# Adversarial Refinement Network for Human Motion Prediction

Xianjin CHAO[1][0000−0002−4020−8223], Yanrui Bin[2][0000−0003−2845−3928], Wenqing Chu[3], Xuan Cao[3], Yanhao Ge[3], Chengjie Wang[3], Jilin Li[3], Feiyue Huang[3], and Howard Leung[1][0000−0002−2633−2965]

[1] City University of Hong Kong, Hong Kong, China
`xjchao2-c@my.cityu.edu.hk`
`howard@cityu.edu.hk`
[2] Key Laboratory of Image Processing and Intelligent Control, School of Artificial Intelligence and Automation, Huazhong University of Science and Technology, Wuhan, China
`yrbin@hust.edu.cn`
[3] Tencent Youtu Lab, Shanghai, China
{`wenqingchu, marscao, halege, jasoncjwang, jerolinli, garyhuang`}`@tencent.com`

**Abstract.** Human motion prediction aims to predict future 3D skeletal sequences by giving a limited human motion as inputs. Two popular methods, recurrent neural networks and feed-forward deep networks, are able to predict rough motion trend, but motion details such as limb movement may be lost. To predict more accurate future human motion, we propose an Adversarial Refinement Network (ARNet) following a simple yet effective coarse-to-fine mechanism with novel adversarial error augmentation. Specifically, we take both the historical motion sequences and coarse prediction as input of our cascaded refinement network to predict refined human motion and strengthen the refinement network with adversarial error augmentation. During training, we deliberately introduce the error distribution by learning through the adversarial mechanism among different subjects. In testing, our cascaded refinement network alleviates the prediction error from the coarse predictor resulting in a finer prediction robustly. This adversarial error augmentation provides rich error cases as input to our refinement network, leading to better generalization performance on the testing dataset. We conduct extensive experiments on three standard benchmark datasets and show that our proposed ARNet outperforms other state-of-the-art methods, especially on challenging aperiodic actions in both short-term and long-term predictions.

## 1 Introduction

Given the observed human 3D skeletal sequences, the goal of human motion prediction is to predict plausible and consecutive future human motion which convey abundant clues about the person's intention, emotion and identity.

Effectively predicting the human motion plays an important role in wide visual computing applications such as human-machine interfaces [1], smart surveillance [2], virtual reality [3], healthcare applications [4], autonomous driving [5] and visual human-object tracking [6]. However, predicting plausible future human motion is a very challenging task due to the non-linear and highly spatial-temporal dependencies of human body parts during movements [7].Considering the time-series property of human motion sequence, recent deep learning based methods formulated the human motion prediction task as a sequence-to-sequence problem and achieved remarkable progresses by using chain-structured Recurrent Neural Networks (RNNs) to capture the temporal dependencies frame-by-frame among motion sequence. However, recent literature [8] indicated that the chain-structured RNNs suffer from error accumulation in temporal modeling and deficiency in spatial dynamic description, leading to problems such as imprecise pose and mean pose in motion prediction.

Feed-forward deep networks  [9] are regarded as alternative solutions for human motion prediction task by learning rich representation from all input motion sequences at once. The holistic reasoning of the human motion sequences leads to more consecutive and plausible predictions than chain-structured RNNs.

Unfortunately, current feed-forward deep networks adopt singe-stage architecture and tend to generate the predicted motion coarsely thus yielding unsatisfactory performance, especially for complex aperiodic actions (e.g., Direction or Greeting in H3.6m dataset). The reason is that it is difficult to guide the network to focus more on detailed information when directly predicting the future human motion from limited input information.

To address the above issues, we propose a novel Adversarial Refinement Network (ARNet) which resorts to a coarse-to-fine framework. We decompose the human motion prediction problem into two stages: coarse motion prediction and finer motion refinement. By joint reasoning of the input-output space of the coarse predictor, we achieve to take both the historical motion sequences and coarse future prediction as input not just one-sided information to polish the challenging human motion prediction task. The coarse-to-fine design allows the refinement module to concentrate on the complete motion trend brought by the historical input and coarse prediction, which are ignored in previous feed-forward deep networks used for human motion prediction.

Given different actions performed by diverse persons fed to the refinement network in training and testing, the coarse prediction results tend to be influenced by generalization error, which makes it difficult for the refinement network to obtain the fine prediction robustly. We therefore enhance the refinement network with adversarial error distribution augmentation. During training, we deliberately introduce the error distribution by learning through the adversarial mechanism among different subjects based on the coarse prediction. In testing, our cascaded refinement network alleviates the prediction error from the coarse predictor resulting in a finer prediction. Our adversarial component acts as regularization to let our network refine the coarse prediction well. Different from the previous work [10] which casts the predictor as a generator and introduces

discriminator to validate the prediction results, our adversarial training strategy aims to generate error distribution which acts as implicit regularization for better refinement instead of directly generating the skeleton data as prediction. The error augmentation is achieved by a pair of adversarial learning based generator and discriminator.

Consequently, the proposed ARNet achieves state-of-the-art results on several standard human motion prediction benchmarks over diverse actions categories, especially over the complicated aperiodic actions as shown in Figure 2.

Our contributions are summarized as follows:

- We propose a coarse-to-fine framework to decompose the difficult prediction problem into coarse prediction task and refinement task for more accurate human motion prediction.
- We design an adversarial learning strategy to produce reasonable error distribution rather than random noise to optimize the refinement network.
- The proposed method is comprehensively evaluated on multiple challenging benchmark datasets and outperforms state-of-the-art methods especially on complicated aperiodic actions.

## 2   Related Work

### 2.1   Human Motion Prediction

With the emergence of large scale open human motion capture (mocap) datasets, exploring different deep learning architectures to improve human motion prediction performance on diverse actions has become a new trend. Due to the inherent temporal-series nature of motion sequence, the chain-structured Recurrent Neural Networks (RNNs) are natively suitable to process motion sequences. The Encoder-Recurrent-Decoder (ERD) model [11] simultaneously learned the representations and dynamics of human motion. The spatial-temporal graph is later employed in [12] to construct the Structural-RNNs (SRNN) model for human motion prediction. The residual connections in RNN model (RRNN) [13] helped the decoder model prior knowledge of input human motion. Tang et al. [8] adopted the global attention and Modified Highway Unit (MHU) to explore motion contexts for long-term dependencies modeling. However, these chain-structured RNNs suffer from either frozen mean pose problems or unnatural motion in predicted sequences because of the weakness of RNNs in both long-term temporal memory and spatial structure description. Feed-forward deep network as an emerging framework has shown the superiority over chain-structured RNNs. Instead of processing input frame by frame like chain-structured RNNs, feed-forward deep networks feed all the frames at once, which is a promising alternative for feature extraction to guarantee the integrity and smoothness of long-term temporal information in human motion prediction [8]. In this paper, our ARNet is on the basis of feed-forward deep network.

## 2.2    Prediction Refinement

Refinement approaches learn good feature representation from the coarse results in output space and infer the precise location of joints in a further step by recovering from the previous error, which have achieved promisingly improvement in human pose related work. Multi-stage refinement network [14] associated the coarse pose estimation and refinement in one go to improve the accuracy of 3D human pose estimation by jointly processing the belief maps of 2D joints and projected 3D joints as the inputs to the next stage. Cascaded Pyramid Network (CPN) [15] introduced refinement after the pyramid feature network for sufficient context information mining to handle the occluded and invisible joints estimation problems. Another trend of refinement mechanism performed coarse pose estimation and refinement separately. PoseRefiner[16] refined the given pose estimation by modelling hard pose cases. Posefix [17] proposed an independent pose refinement network for arbitrary human pose estimator and refined the predicted keypoints based on error statistics prior. Patch-based refinement [18] utilised the retain fine details from body part patches to improve the accuracy of 3D pose estimation. In contrast to the previous work, we further adopt the benefits of refinement network to deal with the problems in 3D human motion prediction via a creative coarse-to-fine manner.

## 2.3    Adversarial Learning

Inspired by the minimax mechanism of Generative Adversarial Networks (GANs) [19], adversarial learning has been widely adopted to train neural networks [20,21,22]. Several attempts have been proposed to perform data augmentation in the way of adversarial learning, which mainly rely on the pixel manipulation through image synthesis [23] or a serious of specific image operations [24]. The adversarial learning based data augmentation shows powerful potential for model performance improvement. In [25], the results of image recognition achieved promising improvement due to the image synthesis data augmentation. In human motion prediction, [10] adopted a predictor with two discriminators to keep the fidelity and continuity of human motion predicted sequences by adversarial training. In this work, we introduce an online data augmentation scheme in the motion space to improve generalization and optimize the refinement network.

# 3    Methodology

## 3.1    Overall Framework

The overall framework of our ARNet is shown in Figure 1. The coarse-to-fine module consists of a coarse predictor $\mathcal{P}$ and a refinement network $\mathcal{R}$. In the context of human motion prediction, given N frames of observed human motion at once, the coarse predictor $\mathcal{P}$ aims to forecast the following T frames of human motion. The input human motion sequences $X = \{x_1, x_2, ..., x_n\}$ are first fed into the predictor $\mathcal{P}$ to obtain coarse future human motion $Y = \{y_1, y_2, ..., y_n\}$,
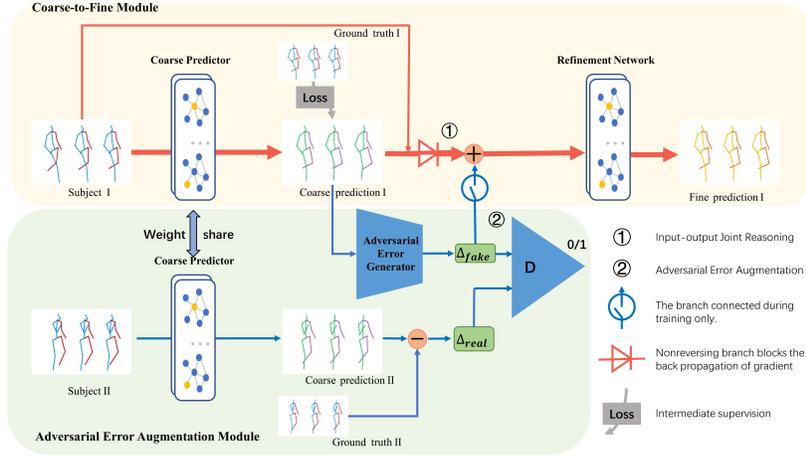
**Fig. 1. The overall framework of our ARNet.** The proposed coarse-to-fine module consists of coarse predictor and refinement network as shown in the top part. The bottom part illustrates the dedicated adversarial error augmentation module which consists of coarse predictor with a pair of error generator and discriminator. The observed human motion sequence of Subject I and Subject II are separately fed to the weight-shared coarse predictors to obtain corresponding coarse human motion prediction. Then the generator in the adversarial error augmentation module adopts the coarse prediction of Subject I as the conditional information to generate fake motion error of Subject II in an adversarial manner. After that, the augmented error distribution and the real coarse prediction are both utilised to optimize the refinement network for fine human motion prediction

where $x_i, y_i \in \mathbb{R}^K$ are $K$ dimensional joint features represented as exponential map of joint angle in each frame. Then in the adversarial error augmentation module, we adopt a pair of generator and discriminator to produce fake motion error calculated from the coarse prediction from a person (subject I) and the real motion error from another person (subject II) as the conditional information for the next stage fine prediction.

During training, we deliberately introduce the error distribution by learning through the adversarial mechanism among different subjects based on the coarse prediction. In testing, our cascaded refinement network alleviates the prediction error from the coarse predictor resulting in a finer prediction.

### 3.2 Refinement Network

Given the input motion sequence, we adopt a Graph Convolutional Network (GCN) [26,27], a popular feed-forward deep network which is specialized in dealing with the graph structured data, to initially model the spatial-temporal dependencies among the human poses and obtain the coarse human motion prediction. We construct a $K$ nodes graph $G = (V, E)$, where $V = \{v_i | i = 1, ..., K\}$

denotes the node set and $E = \{e_{i,j}|i,j = 1,...,K\}$ denotes the edge set. The main idea of Graph Convolutional Network is that, each $d$ dimensional node representations $H_v^l \in \mathbb{R}^d$ is updated by feature aggregation of all its neighbors defined by the weighted adjacency matrix $A^l \in \mathbb{R}^{K \times K}$ on the $l$-th Graph Convolutional layer. Therefore, the spatial structure relationships between the nodes could be fully encoded and the $l$-th Graph Convolution layer outputs a $K \times d$ matrix $H^{l+1} \in \mathbb{R}^d$:

$$H^{l+1} = \sigma(A^l H^l W^l) \tag{1}$$

where $\sigma(\cdot)$ denotes an activation function and $W^l \in \mathbb{R}^{d \times \hat{d}}$ denotes the trainable weight matrix. The network architecture of our predictor is similar to [9], which is the state-of-the-art feed-forward baseline on human motion prediction.

In order to improve the human motion prediction performance in a further step, we construct a coarse-to-fine framework, which cascades N-stage refinement network on top of the preliminary predictor, to process the complete future information of the output from human motion predictor iteratively. Given the input human motion sequences $H_I$, we initially obtain the coarse human motion prediction sequences $H_\mathcal{P} = f_p(H_I)$ from the preliminary predictor and forward the fusion of historical and future sequences as the inputs to the refinement network. As a result, we output the final refined human motion prediction sequences by error correction of initially coarse prediction $H_\mathcal{R} = f_r(H_\mathcal{P} + H_I)$.

### 3.3   Adversarial Learning Enhanced Refinement Network

Considering that the human motion sequences collected by different actors in datasets contain variations, especially for complicated aperiodic actions, various error scenarios will occur. To improve the error-correction ability and robustness of our refinement network, we additionally introduce an adversarial learning mechanism to generate challenging error cases which are fed to the refinement network together with the coarse prediction. We randomly choose 1 person's actions sequences (Subject II) from the 6 subjects' actions sequences in the training dataset and feed it to the predictor in another branch to get the independent coarse prediction sequences for every epoch as shown in Figure 1. Then the real error is able to be computed from this person's coarse prediction sequences and the corresponding ground-truth. To augment this person's error cases to the other 5 people, we utilise a generator that produces fake human motion error to fool the discriminator. The discriminator constantly tries to distinguish between real error cases and fake error cases so as to transfer different persons' error to other subjects in the mocap dataset. This augmentation provides rich error cases as input to our refinement network, leading to better generalization performance on the testing dataset.

We train the networks following the standard GAN pipeline. During training, the adversarial error generator generates error bias which will be added on the coarse prediction and then fed to refinement network. The adversarial refinement network effectively learns from the coarse prediction with adversarial error augmentation. During testing, the coarse prediction without added error is fed

directly to the adversarial refinement network and get finer prediction as final results.

### 3.4 Training Loss

In this section, we describe the training loss functions for different modules. Notably, in order to achieve joint reasoning of the input-output space of the coarse predictor, our ARNet defines the loss function in predictor and refinement network separately to achieve simultaneous supervision. Following [9], we optimize the coarse predictor network parameters with the mean-squared loss, which is denoted as the prediction loss $\mathcal{L}_{\mathcal{P}}$. Suppose $K$ is the number of joints in each frame, $N$ is the number of input frames and $T$ is the number of predicted frames, then $\mathcal{L}_{\mathcal{P}}$ can be written as:

$$\mathcal{L}_{\mathcal{P}} = \frac{1}{(N+T)K} \sum_{n=1}^{N+T} \sum_{k=1}^{K} ||h_{k,n}^{'} - h_{k,n}|| \tag{2}$$

where $h_{k,n}$ and $h_{k,n}^{'}$ respectively represent the ground-truth and predicted joint $k$ in frame n.

For the refinement network to produce the refined human motion sequences, we also adopt the mean-squared loss to optimize the network parameters. The mean-squared loss $\mathcal{L}_{\mathcal{R}}$ can be written as:

$$\mathcal{L}_{\mathcal{R}} = \frac{1}{(N+T)K} \sum_{n=1}^{N+T} \sum_{k=1}^{K} ||h_{k,n}^{''} - h_{k,n}|| \tag{3}$$

where $h_{k,n}$ indicates the ground-truth joint in frame $n$, $h_{k,n}^{''}$ is the refined corresponding joint. Our refiner is trained by minimizing the loss function.

The goal of our refinement network is to refine the coarse human motion prediction by utilizing the sequence-level refinement with the adversarial learning based error distribution augmentation. We utilise the minimax mechanism of adversarial loss to train the GAN:

$$\mathcal{L}_{\mathcal{D}} = \boldsymbol{E}[log\mathcal{D}(\delta_{real})] + \boldsymbol{E}[log(1 - \mathcal{D}(\mathcal{G}(\delta_{fake})] \tag{4}$$

$$\mathcal{L}_{\mathcal{G}} = \boldsymbol{E}[log(1 - \mathcal{D}(\mathcal{G}(\delta_{fake})] \tag{5}$$

where $\mathcal{L}_{\mathcal{D}}$ denotes the discriminator loss, $\mathcal{L}_{\mathcal{G}}$ is the generator loss, and $\boldsymbol{\delta}$ represents the error distribution.

In summary, we gather the predictor and refinement network together to train the whole network in an end-to-end way. As we adopt the adversarial refinement network behind the coarse predictor, the objective function consists of two parts:

$$\mathcal{L} = \mathcal{L}_{\mathcal{P}} + \boldsymbol{s} * \mathcal{L}_{\mathcal{R}} \tag{6}$$

where $\mathcal{L}_{\mathcal{P}}$ denotes the prediction loss, $\mathcal{L}_{\mathcal{R}}$ denotes the refinement loss, and the number of refinement stage $\boldsymbol{s}$ used in our adversarial refinement network will be shown in the ablation studies.

## 4    Experiments

### 4.1    Datasets and Evaluation Metrics

**H3.6m Dataset.** Human 3.6 Million (H3.6m) dataset [28] is the largest and most challenging mocap dataset which has 15 different daily actions performed by 7 males and females, including not only simple periodic actions such as walking and eating, but also complex aperiodic actions such as discussion and purchase. Following previous methods [29,9], the proposed algorithm is trained on subject 1,6,7,8,9,11 and tested on subject 5. There are 25 frames per second and each frame consists of a skeleton of 32 joints. Except for removing the global translations and rotations, some of the joints that do not move (*i.e.*, joints that do not bend) will be ignored as previous work [9].

**CMU-Mocap Dataset.** To be more convincing, we also conduct experiments on the CMU-Mocap dataset [29]. In order to achieve fair comparisons, we employ the same experimental settings as [29,9], including the pre-processing, data representation and training/testing splits.

**3DPW Dataset.** Recently, the 3D Pose in the Wild dataset (3DPW) [30] is released which contains around 51k frames with 3D annotations. The dataset is challenging as the scenarios are composed of indoor and outdoor activities. We follow [30,9] to split the dataset for comparable experimental results.

**Evaluation Metrics.** In order to make fair and comprehensive comparisons with previous work, we adopt the Mean Angle Error (MAE) between the predicted frames and the ground-truth frames in the angle space as the quantitative evaluation and visualize the prediction as the qualitative evaluation, which are the common evaluation metrics in human motion prediction [9].

### 4.2    Implementation Details

The proposed algorithm is implemented on Pytorch [31] and trained on a NVIDIA Tesla V100 GPU. We adopted the Adam [32] optimizer to train our model for about 50 epochs. The learning rate was set to 0.002 and the batch size was 256. To tackle the long-term temporal memory problems, we encode the complete time series by using Discrete Cosine Transform (DCT) [33] and discard the high-frequency jittering to maintain complete expression and smooth consistency of temporal domain information [9] at one time.

### 4.3    Quantitative Comparisons

We conduct quantitative comparisons on three human mocap datasets including H3.6m, 3DPW and CMU-Mocap between our ARNet and the state-of-the-art baselines. For fair comparisons with previous work [34,10,29,9,13], we feed 10 frames as inputs to predict the future 10 frames (400ms) for short-term prediction and the future 25 frames (1000ms) for long-term prediction.

**Table 1.** Short-term (80ms,160ms,320ms,400ms) human motion prediction measured in mean angle error (MAE) over 15 actions on H3.6m dataset

| | Walking | | | | Eating | | | | Smoking | | | | Discussion | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| milliseconds | 80 | 160 | 320 | 400 | 80 | 160 | 320 | 400 | 80 | 160 | 320 | 400 | 80 | 160 | 320 | 400 |
| Zero-velocity [13] | 0.39 | 0.68 | 0.99 | 1.15 | 0.27 | 0.48 | 0.73 | 0.86 | 0.26 | 0.48 | 0.97 | 0.95 | 0.31 | 0.67 | 0.94 | 1.04 |
| Residual sup. [13] | 0.28 | 0.49 | 0.72 | 0.81 | 0.23 | 0.39 | 0.62 | 0.76 | 0.33 | 0.61 | 1.05 | 1.15 | 0.31 | 0.68 | 1.01 | 1.09 |
| convSeq2Seq [29] | 0.33 | 0.54 | 0.68 | 0.73 | 0.22 | 0.36 | 0.58 | 0.71 | 0.26 | 0.49 | 0.96 | 0.92 | 0.32 | 0.67 | 0.94 | 1.01 |
| Retrospec [34] | 0.28 | 0.45 | 0.62 | 0.68 | 0.21 | 0.34 | 0.53 | 0.68 | 0.26 | 0.50 | 0.96 | 0.93 | 0.29 | 0.64 | 0.90 | 0.96 |
| AGED [10] | 0.22 | 0.36 | 0.55 | 0.67 | 0.17 | **0.28** | 0.51 | 0.64 | 0.27 | 0.43 | **0.82** | 0.84 | 0.27 | 0.56 | **0.76** | **0.83** |
| LTraiJ [9] | **0.18** | **0.31** | **0.49** | 0.56 | **0.16** | 0.29 | 0.50 | 0.62 | **0.22** | **0.41** | 0.86 | **0.80** | **0.20** | **0.51** | 0.77 | 0.85 |
| ARNet (Ours) | **0.18** | **0.31** | **0.49** | **0.55** | **0.16** | **0.28** | **0.49** | **0.61** | **0.22** | 0.42 | 0.86 | 0.81 | **0.20** | **0.51** | 0.81 | 0.89 |

| | Direction | | | | Greeting | | | | Phoning | | | | Posing | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| milliseconds | 80 | 160 | 320 | 400 | 80 | 160 | 320 | 400 | 80 | 160 | 320 | 400 | 80 | 160 | 320 | 400 |
| Zero-velocity [13] | 0.39 | 0.59 | 0.79 | 0.89 | 0.54 | 0.89 | 1.30 | 1.49 | 0.64 | 1.21 | 1.65 | 1.83 | 0.28 | 0.57 | 1.13 | 1.37 |
| Residual sup. [13] | 0.26 | 0.47 | 0.72 | 0.84 | 0.75 | 1.17 | 1.74 | 1.83 | 0.23 | 0.43 | 0.69 | 0.82 | 0.36 | 0.71 | 1.22 | 1.48 |
| convSeq2Seq [29] | 0.39 | 0.60 | 0.80 | 0.91 | 0.51 | 0.82 | 1.21 | 1.38 | 0.59 | 1.13 | 1.51 | 1.65 | 0.29 | 0.60 | 1.12 | 1.37 |
| Retrospec [34] | 0.40 | 0.61 | 0.77 | 0.86 | 0.52 | 0.86 | 1.26 | 1.43 | 0.59 | 1.11 | 1.47 | 1.59 | 0.26 | 0.54 | 1.14 | 1.41 |
| AGED [10] | **0.23** | **0.39** | **0.63** | **0.69** | 0.56 | 0.81 | 1.30 | 1.46 | **0.19** | **0.34** | **0.50** | **0.68** | 0.31 | 0.58 | 1.12 | 1.34 |
| LTraiJ [9] | 0.26 | 0.45 | 0.71 | 0.79 | 0.36 | 0.60 | 0.95 | 1.13 | 0.53 | 1.02 | 1.35 | 1.48 | 0.19 | 0.44 | 1.01 | 1.24 |
| ARNet (Ours) | **0.23** | 0.43 | 0.65 | 0.75 | **0.32** | **0.55** | **0.90** | **1.09** | 0.51 | 0.99 | 1.28 | 1.40 | **0.17** | **0.43** | **0.97** | **1.20** |

| | Purchases | | | | Sitting | | | | Sitting Down | | | | Taking Photo | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| milliseconds | 80 | 160 | 320 | 400 | 80 | 160 | 320 | 400 | 80 | 160 | 320 | 400 | 80 | 160 | 320 | 400 |
| Zero-velocity [13] | 0.62 | 0.88 | 1.19 | 1.27 | 0.40 | 1.63 | 1.02 | 1.18 | 0.39 | 0.74 | 1.07 | 1.19 | 0.25 | 0.51 | 0.79 | 0.92 |
| Residual sup. [13] | 0.51 | 0.97 | 1.07 | 1.16 | 0.41 | 1.05 | 1.49 | 1.63 | 0.39 | 0.81 | 1.40 | 1.62 | 0.24 | 0.51 | 0.90 | 1.05 |
| convSeq2Seq [29] | 0.63 | 0.91 | 1.19 | 1.29 | 0.39 | 0.61 | 1.02 | 1.18 | 0.41 | 0.78 | 1.16 | 1.31 | 0.23 | 0.49 | 0.88 | 1.06 |
| Retrospec [34] | 0.59 | 0.84 | 1.14 | 1.19 | 0.40 | 0.64 | 1.04 | 1.22 | 0.41 | 0.77 | 1.14 | 1.29 | 0.27 | 0.52 | 0.80 | 0.92 |
| AGED [10] | 0.46 | 0.78 | 1.01 | **1.07** | 0.41 | 0.76 | 1.05 | 1.19 | 0.33 | 0.62 | 0.98 | 1.10 | 0.23 | 0.48 | 0.81 | 0.95 |
| LTraiJ [9] | 0.43 | 0.65 | 1.05 | 1.13 | 0.29 | 0.45 | **0.80** | **0.97** | 0.30 | **0.61** | 0.90 | 1.00 | 0.14 | 0.34 | 0.58 | 0.70 |
| ARNet (Ours) | **0.36** | **0.60** | **1.00** | 1.11 | **0.27** | **0.44** | **0.80** | **0.97** | **0.29** | **0.61** | **0.87** | **0.97** | **0.13** | **0.33** | **0.55** | **0.67** |

| | Waiting | | | | Walking Dog | | | | Walking Together | | | | Average | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| milliseconds | 80 | 160 | 320 | 400 | 80 | 160 | 320 | 400 | 80 | 160 | 320 | 400 | 80 | 160 | 320 | 400 |
| Zero-velocity [13] | 0.34 | 0.67 | 1.22 | 1.47 | 0.60 | 0.98 | 1.36 | 1.50 | 0.33 | 0.66 | 0.94 | 0.99 | 0.40 | 0.78 | 1.07 | 1.21 |
| Residual sup. [13] | 0.28 | 0.53 | 1.02 | 1.14 | 0.56 | 0.91 | 1.26 | 1.40 | 0.31 | 0.58 | 0.87 | 0.91 | 0.36 | 0.67 | 1.02 | 1.15 |
| convSeq2Seq [29] | 0.30 | 0.62 | 1.09 | 1.30 | 0.59 | 1.00 | 1.32 | 1.44 | 0.27 | 0.52 | 0.71 | 0.74 | 0.38 | 0.68 | 1.01 | 1.13 |
| Retrospec [34] | 0.33 | 0.65 | 1.12 | 1.30 | 0.53 | 0.87 | 1.16 | 1.33 | 0.28 | 0.52 | 0.68 | 0.71 | 0.37 | 0.66 | 0.98 | 1.10 |
| AGED [10] | 0.24 | 0.50 | 1.02 | **1.13** | 0.50 | 0.81 | 1.15 | **1.27** | 0.23 | 0.41 | 0.56 | 0.62 | 0.31 | 0.54 | 0.85 | 0.97 |
| LTraiJ [9] | 0.23 | 0.50 | 0.91 | 1.14 | 0.46 | 0.79 | 1.12 | 1.29 | 0.15 | 0.34 | **0.52** | **0.57** | 0.27 | 0.51 | 0.83 | 0.95 |
| ARNet (Ours) | **0.22** | **0.48** | **0.90** | **1.13** | **0.45** | **0.78** | **1.11** | **1.27** | **0.13** | **0.33** | 0.53 | 0.58 | **0.25** | **0.49** | **0.80** | **0.92** |

**Short-term Prediction on H3.6m.** H3.6m is the most challenging dataset for human motion prediction. Table 1 shows the quantitative comparisons for short-term human motion prediction between our ARNet and a series of baselines including Zero-velocity [13], RRNN[13], convSeq2Seq[29], Retrospec[34], AGED [10] and LTraiJ [9] on H3.6m dataset. We computed the mean angle error (MAE) on 15 actions by measuring the euclidean distance between the ground-truth and prediction at 80ms, 160ms, 320ms, 400ms for short-term evaluation. The results in bold show that our method outperforms both of the state-of-the-art chain-structured baseline AGED and the feed-forward baseline LTraiJ.

Compared with the state-of-the-art feed-forward baseline LTraiJ [9], in Table 1, the proposed ARNet clearly outperforms the feed-forward baseline LTraiJ on average for short-term human motion prediction. Different from LTraiJ which
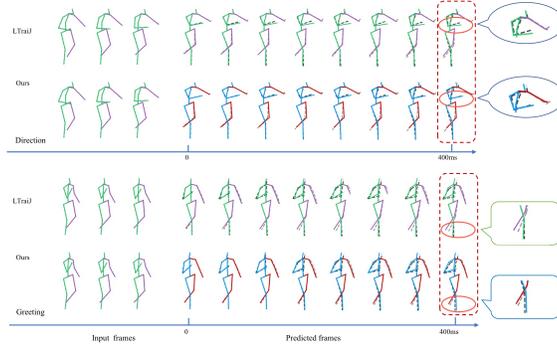
**Fig. 2. Visual comparisons for short-term human motion prediction on H3.6m dataset.** We compare our proposed ARNet with the state-of-the-art feed-forward baseline LTraiJ [9] which is the best performing method for short-term prediction (400ms). The left few frames represent the input human motion sequence. From top to bottom, we show the final predictions obtained by the feed-forward baseline LTraiJ represented as green-purple skeletons and our proposed ARNet represented as red-blue skeletons respectively on two challenging aperiod actions (e.g.,Direction and Greeting). Marked in red circles, our predictions better match the ground-truth shown as the gray dotted skeletons

adopts the single-stage predictor without refinement network, our ARNet obtains better performance especially on aperiodic actions (e.g. Directions, Greeting, Phoning and so on). It is difficult to model this type of actions which involved multiple small movements and high acceleration during human motion especially at the end of human limbs. In addition, due to the stable change of periodic behavior, the traditional feed-forward deep network can also achieve competitive results on periodic actions (such as walking, eating and smoking), but we note that our ARNet further improves the accuracy of prediction. The results validate that the coarse-to-fine design enables our ARNet to correct the error joints in human motion prediction and outperform the existing feed-forward baseline on almost all actions.

Compared with the state-of-the-art chain-structured baseline AGED [10], which utilises chain-structured RNNs as the predictor with two different discriminators, our ARNet still outperforms it on almost all action categories for short-term human motion prediction within 400ms as shown in Table 1. The results show the superiority of our ARNet over the best performing chain-structured methods for short-term human motion prediction tasks.

**Long-term Prediction on H3.6m.** Additionlly, we also quantitatively evaluate the long-term prediction performance of our proposed ARNet at 560ms and 1000ms as shown in Table 2. The results measured in MAE demonstrate that our method still outperforms the state-of-art feed-forward baseline LTraiJ [9]

**Table 2.** Long-term (560ms, 1000ms) human motion prediction on H3.6m dataset

| | Walking | | Eating | | Smoking | | Discussion | | Average | |
|---|---|---|---|---|---|---|---|---|---|---|
| milliseconds | 560 | 1000 | 560 | 1000 | 560 | 1000 | 560 | 1000 | 560 | 1000 |
| Zero-velocity [13] | 1.35 | 1.32 | 1.04 | 1.38 | 1.02 | 1.69 | 1.41 | 1.96 | 1.21 | 1.59 |
| Residual sup. [13] | 0.93 | 1.03 | 0.95 | 1.08 | 1.25 | 1.50 | 1.43 | 1.69 | 1.14 | 1.33 |
| AGED [10] | 0.78 | 0.91 | 0.86 | **0.93** | 1.06 | **1.21** | **1.25** | **1.30** | 0.99 | **1.09** |
| Retrospec [34] | NA | 0.79 | NA | 1.16 | NA | 1.71 | NA | 1.72 | NA | 1.35 |
| LTraiJ [9] | **0.65** | **0.67** | 0.76 | 1.12 | 0.87 | 1.57 | 1.33 | 1.70 | 0.90 | 1.27 |
| ARNet (Ours) | **0.65** | 0.69 | **0.72** | 1.07 | **0.86** | 1.51 | **1.25** | 1.68 | **0.88** | 1.24 |

in long-term human motion prediction on almost action categories as shown in bold. Nevertheless, the MAE of the chain-structured AGED [10] is lower than ours in 1000 milliseconds. We will further examine the results by visualizing the motion sequences obtained by our proposed ARNet and the chain-structured baseline AGED in the later section to provide a qualitative comparison.

**3DPW & CMU-Mocap.** We also conduct experiments on other two human mocap datasets to prove the robustness of our method. Table 3 shows that our method consistently achieves promising improvements compared with other baselines on 3DPW dataset which contains indoor and outdoor activities for both short-term and long-term human motion predictions. As for CMU-Mocap dataset, the results in Table 4 illustrate that our method has better performance on almost action types and outperforms the state-of-the-art methods on average.

**Table 3.** Short-term and long-term human motion predictions on 3DPW dataset

| milliseconds | 200 | 400 | 600 | 800 | 1000 |
|---|---|---|---|---|---|
| Residual sup. [13] | 1.85 | 2.37 | 2.46 | 2.51 | 2.53 |
| convSeq2Seq [29] | 1.24 | 1.85 | 2.13 | 2.23 | 2.26 |
| LTraiJ [9] | 0.64 | **0.95** | 1.12 | 1.22 | 1.27 |
| ARNet (Ours) | **0.62** | **0.95** | **1.11** | **1.20** | **1.25** |

### 4.4   Qualitative Visualizations

**Short-term Prediction on H3.6m.** To evaluate our method qualitatively, we firstly visualize the representative comparisons on Directions and Greeting which belong to challenging aperiodic actions in H3.6m dataset as shown in Figure 2. Given 10 observed frames for each action as motion seeds, which are represented as green-purple skeletons at the left part, we compare our ARNet represented as red-blue skeletons with the best quantitatively performing feed-forward baseline LTraiJ [9] shown as green-purple skeletons for short-term prediction (400 million seconds) as illustrated in Table 1. The dotted rectangles mark that our

**Table 4.** Short-term and long-term human motion predictions on CMU-Mocap dataset

| | Basketball | | | | | Basketball Signal | | | | | Directing Traffic | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| milliseconds | 80 | 160 | 320 | 400 | 1000 | 80 | 160 | 320 | 400 | 1000 | 80 | 160 | 320 | 400 | 1000 |
| Residual sup. [13] | 0.50 | 0.80 | 1.27 | 1.45 | 1.78 | 0.41 | 0.76 | 1.32 | 1.54 | 2.15 | 0.33 | 0.59 | 0.93 | 1.10 | 2.05 |
| convSeq2Seq [29] | 0.37 | 0.62 | 1.07 | 1.18 | 1.95 | 0.32 | 0.59 | 1.04 | 1.24 | 1.96 | 0.25 | 0.56 | 0.89 | 1.00 | 2.04 |
| LTraiJ [9] | 0.33 | 0.52 | 0.89 | **1.06** | **1.71** | 0.11 | 0.20 | 0.41 | 0.53 | 1.00 | 0.15 | 0.32 | 0.52 | 0.60 | 2.00 |
| ARNet (Ours) | **0.31** | **0.48** | **0.87** | 1.08 | **1.71** | **0.10** | **0.17** | **0.35** | **0.48** | **1.06** | **0.13** | **0.28** | **0.47** | **0.58** | **1.80** |

| | Jumping | | | | | Running | | | | | Soccer | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| milliseconds | 80 | 160 | 320 | 400 | 1000 | 80 | 160 | 320 | 400 | 1000 | 80 | 160 | 320 | 400 | 1000 |
| Residual sup. [13] | 0.33 | 0.50 | 0.66 | 0.75 | 1.00 | 0.29 | 0.51 | 0.88 | 0.99 | 1.72 | 0.56 | 0.88 | 1.77 | 2.02 | 2.4 |
| convSeq2Seq [29] | **0.28** | **0.41** | **0.52** | **0.57** | **0.67** | 0.26 | 0.44 | 0.75 | 0.87 | 1.56 | 0.39 | 0.6 | 1.36 | 1.56 | 2.01 |
| LTraiJ [9] | 0.33 | 0.55 | 0.73 | 0.74 | 0.95 | 0.18 | 0.29 | 0.61 | 0.71 | 1.40 | 0.31 | 0.49 | 1.23 | 1.39 | 1.80 |
| ARNet (Ours) | 0.30 | 0.50 | 0.60 | 0.61 | 0.72 | **0.16** | **0.26** | **0.57** | **0.67** | **1.22** | **0.29** | **0.47** | **1.21** | **1.38** | **1.70** |

| | Walking | | | | | Washwindow | | | | | Average | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| milliseconds | 80 | 160 | 320 | 400 | 1000 | 80 | 160 | 320 | 400 | 1000 | 80 | 160 | 320 | 400 | 1000 |
| Residual sup. [13] | 0.35 | 0.47 | 0.60 | 0.65 | 0.88 | 0.30 | 0.46 | 0.72 | 0.91 | 1.36 | 0.38 | 0.62 | 1.02 | 1.18 | 1.67 |
| convSeq2Seq [29] | 0.35 | 0.44 | 0.45 | 0.50 | 0.78 | 0.30 | 0.47 | 0.80 | 1.01 | 1.39 | 0.32 | 0.52 | 0.86 | 0.99 | 1.55 |
| LTraiJ [9] | 0.33 | 0.45 | 0.49 | 0.53 | 0.61 | 0.22 | 0.33 | 0.57 | 0.75 | 1.20 | 0.25 | 0.39 | 0.68 | 0.79 | 1.33 |
| ARNet (Ours) | **0.32** | **0.41** | **0.39** | **0.41** | **0.56** | **0.20** | **0.27** | **0.51** | **0.69** | **1.07** | **0.23** | **0.37** | **0.65** | **0.77** | **1.29** |

predictions better match the ground-truth which is represented as gray dotted skeletons. The qualitative comparison further demonstrates that our ARNet possesses the ideal error-correction ability to generate high-quality prediction, especially for the joints at the end of body which contain multiple small movements on aperiodic actions.

**Long-term Prediction on H3.6m.** Figure 3 visualizes the comparisons between chain-structured baselines RRNN [13] and AGED [13] on Phoning, which belongs to aperiodic actions in H3.6m dataset for long-term prediction (4 seconds). As marked by the red rectangles, our proposed ARNet is still able to predict the motion dynamics when the RRNN converges to mean pose. Meanwhile, the AGED drifts away on the foot joints compared with the ground-truth. The visualised results demonstrate that our ARNet outperforms the chain-structured baselines in long-term prediction.

## 5   Ablation Studies

### 5.1   Different Components in Our ARNet

In order to verify the effectiveness of the different components in our model, we perform comprehensive ablation studies as shown in Table 5. Specifically, we compare our ARNet with three baselines: the 1-stage CoarseNet, the 2-stage CoarseNet without future information as refinement and the 2-stage RefineNet with future information and traditional training strategy. The 1-stage CoarseNet denotes that there only exists single coarse predictor module without other components in the whole framework. We utilize the LTraiJ network [9] as our coarse
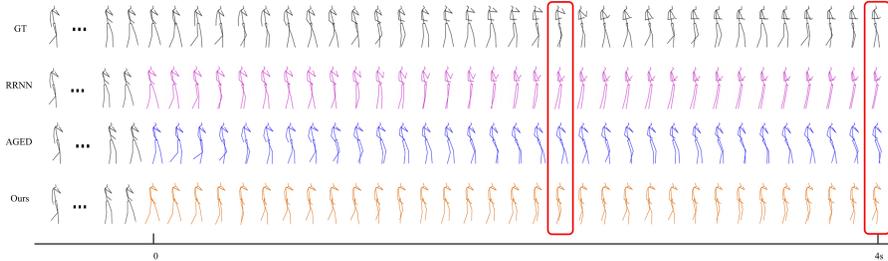
**Fig. 3. Visual comparisons for long-term human motion prediction on H3.6m dataset.** From top to bottom,we show the corresponding ground-truth shown in grey skeletons, the final predictions obtained by RRNN [13] , AGED [10] and our approach on Phoning which belongs to the aperiodic action. The left gray skeletons represent the input motion sequences. Marked in red rectangles, the baseline RRNN converges to mean pose and the baseline AGED drifts away on the foot joints compared with the ground-truth. Our ARNet generates more accurate long-term human motion prediction relatively. Best viewed in color with zoom

predictor. Due to the coarse-to-fine 2-stage structure of our ARNet, the inference time of our ARNet is 56.2ms, which is slightly longer than the 45.4ms of 1-stage CoarseNets on GPU V100. Moreover, another baseline is the 2-stage CoarseNet without future information refinement, which increase the number of layers by simply cascading two 1-stage CoarseNets, utilise the same training strategy as the single coarse predictor by back-propagating the gradient all the way to the beginning. Although the parameters of our ARNet is same as the 2-stage CoarseNet which is twice that of 1-stage CoarseNets, the results show that stacking multi-layers with traditional training strategy fails to improve the performance in a further step and even achieved worse prediction due to over-fitting occurred in stacked feed-forward deep network. Then, the 2-stage RefineNet without adversarial error augmentation leads to improvement over the previous two baselines. Our adversarial refinement network shows the superior performance compared with single-stage model, 2-stage model without refinement and refinement network without adversarial training strategy.

### 5.2   Multi-stage Analysis

We also evaluate the impact of number of stages adopted in our adversarial refinement model by calculating the MAE over 15 actions. The foregoing results in the Table 6 indicate that the 2-stage refined model design, in general, utilising the output space of previous stage, is simple enough to learn the rich representation and achieves superior results in most cases. The reason is that concatenating more than 2 stages refinement module faces up over-fitting problems and fails to further improve the human motion prediction performance. Taking the efficiency and simplicity into account, we employ the 2-stage adversarial refinement network as the final model design.

**Table 5.** Ablation study for refined model design and adversarial training strategy. We compared the results measured in MAE of our model with the 1-stage CoarseNet, the 2-stage CoarseNet without future information as refinement and the 2-stage RefineNet with traditional training strategy on H3.6m dataset

| | Direction | | | | | | Posing | | | | | | Greeting | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| milliseconds | 80 | 160 | 320 | 400 | 560 | 1000 | 80 | 160 | 320 | 400 | 560 | 1000 | 80 | 160 | 320 | 400 | 560 | 1000 |
| 1-stage CoarseNet | 0.26 | 0.45 | 0.71 | 0.79 | 0.88 | 1.29 | 0.19 | 0.44 | 1.01 | 1.24 | 1.44 | 1.64 | 0.36 | 0.60 | 0.95 | 1.13 | 1.51 | 1.70 |
| 2-stage CoarseNet | 0.25 | 0.45 | 0.67 | 0.78 | 0.88 | 1.30 | 0.19 | 0.46 | 1.01 | 1.26 | 1.42 | 1.68 | 0.34 | 0.60 | 0.94 | 1.11 | 1.66 | 1.92 |
| 2-stage RefineNet | 0.25 | 0.44 | 0.67 | 0.77 | 0.86 | 1.27 | 0.19 | 0.43 | 0.99 | 1.23 | 1.42 | 1.63 | 0.34 | 0.58 | 0.92 | 1.10 | 1.49 | 1.63 |
| ARNet | **0.23** | **0.43** | **0.65** | **0.75** | **0.85** | **1.23** | **0.17** | **0.43** | **0.97** | **1.20** | **1.41** | **1.60** | **0.31** | **0.55** | **0.90** | **1.08** | **1.46** | **1.56** |

| | Greeting | | | | | | Phoning | | | | | | Average(on 15 actions) | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| milliseconds | 80 | 160 | 320 | 400 | 560 | 1000 | 80 | 160 | 320 | 400 | 560 | 1000 | 80 | 160 | 320 | 400 | 560 | 1000 |
| 1-stage CoarseNet | 0.36 | 0.60 | 0.95 | 1.13 | 1.51 | 1.70 | 0.53 | 1.02 | 1.35 | 1.48 | 1.45 | 1.68 | 0.27 | 0.51 | 0.83 | 0.95 | 1.18 | 1.59 |
| 2-stage CoarseNet | 0.34 | 0.60 | 0.94 | 1.11 | 1.66 | 1.92 | 0.53 | 1.02 | 1.34 | 1.48 | 1.58 | 1.98 | 0.27 | 0.52 | 0.83 | 0.95 | 1.20 | 1.61 |
| 2-stage RefineNet | 0.34 | 0.58 | 0.94 | 1.10 | 1.48 | 1.64 | 0.52 | 1.01 | 1.33 | 1.46 | 1.42 | 1.65 | 0.27 | 0.50 | 0.82 | 0.94 | 1.17 | 1.58 |
| ARNet | **0.31** | **0.55** | **0.90** | **1.08** | **1.46** | **1.56** | **0.50** | **0.99** | **1.28** | **1.40** | **1.41** | **1.60** | **0.25** | **0.49** | **0.80** | **0.92** | **1.16** | **1.57** |

**Table 6.** Ablation study of adversarial refinement network with different number of stages. We compared the results measured in MAE on H3.6m dataset

| | Direction | | | | | | Posing | | | | | | Greeting | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| milliseconds | 80 | 160 | 320 | 400 | 560 | 1000 | 80 | 160 | 320 | 400 | 560 | 1000 | 80 | 160 | 320 | 400 | 560 | 1000 |
| 2-stage | **0.23** | **0.43** | **0.65** | **0.75** | 0.85 | **1.23** | **0.17** | **0.43** | **0.97** | **1.20** | **1.41** | **1.60** | **0.31** | 0.55 | 0.90 | **1.08** | **1.46** | **1.56** |
| 3-stage | 0.25 | 0.46 | 0.64 | 0.75 | **0.84** | 1.50 | 0.18 | 0.44 | 1.00 | 1.25 | 1.71 | 2.64 | 0.32 | **0.54** | **0.89** | 1.12 | 1.52 | 1.75 |
| 4-stage | 0.25 | 0.46 | 0.68 | 0.77 | 1.02 | 1.70 | 0.19 | 0.46 | 1.05 | 1.28 | 1.86 | 3.03 | 0.33 | 0.56 | 0.93 | 1.15 | 1.56 | 1.82 |

| | Greeting | | | | | | Phoning | | | | | | Average(on 15 actions) | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| milliseconds | 80 | 160 | 320 | 400 | 560 | 1000 | 80 | 160 | 320 | 400 | 560 | 1000 | 80 | 160 | 320 | 400 | 560 | 1000 |
| 2-stage | **0.31** | 0.55 | 0.90 | **1.08** | **1.46** | **1.56** | **0.50** | **0.99** | **1.28** | **1.40** | **1.41** | **1.60** | **0.25** | **0.49** | **0.80** | **0.92** | **1.16** | **1.57** |
| 3-stage | 0.32 | **0.54** | **0.89** | 1.12 | 1.52 | 1.75 | 0.52 | 1.02 | 1.36 | 1.45 | 1.49 | 1.80 | 0.25 | 0.49 | 0.83 | 0.95 | 1.17 | 1.58 |
| 4-stage | 0.33 | 0.56 | 0.93 | 1.15 | 1.56 | 1.82 | 0.52 | 0.99 | 1.33 | 1.48 | 1.49 | 1.76 | 0.27 | 0.50 | 0.83 | 0.95 | 1.17 | 1.58 |

# 6   Conclusions

In this paper, we introduce an Adversarial Refinement Network (ARNet) to forecast more accurate human motion sequence in a coarse-to-fine manner. We adopt a refinement network behind the single-stage coarse predictor to generate finer human motion. Meanwhile, we utilise an adversarial learning strategy to enhance the generalization ability of the refinement network. Experimental results on the challenging benchmark H3.6m, CMU-Mocap and 3DPW datasets show that our proposed ARNet outperforms the state-of-the-art approaches in both short-term and long-term predictions especially on the complex aperiodic actions. Our adversarial refinement network shows promising potential for feed-forward deep network to deal with rich representation in a further step on other areas.

# References

1. Koppula, H.S., Saxena, A.: Anticipating human activities for reactive robotic response. In: IROS. (2013)
2. Saquib Sarfraz, M., Schumann, A., Eberle, A., Stiefelhagen, R.: A pose-sensitive embedding for person re-identification with expanded cross neighborhood re-ranking. In: CVPR. (2018)
3. Elhayek, A., Kovalenko, O., Murthy, P., Malik, J., Stricker, D.: Fully automatic multi-person human motion capture for vr applications. In: International Conference on Virtual Reality and Augmented Reality. (2018)
4. Yuminaka, Y., Mori, T., Watanabe, K., Hasegawa, M., Shirakura, K.: Non-contact vital sensing systems using a motion capture device: medical and healthcare applications. In: Key engineering materials. (2016)
5. Paden, B., Čáp, M., Yong, S.Z., Yershov, D., Frazzoli, E.: A survey of motion planning and control techniques for self-driving urban vehicles. IEEE Transactions on intelligent vehicles (2016)
6. Gong, H., Sim, J., Likhachev, M., Shi, J.: Multi-hypothesis motion planning for visual object tracking. In: ICCV. (2011)
7. Wang, J.M., Fleet, D.J., Hertzmann, A.: Gaussian process dynamical models for human motion. IEEE transactions on pattern analysis and machine intelligence **30** (2007) 283–298
8. Tang, Y., Ma, L., Liu, W., Zheng, W.: Long-term human motion prediction by modeling motion context and enhancing motion dynamic. In: IJCAI. (2018)
9. Mao, W., Liu, M., Salzmann, M., Li, H.: Learning trajectory dependencies for human motion prediction. In: ICCV. (2019)
10. Gui, L.Y., Wang, Y.X., Liang, X., Moura, J.M.: Adversarial geometry-aware human motion prediction. In: ECCV. (2018)
11. Fragkiadaki, K., Levine, S., Felsen, P., Malik, J.: Recurrent network models for human dynamics. In: ICCV. (2015)
12. Jain, A., Zamir, A.R., Savarese, S., Saxena, A.: Structural-rnn: Deep learning on spatio-temporal graphs. In: CVPR. (2016)
13. Martinez, J., Black, M.J., Romero, J.: On human motion prediction using recurrent neural networks. In: CVPR. (2017)
14. Tome, D., Russell, C., Agapito, L.: Lifting from the deep: Convolutional 3d pose estimation from a single image. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2017) 2500–2509
15. Chen, Y., Wang, Z., Peng, Y., Zhang, Z., Yu, G., Sun, J.: Cascaded pyramid network for multi-person pose estimation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. (2018) 7103–7112
16. Fieraru, M., Khoreva, A., Pishchulin, L., Schiele, B.: Learning to refine human pose estimation. In: CVPR-W. (2018)
17. Moon, G., Chang, J.Y., Lee, K.M.: Posefix: Model-agnostic general human pose refinement network. In: CVPR. (2019)
18. Wan, Q., Qiu, W., Yuille, A.L.: Patch-based 3d human pose refinement. arXiv preprint arXiv:1905.08231 (2019)
19. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: NIPS. (2014)
20. Deng, K., Fei, T., Huang, X., Peng, Y.: Irc-gan: introspective recurrent convolutional gan for text-to-video generation. In: IJCAI. (2019)

21. Balaji, Y., Min, M.R., Bai, B., Chellappa, R., Graf, H.P.: Conditional gan with discriminative filter generation for text-to-video synthesis. In: IJCAI. (2019)
22. Vankadari, M., Kumar, S., Majumder, A., Das, K.: Unsupervised learning of monocular depth and ego-motion using conditional patchgans. In: IJCAI. (2019)
23. Chu, W., Hung, W.C., Tsai, Y.H., Cai, D., Yang, M.H.: Weakly-supervised caricature face parsing through domain adaptation. In: ICIP. (2019)
24. Zhang, X., Wang, Q., Zhang, J., Zhong, Z.: Adversarial autoaugment. In: ICLR. (2020)
25. Frid-Adar, M., Klang, E., Amitai, M., Goldberger, J., Greenspan, H.: Gan-based data augmentation for improved liver lesion classification. (2018)
26. Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. (2017)
27. Yan, S., Xiong, Y., Lin, D.: Spatial temporal graph convolutional networks for skeleton-based action recognition. In: Thirty-second AAAI conference on artificial intelligence. (2018)
28. Ionescu, C., Papava, D., Olaru, V., Sminchisescu, C.: Human3.6m: Large scale datasets and predictive methods for 3d human sensing in natural environments. TPAMI **36** (2014) 1325–1339
29. Li, C., Zhang, Z., Lee, W.S., Lee, G.H.: Convolutional sequence to sequence model for human dynamics. In: CVPR. (2018)
30. von Marcard, T., Henschel, R., Black, M., Rosenhahn, B., Pons-Moll, G.: Recovering accurate 3d human pose in the wild using imus and a moving camera. In: ECCV. (2018)
31. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in pytorch. In: NIPS-W. (2017)
32. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
33. Akhter, I., Sheikh, Y., Khan, S., Kanade, T.: Nonrigid structure from motion in trajectory space. In: NIPS. (2009)
34. Dong, M., Xu, C.: On retrospecting human dynamics with attention. In: IJCAI. (2019)