non-class-dependent visualisation with the low-resolution class-dependent activation heatmap to generate a high-resolution class-dependent activation map that has the same shape as the original input. Retaining the original input shape is crucial for us because, in our case, we cannot feasibly interpolate from a very low-resolution activation map to 45 independent highly variable time series while retaining enough detail to detect changes in activations at specific gait cycle events. With Guided Grad-CAM, on the other hand, we can see changes in neuron activations at a frame-level accuracy allowing us to visualise what happens at different stages of a gait cycle.

H. From Local to Global Explanations

We are interested in evaluating interpretability across class samples in order to understand global properties of the pre-

diction model. We first split the test data by class so that each sample can be attributed towards its respective class, which in turn gives us an attribution map of how strongly each joint coordinate contributed towards the classification throughout the gait cycle. To determine a ranking of joint contributions we min-max normalise the absolute attribution values per data sample such that we can determine relative joint importance between all the test gait cycles. Since each joint is made up of three temporal signals that correspond to x, y, and z coordinates of each of the joint. We use the mean of all the signals to represent the overall activation for each of the 15 joints, giving us a single value for each joint in each gait cycle corresponding to the mean activations of that joint throughout. Then by summing the joint activations between gait cycles and dividing by the total number of activations of all joints we can also measure a percentage of how much each joint contributes towards the total activation.

IV. EVALUATION METHODOLOGY

We evaluated our proposed model, referred to as **AE-xyz** against a CNN model, referred to as **CNN-xyz**, on the raw joint time-series data, and three shallow baseline models that have been used in affect recognition based on human motion data [19]. These are a K Nearest Neighbours (KNN) classifier and a Support Vector Machine classifier trained on the manually engineered angular and velocity features, referred to as **KNN-man** and **SVM-man**, as well as an **SVM-xyz** model that was trained on the raw data. As a baseline the set of raw features extracted were the mean, std, min and max of the key angles such as knees or elbows, as well as velocities of key joints in the arms and legs, all of which have been shown to be useful features for affect recognition [19].

We performed an intra-subject stratified K-fold cross validation, where the entire dataset was split evenly into 5 different folds, ensuring an even split of subject/affect pairs between each of the folds. This kind of evaluation gives us an insight into how capable the model is at generalising to unseen gait cycles generated by known subjects.

A. Datasets

To test our classifier models' capabilities, we are making use of a publicly available motion capture library [38]. The dataset contains a set of motion captured movements of 30 non-professional subjects (15 male, 15 female, mean age 22, ranging from 17 to 29 years) performing various motions such as walking or throwing, labelled based on their identity, gender and emotion. The subjects acted out scripts specially prepared to stimulate target emotional responses. Our work's focal area will be the dataset's subject and emotion labels, specifically relating to gait motions, with a mean and standard deviation of 107 ± 15.5 gait cycle samples, respectively, per subject. The subjects acted out scripts specially prepared to stimulate target emotional responses. The motions were captured using retroreflective markers and a state-of-the-art motion capture system. The raw motion capture data was post-processed through a 3D animation software where the key joint positions

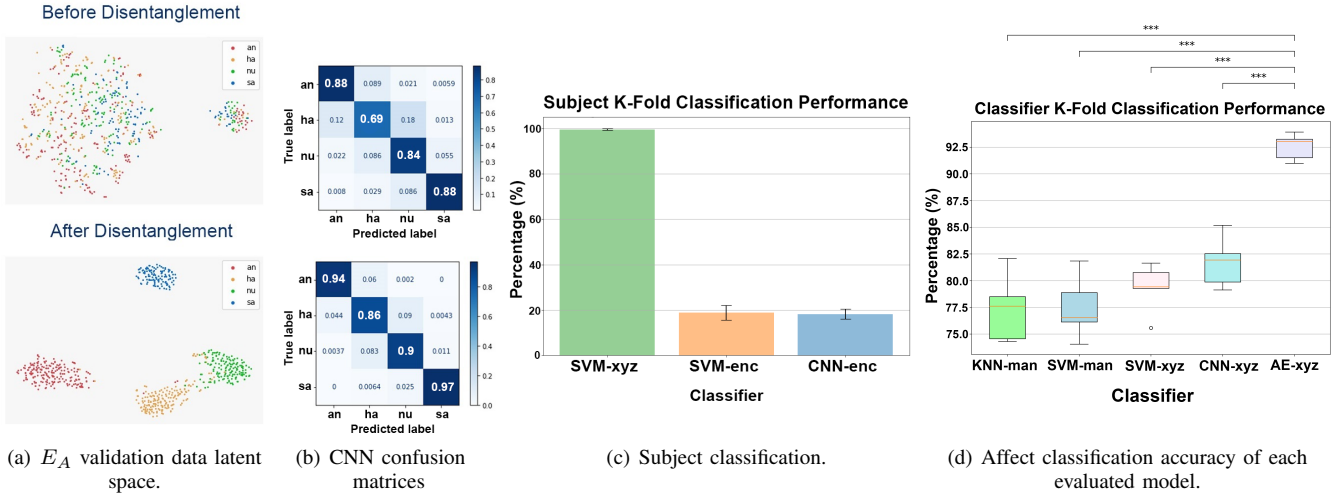


Fig. 4. (a) A T-SNE visualisation of the E_A latent space for the validation data, before training on the top and after on the bottom, (b) Confusion matrices of the CNN-xyz classifier, with labels an, ha, nu, sa, corresponding to angry, happy, neutral, sad, respectively. The top model trained and validated on the raw xyz data, and the bottom the AE-xyz classifier where the raw xyz validation data is first disentangled by the E_A encoder, and then classified using a CNN classifier trained on the E_A encodings of the train data. (c) Classification accuracy of subjects using a baseline SVM-xyz that is able to identify subjects with almost 100% accuracy, compared to models SVM-enc and CNN-enc, trained on the affect latent space rather than raw data. (d) Classification accuracy of each model represented as boxplots, plotted from the classification accuracy of each fold in the Stratified K-Fold cross-validation experiment, with a pairwise statistical comparison of each model with the AE-xyz model indicated above the figure.

from the motion data were projected to outline the human skeleton.

B. Preprocessing

The preprocessing pipeline is summarised in Figure 3. Firstly, it requires the segmentation of the data into gait cycles, defined as two consecutive heel strikes of the same foot. We extract the gait cycles by detecting points of minima in the acceleration of the right ankle. We use Single Spectrum Analysis (SSA) [39], [40], to smooth the acceleration signal for consistent and clean peak detection. To normalise the gait cycles four more steps are required: a) The removal of the global displacement and rotation, to ensure consistency in the coordinate space between the data samples. b) The bone lengths between the subjects are normalised such that the model is encouraged to learn the difference in subject gait styles rather than differences in height, and c) a z-score normalisation is performed such that all of the gait cycles have 0 mean and 1 standard deviation.

C. Displacement Removal

To remove the global displacement of the subject during the gait cycle we take the position of a root joint, in our case the center hip joint, and subtract its position from all the other joints for each frame of the gait cycle. Let $x \in \mathbb{R}^{T \times 3J}$ be a motion, where x is a gait cycle, T is the temporal duration of the motion, and J is the number of joints each specified by $(x\hat{i}, y\hat{j}, z\hat{k})$ coordinates. Then x^{jt} , is a joint ($j \in J$) at time frame ($t \in T$), of the gait cycle. Given a root joint ($r \in J$), then the new position \hat{x}^{jt} of the joint without global displacement is given by:

$$\hat{x}^{jt} = x^{jt} - x^{rt} \quad (5)$$

To remove rotation effects from motion \hat{x} , we need to project the coordinates of each joint to a new set of axes, with the facing direction of the skeleton as one of the axes. To extract the forward facing vector we exploit the fact that in general the body's facing is perpendicular to the hip, and as such we define a limb vector $l_{hp} \in L$, as a valid connection between two joints, representing the skeleton's bone, where L is the set of all valid joint connections, and $\{h, p\} \in J$ representing the left and right hip joints respectively. We can derive l_{hp} from motion \hat{x} at any given point in time t using:

$$l_{hp} = \frac{x^{ht} - x^{pt}}{\|x^{ht} - x^{pt}\|} \quad (6)$$

From which we can find the direction vector d by finding the vector perpendicular to l_{hp} and \hat{j} , a unit vector representing the up direction.

$$d = l_{hp} \times \hat{j} \quad (7)$$

Using all this we can define a change of basis matrix C , which will transform the joint positions to a new set of axes, with the skeleton facing the y-axis, effectively removing the effects of the z-axis rotation of the motion. We define C as:

$$C = \begin{bmatrix} l_{hp} \\ d \\ \hat{j} \end{bmatrix}, \quad (8)$$

and then using it we can define a new position for a joint, $\hat{\hat{x}}^{jt}$, where the effects of global displacement and rotation are removed, by solving the following linear equation:

$$C^{-1}\hat{\hat{x}}^{jt} = \hat{x}^{jt} \quad (9)$$

D. Bone Normalisation

To normalize the skeleton bone lengths between the subjects, we define each bone of the skeleton as limb vectors $\mathbf{l}_{pc} = \ddot{x}^{ct} - \ddot{x}^{pt}$, where $p \in J$ is a parent joint and $c \in J$ the child joint. For the motions to remain unchanged it is crucial that the joints are rescaled in a hierarchical order starting from the base joint, i.e. to rescale the arm, start with the shoulder, then upper arm and finally the forearm. Each limb vector is rescaled to the mean length of the limb within the dataset $\bar{\mathbf{l}}_{pc}$, which can then be used to calculate a scaling ratio α for each limb:

$$\alpha = \frac{\bar{\mathbf{l}}_{pc}}{\|\mathbf{l}_{pc}\|} \quad (10)$$

Using this scaling ratio a joint can then be rescaled to the new position \hat{x}^{jt} using:

$$\hat{x}^{jt} = \alpha \cdot \mathbf{l}_{pc} + \ddot{x}^{pt}, \quad (11)$$

where it is now scale, rotation and displacement independent.

V. RESULTS

This section aims to evaluate the ability of the proposed method to disentangle biometrics and affects as well as to quantify the ability of the algorithm to preserve the privacy of the subjects and explain the network decisions with relation to the gait cycle and the joints contributions.

TABLE I
ACCURACY - QUANTITATIVE RESULTS (%) FOR STRATIFIED K-FOLD CROSS VALIDATION

Model	Angry		Happy		Neutral		Sad	
	Acc	Std	Acc	Std	Acc	Std	Acc	Std
KNN-man	79.71	3.58	68.42	5.52	75.41	4.63	88.17	2.35
SVM-man	80.66	3.61	68.89	5.59	74.96	6.38	87.1	1.83
SVM-xyz	87.03	3.48	69.62	3.36	76.2	3.36	85.19	4.16
CNN-xyz	88.4	6.11	68.77	9.24	83.62	3.97	87.74	3.55
AE-xyz	93.82	2.34	86.17	3.99	90.25	2.89	96.82	2.06

TABLE II
F1 SCORES - QUANTITATIVE RESULTS (%) FOR STRATIFIED K-FOLD CROSS VALIDATION

Model	Angry			Happy			Neutral			Sad		
	Pre	Rec	F1	Pre	Re	F1	Pre	Rec	F1	Pre	Rec	F1
KNN-man	82	80	81	71	68	70	72	75	74	86	88	87
SVM-man	80	81	80	71	69	70	72	75	74	88	87	88
SVM-xyz	83	87	85	72	70	71	76	76	76	88	85	87
CNN-xyz	87	88	88	78	69	73	74	84	78	90	88	89
AE-xyz	95	93	94	88	90	89	92	91	91	96	97	96

A. Affect Classification

The proposed methodology improved the mean classification accuracy more than 7% and the difference in performance between the **CNN-xyz** and **AE-xyz** was statistically significant with ($P \leq 0.005$), achieving a 92.6% mean accuracy and 93 mean F1 score over the CNN's 84.6% mean accuracy and 82

mean F1 score. From the confusion matrix in Figure 4(b) we can see that the models' weakest affect was happy in both cases, however **AE-xyz** misclassified it a lot less than **CNN-xyz**, and it was overall better across the board. This shows that the Autoencoder is an efficient way of extracting low-dimensional representations of the raw human motion data to reduce complexity and improve classification. Table I provides an overview of the accuracy scores for each of the models, and Table II provides an overview of the F1 scores.

in Figure 4(a), a closer inspection of the low dimensional representations formed by the cross-subject Autoencoder reveals that both the Subject-Encoder and Affect-Encoder are forming clearly separated clusters. In fact, before training the low dimensional encodings are randomly intertwined, whereas after training the affects become separated out into their individual clusters with a small amount of overlap between the happy and neutral affects, where the model's predictive ability is slightly lower.

B. Enhancing Privacy Preservation

Finally, to evaluate the effectiveness of our subject identification, we devised an experiment where we first test the subject identification accuracy using the raw joint position data and compare it to the identification accuracy using the disentangled representation. Figure 4(c) demonstrates that subjects classification based on the encoded data is low, which confirms that the proposed method is effective.

C. Joints' Attribution across and within gait cycles

Figure 5(a) shows the mean total attributions of each joint between the test samples for each emotion and a box plot overlaid to show the distribution of joint attributions between subjects. We find that the model generally balances its attention between the upper and lower body, with the left and right sides of the hips standing out the most from the others in terms of importance. The center hip joint has an attribution of 0, as it is used to normalise the displacement of the skeleton and thus it does not move from the starting position. There is a general symmetry between the left- and right-hand sides of the body, which is expected since human posture normally remains symmetric under different affects. We also notice that for the positively classified samples depicted with the green bars, the model demonstrates larger activations, which reflect certainty in the decision compared to samples classified incorrectly, which result in smaller mean attributions (red bars). Overall higher activations are associated with the hips, torso, elbows and head, whereas end effector joints, such as wrists and ankles, attract less attention.

Different phases of the gait cycle are related to the synergistic function of specific groups of muscles and tendons and thus investigating activations across the gait cycle is clinically relevant. Towards this aim we firstly segment the body into three parts, the upper body, mid-body, and lower body, as we found those areas to have the most overlap between joints' activation. The mid-body was specifically separated from the lower body as the hips have shown themselves to be a strong

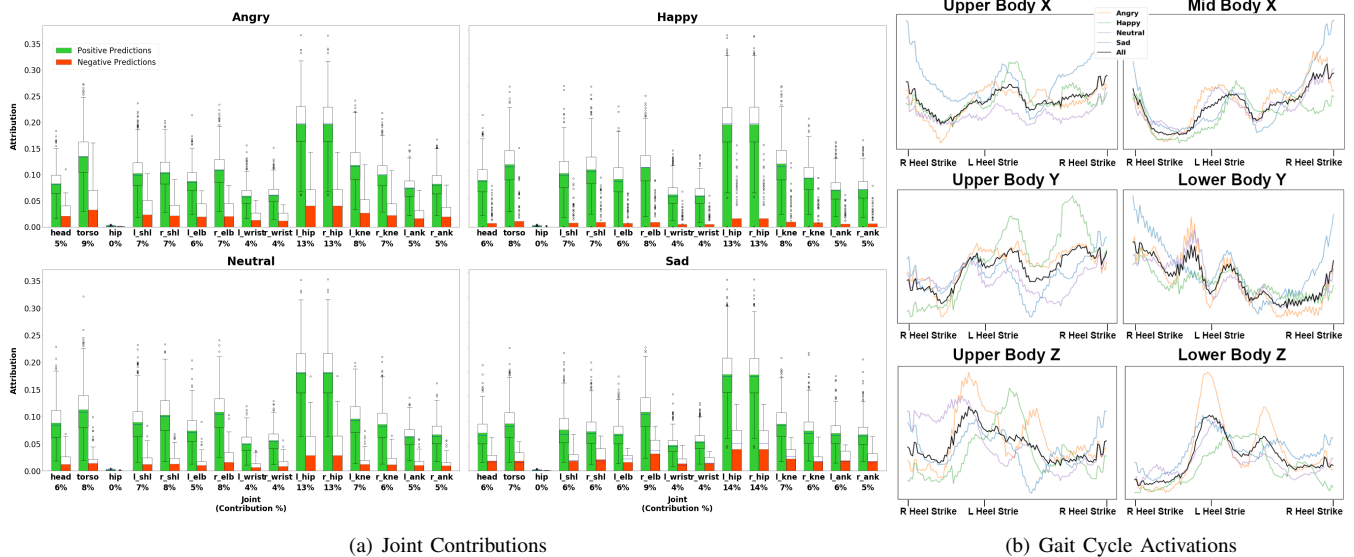


Fig. 5. (a) A bar chart with the mean total attributions of each joint between the test samples for each emotion and a box plot overlaid to show the distributions of joint attributions between subjects. Prefixes $r_$ and $l_$ represent the right and left side of the subject’s body and joint abbreviations shl, elb, wrist, hip, kne, ank correspond to shoulders, elbow, wrist, hip, knee and ankle respectively. (b) A visualisation of the attributions of joint groups, for each emotion individually as well as the mean, over a gait cycle consisting of 128 frames. The upper body consists of the head, neck, shoulders, elbows and wrists. The mid body comprises the hips and the lower body of the knees and ankles. Only the axes with the majority of joint movements, and as a result, the majority of attributions, are presented.

indicator of emotion and their activation pattern is slightly different from the upper and lower body. The head, torso, shoulder, elbow and wrist joints have been combined for the upper body as their activation patterns closely overlap. Then for the lower body, the knees and ankles were chosen as not only do they tend to follow a slightly different pattern from the upper body, but they also line up more with gait cycle events such as heel strikes which is what we would expect as those events’ definitions are based on the lower body. Since the motions have also been normalised to be rotation-independent, any given joint will generally have the majority of its motions mapped on one or two axes, whereas the projection to the third axes will be minimal. For example, the movements of the lower body are largely constrained to the X and Y-axis as the rotation within the Z-axis has been removed from the motions. This results in the variances within the Z-axis not attributing much towards the final classification of the gait cycle. Figure 5(b) shows the normalised attributions of the body parts (Upper, Mid, Lower) for each emotion and their mean across the gait cycle. For a total of 45 temporal signals of the joints, We present the normalised joint attributions considering axes with significant attributions.

We also notice a strong relation between gait cycle events such as heel strikes in the lower body and attributions in the X and Z axis for the lower body. For example in Figure 5(b) we see a strong dip in the attribution of the lower body’s Y-axis as the left heel strikes the ground. Furthermore, we note a peak in activations on both ends of the heel strikes with the model paying attention to the manner of how the subject places their foot and lifts the next one. We also see a strong peak in the lower body Z-axis as the left heel strikes the ground and the

right toe-off happens, with continuous activations during the swing phase of the right foot. With the upper body we observe peaks in activations around the left heel strike events with more consistent activations throughout, but this time rather than peaking at the point of the heel strike the model activates stronger earlier when one swing phase terminates and another starts. Presumably, this implies that the model focuses on areas of terminal arm swing. In more expressive emotions with more vivid upper body movement like happy and angry the model activations’ are stronger than for less expressive emotions such as neutral and sad, with sad having especially low attributions within the arm swing phase. With the mid-body, we notice more consistent activations between different emotions with strong activations during the cross-over phase where the gait cycle enters into double support and then terminates from it into a right leg swing, presumably the section of the gait cycle where the most expression of the hip tilt happens.

VI. CONCLUSIONS

To our knowledge, our work is the first to disentangle affects and biometrics successfully by exploiting a multi-encoder, decoder architecture to map human motion data onto a low dimensional space. We show that this method not only enhances affect recognition performance but also preserves subjects privacy. Furthermore, based on a gradient-based explainability method, we examine global properties of the model that allow us to understand how different joints contribute to the model decision across gait cycles. Future work will focus on evaluating the method in larger datasets and developing an end-to-end training strategy that does not require explicit gait segmentation.

REFERENCES

- [1] P. Winkielman, P. Niedenthal, J. Wielgosz, J. Eelen, and L. C. Kavanagh, "Embodiment of cognition and emotion." in *APA handbook of personality and social psychology, Volume 1: Attitudes and social cognition.*, M. Mikulincer, P. R. Shaver, E. Borgida, and J. A. Bargh, Eds. American Psychological Association, 2015, pp. 151–175.
- [2] M. A. Hashmi, Q. Riaz, M. Zeeshan, M. Shahzad, and M. M. Fraz, "Motion Reveal Emotions: Identifying Emotions From Human Walk Using Chest Mounted Smartphone," *IEEE Sensors Journal*, vol. 20, no. 22, pp. 13 511–13 522, Nov. 2020.
- [3] F. Deligianni, Y. Guo, and G.-Z. Yang, "From Emotions to Mood Disorders: A Survey on Gait Analysis Methodology," *IEEE Journal of Biomedical and Health Informatics*, vol. 23, no. 6, pp. 2302–2316, Nov. 2019.
- [4] Z. Chen, T. Zhu, P. Xiong, C. Wang, and W. Ren, "Privacy preservation for image data: A GAN-based method," *International Journal of Intelligent Systems*, vol. 36, no. 4, pp. 1668–1685, Apr. 2021.
- [5] M. Gong, J. Liu, H. Li, Y. Xie, and Z. Tang, "Disentangled Representation Learning for Multiple Attributes Preserving Face Deidentification," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–13, 2020.
- [6] X. Gu, Y. Guo, F. Deligianni, B. Lo, and G.-Z. Yang, "Cross-Subject and Cross-Modal Transfer for Generalized Abnormal Gait Pattern Recognition," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–15, 2020.
- [7] K. Aberman, R. Wu, D. Lischinski, B. Chen, and D. Cohen-Or, "Learning Character-Agnostic Motion for Motion Retargeting in 2D," *ACM Transactions on Graphics*, vol. 38, no. 4, pp. 1–14, Jul. 2019.
- [8] A. B. Arrieta, N. Díaz-Rodríguez, J. Del Ser, A. Bennetot, S. Tabik, A. Barbado, S. García, S. Gil-López, D. Molina, R. Benjamins, R. Chatila, and F. Herrera, "Explainable Artificial Intelligence (XAI): Concepts, Taxonomies, Opportunities and Challenges toward Responsible AI," *arXiv:1910.10045 [cs]*, Dec. 2019.
- [9] C. L. Vaughan, C. L. Vaughan, C. L. Vaughan, and C. L. Vaughan, *GaitCD*. Kiboho Publishers, 1999.
- [10] S. Xu, J. Fang, X. Hu, E. Ngai, Y. Guo, V. C. M. Leung, J. Cheng, and B. Hu, "Emotion Recognition From Gait Analyses: Current Research and Future Directions," *arXiv:2003.11461 [cs, stat]*, Aug. 2020.
- [11] M. M. Gross, E. A. Crane, and B. L. Fredrickson, "Effort-Shape and kinematic assessment of bodily expression of emotion during gait," *Human Movement Science*, vol. 31, no. 1, pp. 202–221, Feb. 2012.
- [12] C. L. Roether, L. Omlor, A. Christensen, and M. A. Giese, "Critical features for the perception of emotion from gait," *Journal of Vision*, vol. 9, no. 6, pp. 15–15, Jun. 2009.
- [13] J. M. Montepare, S. B. Goldstein, and A. Clausen, "The identification of emotions from gait information," *Journal of Nonverbal Behavior*, vol. 11, no. 1, pp. 33–42, 1987.
- [14] A. Barliya, L. Omlor, M. A. Giese, A. Berthoz, and T. Flash, "Expression of emotion in the kinematics of locomotion," *Experimental Brain Research*, vol. 225, no. 2, pp. 159–176, Mar. 2013.
- [15] F. Loula, S. Prasad, K. Harber, and M. Shiffrar, "Recognizing people from their movement." *Journal of Experimental Psychology: Human Perception and Performance*, vol. 31, no. 1, pp. 210–220, Feb. 2005.
- [16] D. Das, "An Approach to Emotion Identification Using Human Gait," in *Proceedings of Fourth International Conference on Soft Computing for Problem Solving*, K. N. Das, K. Deep, M. Pant, J. C. Bansal, and A. Nagar, Eds. Springer India, 2015, vol. 336, pp. 165–175.
- [17] Z. Zhang, Y. Song, L. Cui, X. Liu, and T. Zhu, "Emotion recognition based on customized smart bracelet with built-in accelerometer," *PeerJ*, vol. 4, p. e2258, Jul. 2016.
- [18] G. Venture, H. Kadone, T. Zhang, J. Grèzes, A. Berthoz, and H. Hicheur, "Recognizing Emotions Conveyed by Human Gait," *International Journal of Social Robotics*, vol. 6, no. 4, pp. 621–632, Nov. 2014.
- [19] S. Srivastava, V. Rastogi, C. Prakash, and D. Sethi, "Robust approach for emotion classification using gait," in *International Conference on Innovative Computing and Communications*, D. Gupta, A. Khanna, S. Bhattacharyya, A. E. Hassanien, S. Anand, and A. Jaiswal, Eds. Singapore: Springer Singapore, 2021, pp. 885–894.
- [20] A. M. Saleh and T. Hamoud, "Analysis and best parameters selection for person recognition based on gait model using CNN algorithm and image augmentation," *Journal of Big Data*, vol. 8, no. 1, p. 1, Dec. 2021.
- [21] M. Shahrum Md Guntor, R. Sahak, A. Zabidi, N. Md Tahir, I. Mohd Yassin, Z. Ismael Rizman, R. Baharom, and N. Abdul Wahab, "Convolutional Neural Network (CNN) based Gait Recognition System using Microsoft Kinect Skeleton Features," *International Journal of Engineering & Technology*, vol. 7, no. 4.11, p. 202, Oct. 2018.
- [22] Q. Zou, Y. Wang, Q. Wang, Y. Zhao, and Q. Li, "Deep Learning-Based Gait Recognition Using Smartphones in the Wild," *arXiv:1811.00338 [cs, eess, stat]*, Apr. 2020.
- [23] Y. Guo, F. Deligianni, X. Gu, and G.-Z. Yang, "3-D Canonical Pose Estimation and Abnormal Gait Recognition With a Single RGB-D Camera," *IEEE Robotics and Automation Letters*, vol. 4, no. 4, pp. 3617–3624, Oct. 2019.
- [24] Y. Iwasawa, K. Nakayama, I. Yairi, and Y. Matsuo, "Privacy Issues Regarding the Application of DNNs to Activity-Recognition using Wearables and Its Countermeasures by Use of Adversarial Training," in *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence*. International Joint Conferences on Artificial Intelligence Organization, Aug. 2017, pp. 1930–1936.
- [25] I. Y. Jung, "A review of privacy-preserving human and human activity recognition," *International Journal on Smart Sensing and Intelligent Systems*, vol. 13, no. 1, pp. 1–13, 2020.
- [26] C. Song and V. Shmatikov, "Overlearning Reveals Sensitive Attributes," *arXiv:1905.11742 [cs, stat]*, Feb. 2020, arXiv: 1905.11742.
- [27] E. Luger and T. Rodden, *Ethics and consent in the (sociotechnical) wild, ser. Into the Wild: Beyond the Design Research Lab*. Springer, jul 2019, pp. 149–172.
- [28] Y. Kim, B. Yoo, Y. Kwak, C. Choi, and J. Kim, "Deep generative-contrastive networks for facial expression recognition," *arXiv:1703.07140 [cs]*, May 2019.
- [29] X. Wang, K. Wang, and S. Lian, "A Survey on Face Data Augmentation," *Neural Computing and Applications*, vol. 32, no. 19, Oct. 2020.
- [30] X. Liu, B. Vijaya Kumar, P. Jia, and J. You, "Hard negative generation for identity-disentangled facial expression recognition," *Pattern Recognition*, vol. 88, pp. 1–12, Apr. 2019.
- [31] Q. Zhu, L. Gao, H. Song, and Q. Mao, "Learning to disentangle emotion factors for facial expression recognition in the wild," *International Journal of Intelligent Systems*, vol. 36, no. 6, pp. 2511–2527, Jun. 2021.
- [32] H. Wu, J. Jia, L. Xie, G. Qi, Y. Shi, and Q. Tian, "Cross-VAE: Towards Disentangling Expression from Identity For Human Faces," in *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, May 2020, pp. 4087–4091.
- [33] Y. Deng, J. Yang, D. Chen, F. Wen, and X. Tong, "Disentangled and Controllable Face Image Generation via 3D Imitative-Contrastive Learning," *arXiv:2004.11660 [cs]*, Sep. 2020.
- [34] B. Chen, Y. Zhang, H. Tan, B. Yin, and X. Liu, "PMAN: Progressive Multi-Attention Network for Human Pose Transfer," *IEEE Transactions on Circuits and Systems for Video Technology*, pp. 1–1, 2021.
- [35] Y. Sun, Y. Ye, W. Liu, W. Gao, Y. Fu, and T. Mei, "Human Mesh Recovery From Monocular Images via a Skeleton-Disentangled Representation," in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE, Oct. 2019, pp. 5348–5357.
- [36] K. Mangalam, E. Adeli, K.-H. Lee, A. Gaidon, and J. C. Niebles, "Disentangling Human Dynamics for Pedestrian Locomotion Forecasting with Noisy Supervision," in *2020 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, Mar. 2020, pp. 2773–2782.
- [37] G. Chen, P. Chen, Y. Shi, C.-Y. Hsieh, B. Liao, and S. Zhang, "Rethinking the Usage of Batch Normalization and Dropout in the Training of Deep Neural Networks," *arXiv:1905.05928 [cs, stat]*, May 2019.
- [38] Y. Ma, H. M. Paterson, and F. E. Pollick, "A motion capture library for the study of identity, gender, and emotion perception from biological motion," *Behavior Research Methods*, vol. 38, no. 1, pp. 134–141, Feb. 2006.
- [39] N. Golyandina, A. Korobeynikov, and A. Zhigljavsky, *Singular Spectrum Analysis with R, ser. Use R!* Springer Berlin Heidelberg, 2018.
- [40] X. Gu, F. Deligianni, B. Lo, W. Chen, and G. Yang, "Markerless gait analysis based on a single RGB camera," in *2018 IEEE 15th International Conference on Wearable and Implantable Body Sensor Networks (BSN)*. IEEE, Mar. 2018, pp. 42–45.