

SSHFD: Single Shot Human Fall Detection with Occluded Joints Resilience

Umar Asif¹ and Stefan Von Cavallar² and Jianbin Tang³ and Stefan Harrer⁴

Abstract. Falling can have fatal consequences for elderly people especially if the fallen person is unable to call for help due to loss of consciousness or any injury. Automatic fall detection systems can assist through prompt fall alarms and by minimizing the fear of falling when living independently at home. Existing vision-based fall detection systems lack generalization to unseen environments due to challenges such as variations in physical appearances, different camera viewpoints, occlusions, and background clutter. In this paper, we explore ways to overcome the above challenges and present Single Shot Human Fall Detector (SSHFD), a deep learning based framework for automatic fall detection from a single image. This is achieved through two key innovations. First, we present a human pose based fall representation which is invariant to appearance characteristics. Second, we present neural network models for 3d pose estimation and fall recognition which are resilient to missing joints due to occluded body parts. Experiments on public fall datasets show that our framework successfully transfers knowledge of 3d pose estimation and fall recognition learnt purely from synthetic data to unseen real-world data, showcasing its generalization capability for accurate fall detection in real-world scenarios.

1 Introduction

Falling on the ground is considered to be one of the most critical dangers for elderly people living alone at home which can cause serious injuries and restricts normal activities because of the fear of falling again [4]. Automated fall detection systems can produce prompt alerts in hazardous situations. They also allow automatic collection and reporting of fall incidents which can be used to analyse the causes of falls, thus improving the quality of life for people with mobility constraints and limited supervision. Vision-based systems provide a low cost solution to fall detection. They do not cause sensory side effects on the human health and do not affect the normal routines of elderly people as observed in systems using wearable devices [19]. In a typical fall detection approach, human regions are detected from the visual data and used to learn features to distinguish fall from other activities. Existing methods such as [15] learn fall representations using physical appearance based features extracted from video data. However, appearance based features suffer from poor generalization in real-world environments due to large variations in appearance characteristics, different camera viewpoints, and background clutter. Furthermore, due to the unavailability of large-scale public fall datasets, most of the existing fall detectors are trained and evaluated using simulated environments or using

restricted datasets (which cannot be shared publicly due to privacy concerns). Therefore, these methods do not exhibit generalization capabilities for fall detection in unseen real-world environments. In this paper, we explore ways to overcome the above challenges and present a deep learning framework termed “Single Shot Human Fall Detector (SSHFD)” for accurate fall detection in unseen real-world environments. The main contributions of this paper are as follows:

1. We present a human pose based fall representation which is invariant to appearance characteristics, backgrounds, lighting conditions, and spatial locations of people in the scene. Experiments show that neural network models trained on our 2d-pose and 3d-pose based fall representations successfully generalize to unseen real-world environments for fall recognition.
2. We present neural network models for 3d pose estimation and fall recognition which are robust to partial occlusions. Experiments show that our models successfully recover joints information from occluded body parts, and accurately recognize fall poses from incomplete input data.
3. We evaluate our framework on real-world public fall datasets, where we show that our framework when trained using only synthetic data, shows excellent generalization capabilities of fall recognition on unseen real-world data.

2 Related Work

Existing vision-based fall detection approaches detect human regions in the scene and use visual information from the detected regions to learn features for fall recognition. For instance, the method of [14] generated human bounding boxes through background-foreground subtraction and compared the visual content of the boxes in consecutive frames of the videos of the MultiCam fall dataset [2] to detect fall events. The method of [18] compared multiple bounding boxes to distinguish between different events (e.g., standing, sitting, and fall). The work of [10] used a fuzzy neural network classifier for fall detection. The methods of [15] and [8] used motion segmentation to detect human regions in the scene and combined visual appearance and shape information from the detected regions to learn features for fall recognition. However, errors in background-foreground subtraction or motion segmentation (e.g., due to small or no change in the visual content between subsequent image frames) degrade the accuracy of these methods. To overcome this challenge, the method of [9] used cues from multiple cameras and produced fall decisions through voting among different viewpoints. However, this approach requires accurate synchronization between the individual cameras. Other methods such as [5, 13] used Kinect depth maps to learn 3d features for fall recognition. However, these methods are restricted in real-world deployment due to hardware limitations (e.g. limited depth

¹ IBM Research Australia, email: umarasif@au1.ibm.com

² IBM Research Australia, email: svcavallar@au1.ibm.com

³ IBM Research Australia, email: jbtang@au1.ibm.com

⁴ IBM Research Australia, email: sharrer@au1.ibm.com

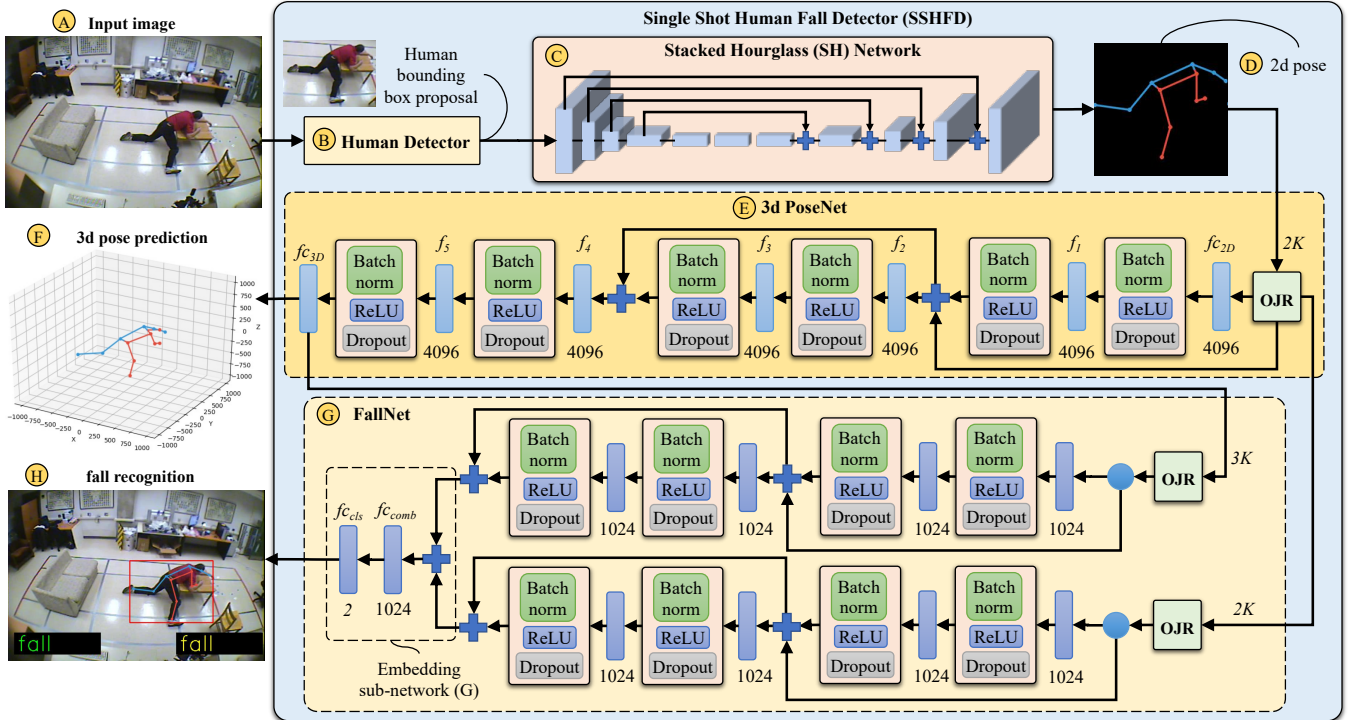


Figure 1: Overview of our Single Shot Human Fall Detector (SSHFD). Given a single RGB image of the scene (A), SSHFD generates human proposals (B) which are fed into a Stacked Hourglass network (C) for 2d pose prediction. Next, the predicted 2d pose (D) is fed into a neural network (E) for 3d pose prediction (F). Finally, the 2d pose and the 3d pose information are fed into a neural network (G) for fall recognition (H). Our models integrate Occluded Joints Resilience (OJR) modules which make the models robust to missing information in the pose data.

sensing range). Compared to existing methods, our work differs in several ways. **First**, our framework learns pose based fall representations which are invariant to appearance characteristics. This enables our framework to successfully transfer fall recognition knowledge learned from pure synthetic data to real-world data with unknown backgrounds and different human actors. **Second**, our framework integrates a 3d pose estimator which predicts 3d pose information from 2d pose. The combined 2d and 3d pose knowledge enables our framework to successfully handle ambiguities in the 2d pose (under different camera viewpoints), without requiring multiple camera setups or depth sensor technologies. **Finally**, our neural network models for 3d pose estimation and fall recognition are resilient to missing information in the pose data. This enables our framework to accurately discriminate between fall and no-fall cases from human poses under occlusions.

3 The Proposed Framework (SSHFD)

Fig. 1 shows the overall architecture of our framework which has three main modules. **i)** 2d pose estimation, which takes an RGB image of the scene as input and produces body joints locations in 2d image space, **ii)** 3d pose estimation, which takes 2d pose as input and predicts joints locations in 3d Cartesian space, and **iii)** Fall recognition, which combines 2d pose and 3d pose data and predicts probabilities with respect to the target classes. In the following, we describe in detail the individual components of our framework.

3.1 The Proposed Fall Representation

Our fall representation is based on joints locations in 2d image space and 3d Cartesian space. We normalize the 2d pose by transforming

the joints estimates (predicted in the scene image) to a fixed reference image of 224×224 dimensions as shown in Fig. 1-D. The normalized 2d pose is then used to predict joints locations in a Cartesian space of size $1000 \times 1000 \times 1000mm^3$ as shown in Fig. 1-F. The 3d predictions are normalized with respect to the hip joint.

3.2 The Proposed 2d Pose Estimation (Fig. 1)

Our 2d pose estimator is composed of two main modules: **i)** a human detector [6], which produces human bounding box proposals from the input image, and **ii)** a Stacked Hourglass (SH) network [16], which predicts body joints 2d locations and their corresponding confidence scores. The SH network is trained using ground truth labelled in terms of $W \times H \times K$ -dimensional heatmaps (\mathcal{H}), where W and H represent the width and height of the heatmap and K represents the number of joints. We used $K = 17$ joint types as per the format used in [6]. The heatmap (\mathcal{H}_k) for a joint $k \in \{1, \dots, K\}$ is generated by centering a Gaussian kernel around the joint's pixel position (x_k, y_k) . It is given by:

$$\mathcal{H}_k(x, y) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{[(x - x_k)^2 + (y - y_k)^2]}{2\sigma^2}\right), \quad (1)$$

where σ is a hyper-parameter for spatial variance. We set $\sigma = 4$ in our experiments. The training objective function of the SH network is defined by:

$$\mathcal{L}_{2d} = \frac{1}{K} \sum_{k=1}^K \|\mathcal{H}_k - \hat{\mathcal{H}}_k\|_2^2, \quad (2)$$

where $\hat{\mathcal{H}}_k$ represents the predicted confidence map for the k th joint.

3.3 The Proposed 3D Pose Estimation (Fig. 1-E)

Here, the goal is to estimate K body joints in 3d Cartesian space $\mathbf{Q} \in \mathbb{R}^{3K}$ given a 2d input $\mathbf{P} \in \mathbb{R}^{2K}$. For this, we learn an objective function $\mathcal{F}^* : \mathbb{R}^{2K} \rightarrow \mathbb{R}^{3K}$ which minimizes the prediction error over a dataset with N poses:

$$\mathcal{F}^* = \min_f \frac{1}{N} \sum_{i=1}^N \mathcal{L}_{3d}(f(\mathbf{p}_i) - \mathbf{q}_i), \quad (3)$$

where \mathcal{L}_{3d} represents an MSE loss. Fig. 1-E shows the structure of our 3d pose estimation model “3d PoseNet” based on the architecture of [12]. It starts with a linear layer $f_{c_{2D}}$ which transforms the $2K$ -dimensional pose to 1024 dimensional features. Next, there are five linear layers $f_1 - f_5$, each with 4096 dimensions followed by Batch normalization, a Rectified Linear Unit and a dropout module. The final layer $f_{c_{3D}}$ produces $3K$ dimensional output. There are two residual connections defined in the network which combine information from lower layers to higher layers and improve model generalization performance.

3.4 The Proposed Fall Recognition (Fig. 1-G)

We present a neural network (*FallNet*) which consists of two sub-networks: a modality-specific network $F_\phi, \phi \in \{\mathbf{P}, \mathbf{Q}\}$, and an embedding network \mathbf{G} as shown in Fig. 1-G. The sub-network F_ϕ has a structure similar to [12] but with fewer linear layers. It produces 1024-dimensional features each from the two input modalities (\mathbf{P} and \mathbf{Q}). The output features are summed and fed into the embedding sub-network \mathbf{G} which uses two linear layers and learns probabilistic distributions with respect to the target classes. Let ρ_i denote the outputs of the last layer (f_{cls}) for the i^{th} input sample. The training objective function is defined over N poses as:

$$\mathcal{L}_{fall} = \sum_{i \in N} \mathcal{L}_{cls}(\rho_i, \rho_i^*), \quad (4)$$

where ρ_i^* represent the ground-truth labels. The term \mathcal{L}_{cls} is a Cross Entropy Loss, given by:

$$\mathcal{L}_{cls}(x, C) = - \sum_{C=1}^{N_C} \mathcal{Y}_{x,C} \log(p_{x,C}), \quad (5)$$

where \mathcal{Y} is a binary indicator if class label C is the correct classification for observation x , and p is the predicted probability of observation x of class C .

3.5 The Proposed Occluded Joints Resilience (OJR)

Pose estimators trained on RGB images inevitably make errors in joint predictions due to factors such as: image imperfections, occlusions, background clutter, and incorrect ground truth annotations. Since, our 3d PoseNet and FallNet models rely on the output of the SH network, errors in 2d pose predictions affect the quality of 3d pose estimation and fall recognition. To overcome this challenge, we present a method termed “Occluded Joints Resilience (OJR)” which increases the robustness of our models to incomplete information in the pose data. To achieve this, the OJR method creates an occlusion pattern \mathcal{M}_i and uses it to transform the original pose data into occluded pose data. The occlusion pattern \mathcal{M}_i is defined as:

$$\mathcal{M}_i = [v_1 J_1, \dots, v_k J_k], v \in \{0, 1\}, \quad (6)$$

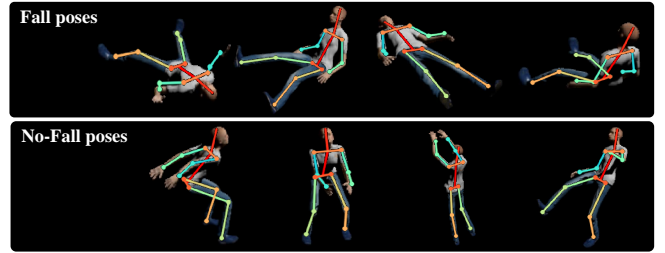


Figure 2: Sample frames from our Synthetic Human Fall dataset showing different poses.

where $J_i = (x_i, y_i)$ represents a body joint and v is a binary variable, indicating the visibility of the k th joint. During training, the OJR method generates a rich library of unique occlusion patterns $\{\mathcal{M}\}$ which vary across training samples, thereby increasing the network’s adaptivity to various occluded situations.

4 Experiments

4.1 Training and Implementation Details

We trained the SH network for 2d pose estimation using the MS COCO Keypoints dataset [11], which contains 64K images and 150K instances with 2d pose ground truth. To train our models for 3d pose estimation and fall recognition, we used the synthetic human fall dataset of [1], which provides 767K samples of human poses with 2d and 3d pose annotations categorized into fall and no-fall body poses. Fig. 2 shows some samples from the synthetic dataset. For training the 3d PoseNet and FallNet models, we initialized the weights of the fully connected layers with zero-mean Gaussian distributions (standard deviations were set to 0.01 and biases were set to 0), and trained each network for 300 epochs. The starting learning rate was set to 0.01 and divided by 10 at 50% and 75% of the total number of epochs. The parameter decay was set to 0.0005 on the weights and biases. The probability of dropout was set to 0.5. Our implementation is based on the framework of Torch library [17]. Training was performed using ADAM optimizer and four Nvidia Tesla K80 GPUs.

4.2 Test Datasets

To evaluate the generalization capability of our SSHFD for fall detection in unseen real-world environments, we trained our models using only synthetic data and tested the models on the public MultiCam fall dataset [2] and the Le2i fall dataset [3]. The MultiCam dataset consists of 24 different scenarios where each scenario is comprised of a video sequence of a person performing a number of activities (such as falling on a mattress, walking, carrying objects). Each scenario is recorded using 8 cameras from 8 different locations. The dataset is challenging for single-shot single camera fall detection because, different camera viewpoints produce occlusions and significant variations in the spatial locations, scale, and orientations of the falls [1]. The Le2i dataset contains 221 videos of different actors performing fall actions and various other normal activities in different environments. The dataset is challenging due to variable lighting conditions and occlusions [1]. To quantify the recognition performance of our SSHFD, we extracted image frames from the target videos at 25 fps resolution and generated 2d poses using the SH network. Next, we computed weighted F1 scores, precision (PRE) and recall

Table 1: Fall recognition results of the proposed SSHFD in terms of its different variants termed Human Fall Detection Models (HDF) on the MultiCam fall dataset and the Le2i fall dataset. The models for 3d pose estimation and fall recognition were trained only on the synthetic data and evaluated on real-world test datasets.

Human Fall Detection (HFD) Models	MultiCam fall dataset			Le2i fall detection database		
	F1Score	Precision	Recall	F1Score	Precision	Recall
SSHFD-A: SH + FallNet2d3d	0.8453	0.8487	0.8431	0.8991	0.9008	0.8992
SSHFD-B: SH + FallNet2d	0.8388	0.8437	0.8358	0.8885	0.8907	0.8887
SSHFD-C: SH + ResNet (RGB)	0.8638	0.8628	0.8658	0.6595	0.7985	0.6912

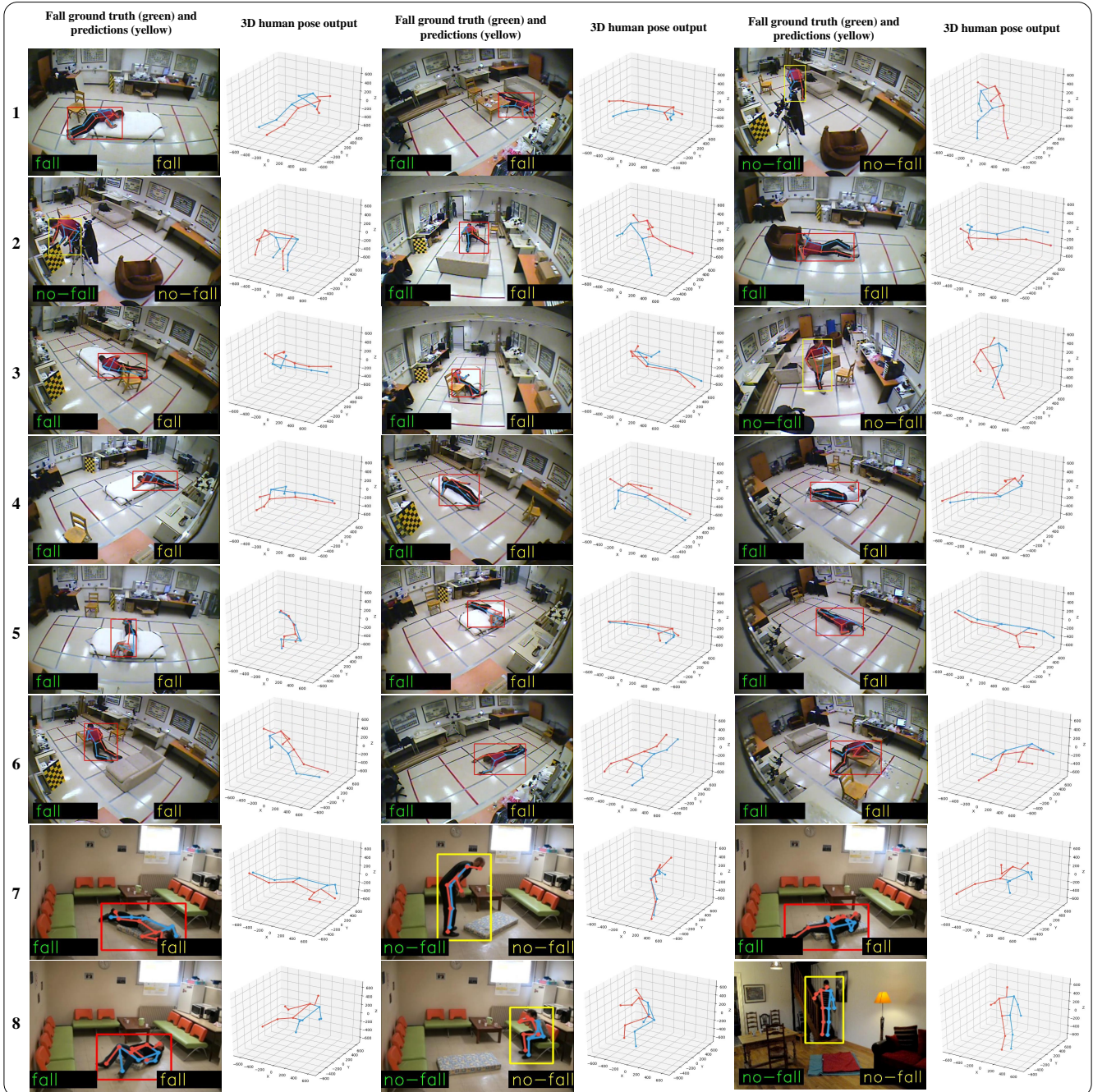


Figure 3: Qualitative results of our framework on the MultiCam fall dataset (rows 1-6) and the Le2i fall dataset (rows 7-8). Ground truth labels and model predictions are shown by text in green and yellow, respectively. Fall and no-fall cases are represented by bounding boxes in red and yellow, respectively.

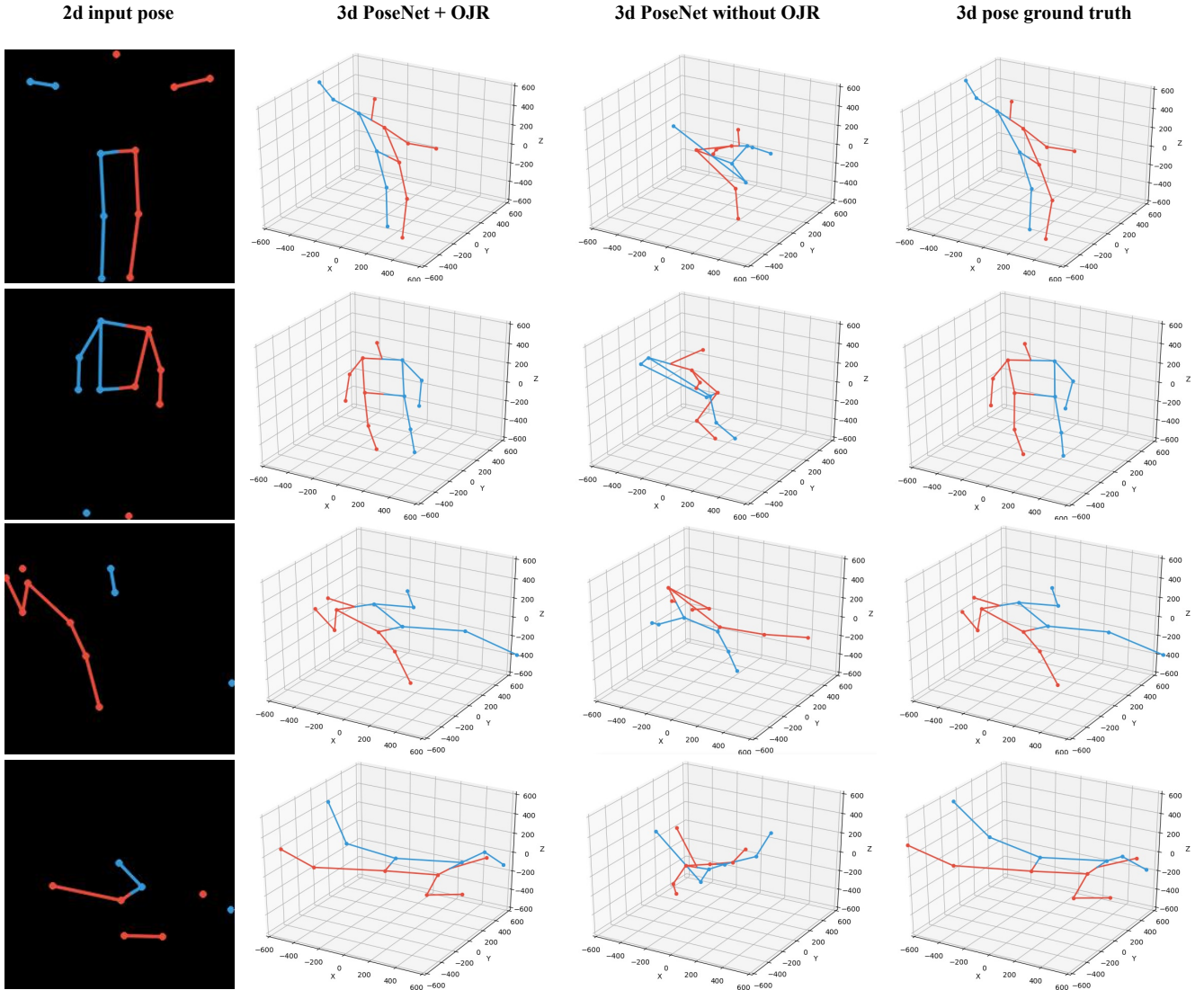


Figure 4: Qualitative comparison of the predictions of our 3d PosNet with and without the proposed OJR using inputs with missing joints on our synthetic dataset. Our OJR-based 3d PoseNet enables the model to successfully recover missing joints in the input pose data.

(RE) scores per image frame with at least one pose detected and averaged the scores over all image frames of the targets datasets. We used the weighted measures as they are not biased by imbalanced class distributions which make them suitable for the target datasets where the number of fall samples are considerably small compared to the number of non-fall samples.

5 Results

Table 1 shows fall recognition results on the test datasets, for different variants of our framework termed “Human Fall Detection Models”. The variants “A” and “B” use neural networks with linear structures which were trained on pose data as shown in Fig. 1 and described in Sec. 3.4. The variant “C” shown in Table 1 uses a ResNet18 [7] based CNN architecture which was trained on RGB appearance information of synthetic human proposals. The results reported in Table 1 show that although the RGB-based fall detector produced higher f1scores compared to the pose-based fall detectors

on the MultiCam dataset, it produced the lowest f1scores on the Le2i dataset. This is because, the RGB-based fall detector trained on color information of synthetic human proposals failed to generalize to the scenes of Le2i dataset with high variations in the appearance characteristics and different backgrounds. Compared to the RGB-based detector, our pose-based fall detector (SSHFD-A) produced competitive f1scores on the MultiCam dataset and superior f1scores on the Le2i dataset as shown in Table 1. Fig. 3 shows qualitative results of our pose-based fall detector on sample images from the test datasets. The results show that our fall detection framework is robust to partial occlusions, and variations in the spatial locations, scale, and orientations of fall poses in real-world scenes. These improvements are attributed to our pose-based fall representation which is invariant to appearance characteristics and makes our framework robust to different human actors and background clutter in real-world scenes. These results demonstrate the generalization capability of our framework in successfully transferring fall recognition knowledge learnt purely from synthetic data to unseen real-world data. Table 1 also shows that

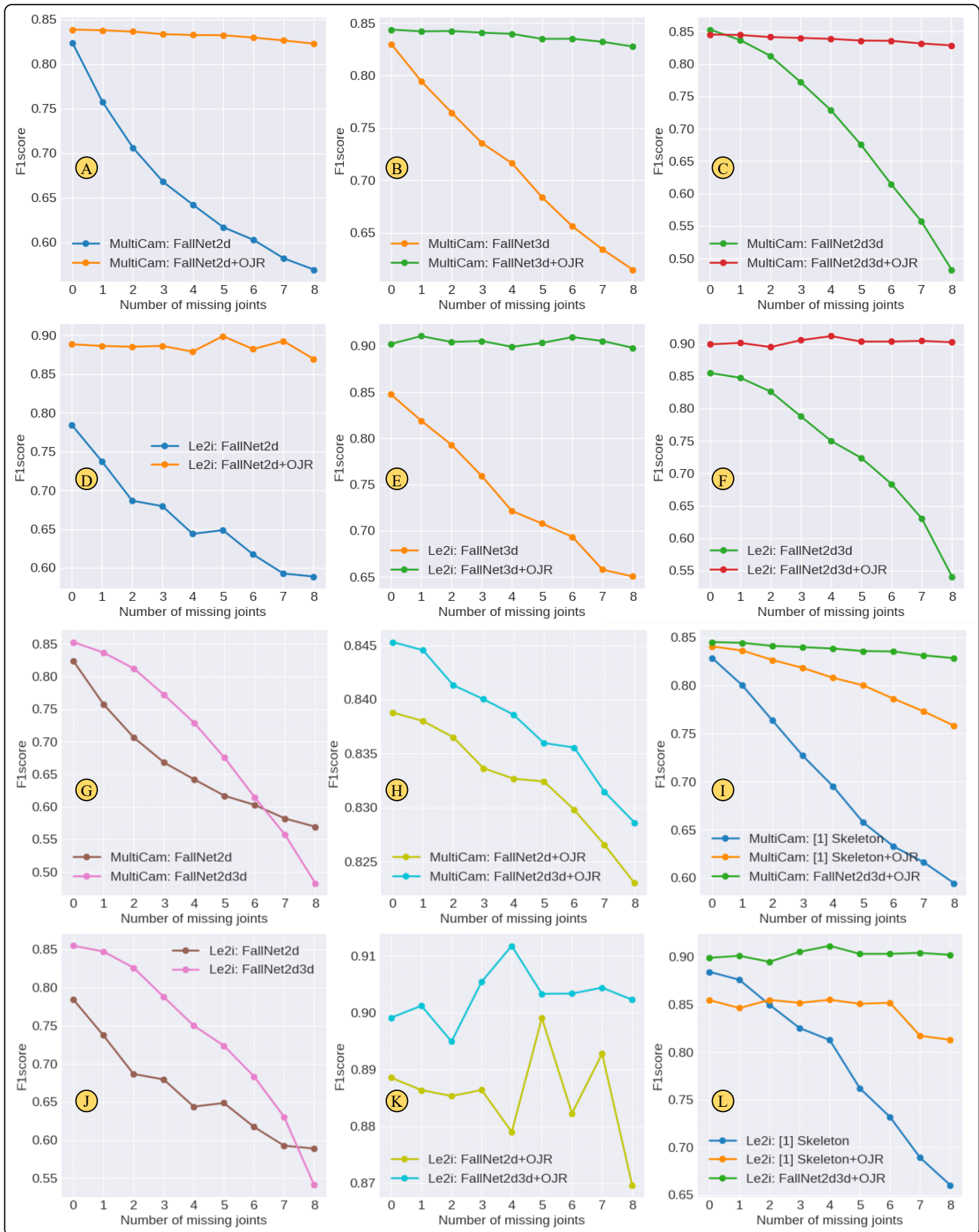


Figure 5: F1scores of our SSHFD on the MultiCam fall dataset and the Le2i fall detection database under different noise levels. The subplots A-F show that the proposed OJR-based models produce considerable higher f1scores for fall recognition under missing joints information compared to the models which were trained without the proposed OJR method. The subplots G, H, J, and K, show comparison between our 2d-pose based model “FallNet2d” and “FallNet2d3d” which uses both 2d and 3d pose for fall recognition. The subplots I and L show a comparison between f1scores of our method and the visual skeleton representation based method of [1] under different noise levels.

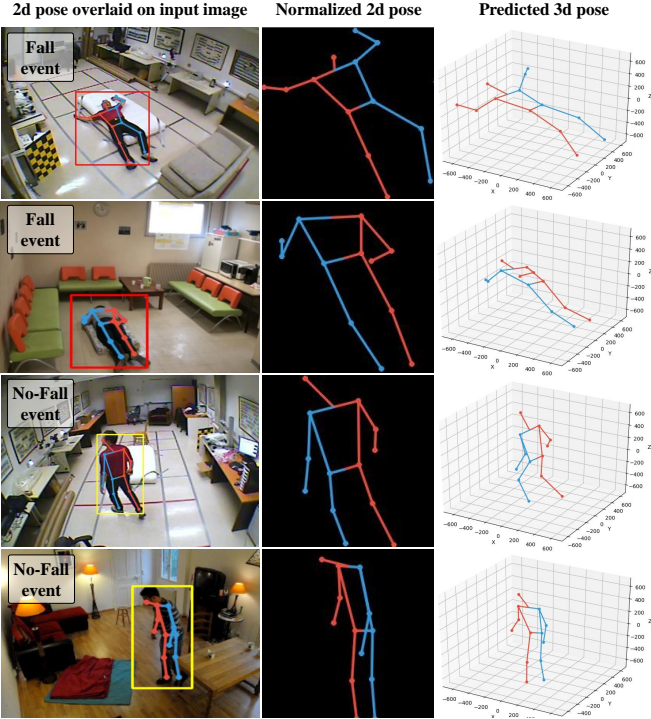


Figure 6: Variations in camera viewpoints cause ambiguities in 2d pose based fall representations (fall and no-fall 2d poses resemble each other as shown in the middle column). In contrast, our 3d PoseNet predictions (as shown in the right column) are more discriminative and reduce inter-class similarities for fall recognition. Fall and no-fall cases are represented by bounding boxes in red and yellow, respectively.

our FallNet2d3d model using combined 2d- and 3d-pose information performed better than the FallNet2d model which used only 2d pose information. This is attributed to the proposed FallNet architecture which uses low-level modality-specific layers to learn discriminative information from the individual pose modalities, and uses high-level fusion layers to learn the complimentary information in the multi-modal input pose data, thereby producing features which are robust to pose ambiguities in the 2d image space under different camera viewpoints as shown in Fig. 6.

5.1 Robustness to Missing Joints

5.1.1 Fall Recognition

Fig. 5 shows comparison of f1scores produced by our models with and without the proposed OJR on the MultiCam dataset and the Le2i fall dataset under different noise levels. The results show that our OJR-based models produced significantly higher f1scores for all the noise levels compared to the models which were trained without using OJR. For instance, using input pose data with 8 missing joints, our OJR-based models improved f1scores by upto 35% and 40% compared to the models without using OJR on the MultiCam and Le2i datasets, respectively (see Fig. 5-C and Fig. 5-F). We also conducted experiments to compare the performance of our FallNet with the method of [1] which uses visual representations of 2d skeletons and segmentation information for fall recognition. Fig. 5-I and Fig. 5-L show the results of these experiments. The results show that our 2d- and 3d-pose based fall representation produces superior fall recogni-

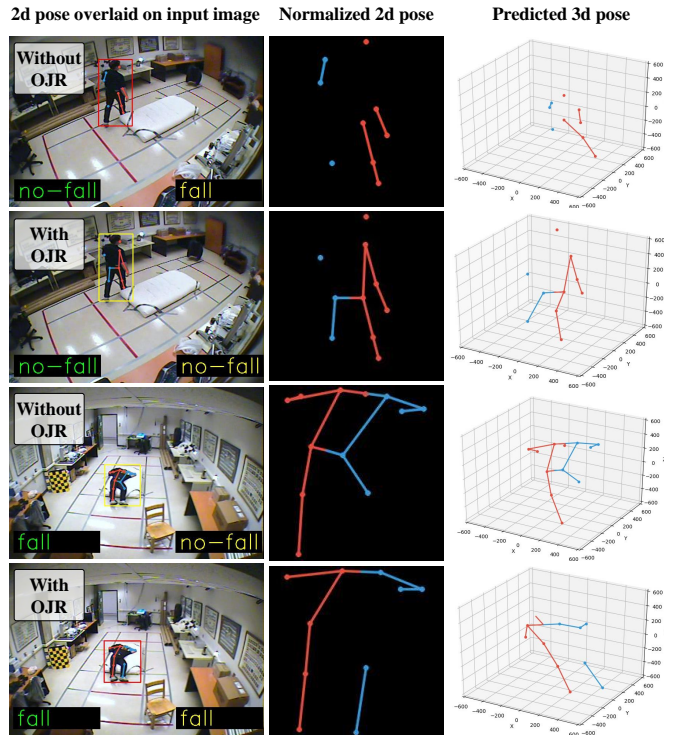


Figure 7: Our OJR-enabled FallNet model produces correct fall predictions in the presence of missing data in 2d pose and 3d pose compared to the model which was trained without the OJR method. Ground truth labels and model predictions are shown by text in green and yellow, respectively. Fall and no-fall cases are represented by bounding boxes in red and yellow, respectively.

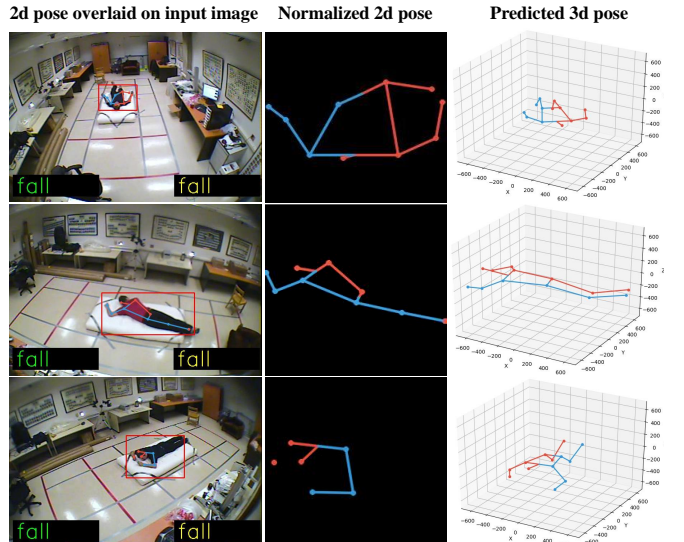


Figure 8: Our 3d PoseNet using the proposed OJR method successfully recovers missing data in the input 2d pose and enables the framework to produce correct fall predictions. Ground truth labels and model predictions are shown by text in green and yellow, respectively.

tion performance especially under missing joints data compared to the skeleton-based visual representation of [1]. Fig. 5-I and Fig. 5-L also show that the proposed OJR improved the performance of the method of [1] under different noise levels, demonstrating the signif-

icance of the proposed OJR for improving the robustness of models under scenarios with occluded joints. Fig. 7 shows qualitative results of our models on the MultiCam dataset using incomplete 2d- and 3d-pose data. The results show that the proposed OJR method makes our FallNet model robust to missing information in the pose data, and enables the model to make correct fall predictions under 2d or 3d pose errors.

5.1.2 3d Pose Estimation

Table 2: Comparison of the performance of our 3d PoseNet with and without the proposed Occluded Joints Resilience (OJR) on our synthetic dataset under different noise levels.

No. of missing joints	mean pose error (mm)	
	with OJR	without OJR
1	17.11	197.73
3	21.13	351.19
5	26.25	420.52
7	34.21	464.16

Here, we tested our 3d PoseNet under different levels of noise (missing joints) on the synthetic data. For this, we randomly split the data into 70% train and 30% test data subsets. Table 2 shows the mean joints position errors in millimeters which is the mean Euclidean distance between predicted joint positions and ground-truth joint positions averaged over all the joints, produced by our models on the test dataset. Table 2 shows that our OJR-based 3d PoseNet consistently produced lower pose errors compared to the model without OJR for all levels of noise on the test dataset. Fig. 4 shows qualitative results of our 3d PoseNet with and without the proposed OJR on our synthetic dataset. The results show that the proposed OJR enables our 3d PoseNet to successfully recover 3d joints information from incomplete 2d pose inputs. This enables our framework to make correct fall predictions under 2d pose errors as shown in Fig. 8.

6 Conclusion and Future Work

In this paper we present Single Shot Human Fall Detector (SSHFD), a deep learning framework for human fall detection from a single image. SSHFD learns fall representations based on human joint locations in 2d image space and 3d Cartesian space. Our fall representation is invariant to physical appearance, background, and enables our framework to successfully transfer fall recognition knowledge from pure synthetic data to unseen real-world data. We also present neural network models for 3d pose estimation and fall recognition which are resilient to occluded body parts. Experiments on real-world datasets demonstrate that our framework successfully handles challenging scenes with occlusions. These capabilities open new possibilities for advancing human pose based fall detection purely from synthetic data. In future, we plan to expand our framework for the recognition of other activities to enhance its potential for general human activity recognition.

REFERENCES

- [1] Umar Asif, Benjamin Mashford, Stefan von Cavallar, Shivanthan Yohanandan, Subhrajit Roy, Jianbin Tang, and Stefan Harrer, ‘Privacy preserving human fall detection using video data’, (2019).
- [2] Edouard Auvinet, Caroline Rougier, Jean Meunier, Alain St-Arnaud, and Jacqueline Rousseau, ‘Multiple cameras fall dataset’, *DIRO-Université de Montréal, Tech. Rep.*, **1350**, (2010).
- [3] Imen Charfi, Johel Miteran, Julien Dubois, Mohamed Atri, and Rached Tourki, ‘Optimized spatio-temporal descriptors for real-time fall detection: comparison of support vector machine and adaboost-based classification’, *Journal of Electronic Imaging*, **22**(4), 041106, (2013).
- [4] Jane Fleming and Carol Brayne, ‘Inability to get up after falling, subsequent time on floor, and summoning help: prospective cohort study in people over 90’, *Bmj*, **337**, a2227, (2008).
- [5] Samuele Gasparrini, Enea Cippitelli, Susanna Spinsante, and Ennio Gambi, ‘A depth-based fall detection system using a kinect sensor’, *Sensors*, **14**(2), 2756–2775, (2014).
- [6] Ross Girshick, Ilija Radosavovic, Georgia Gkioxari, Piotr Dollár, and Kaiming He. Detectron. <https://github.com/facebookresearch/detectron>, 2018.
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, ‘Deep residual learning for image recognition’, in *CVPR*, pp. 770–778, (2016).
- [8] Guang-Bin Huang, Qin-Yu Zhu, and Chee-Kheong Siew, ‘Extreme learning machine: a new learning scheme of feedforward neural networks’, in *Neural Networks, Proceedings. IEEE International Joint Conference on*, volume 2, pp. 985–990, (2004).
- [9] Dao Huu Hung, Hideo Saito, and Gee-Sern Hsu, ‘Detecting fall incidents of the elderly based on human-ground contact areas’, in *2nd IAPR Asian Conference on Pattern Recognition*, pp. 516–521. IEEE, (2013).
- [10] Chia-Feng Juang and Chia-Ming Chang, ‘Human body posture classification by a neural fuzzy network and home care system application’, *IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans*, **37**(6), 984–994, (2007).
- [11] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick, ‘Microsoft coco: Common objects in context’, in *ECCV*, pp. 740–755, (2014).
- [12] Julieta Martinez, Rayat Hossain, Javier Romero, and James J Little, ‘A simple yet effective baseline for 3d human pose estimation’, in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2640–2649, (2017).
- [13] Georgios Mastorakis and Dimitrios Makris, ‘Fall detection system using kinect’s infrared sensor’, *Journal of Real-Time Image Processing*, **9**(4), 635–646, (2014).
- [14] S-G Miaou, Pei-Hsu Sung, and Chia-Yuan Huang, ‘A customized human fall detection system using omni-camera images and personal information’, in *Distributed Diagnosis and Home Healthcare, 1st Transdisciplinary Conference on*, pp. 39–42. IEEE, (2006).
- [15] Behzad Mirmahboub, Shadrokh Samavi, Nader Karimi, and Shahram Shirani, ‘Automatic monocular system for human fall detection based on variations in silhouette area’, *IEEE Transactions on Biomedical Engineering*, **60**(2), 427–436, (2013).
- [16] Alejandro Newell, Kaiyu Yang, and Jia Deng, ‘Stacked hourglass networks for human pose estimation’, in *ECCV*, pp. 483–499, (2016).
- [17] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer, ‘Automatic differentiation in pytorch’, (2017).
- [18] Huimin Qian, Yaobin Mao, Wenbo Xiang, and Zhiqian Wang, ‘Home environment fall detection system based on a cascaded multi-svm classifier’, in *ICARCV*, pp. 1567–1572. IEEE, (2008).
- [19] Guoru Zhao, Zhanyong Mei, Ding Liang, Kamen Ivanov, Yanwei Guo, Yongfeng Wang, and Lei Wang, ‘Exploration and implementation of a pre-impact fall recognition method based on an inertial body sensor network’, *Sensors*, **12**(11), 15338–15355, (2012).