# In-Bed Pose Estimation: A Review

Ziya Ata Yazıcı
*Istanbul Technical University*
Istanbul, Turkey
yaziciz21@itu.edu.tr

Sara Colantonio
*ISTI-CNR*
Pisa, Italy
sara.colantonio@isti.cnr.it

Hazım Kemal Ekenel
*Istanbul Technical University*, Istanbul, Turkey
*Qatar University*, Doha, Qatar
ekenel@itu.edu.tr, hekenel@qu.edu.qa

*Abstract*—Human pose estimation, the process of identifying joint positions in a person's body from images or videos, represents a widely utilized technology across diverse fields, including healthcare. One such healthcare application involves in-bed pose estimation, where the body pose of an individual lying under a blanket is analyzed. This task, for instance, can be used to monitor a person's sleep behavior and detect symptoms early for potential disease diagnosis in homes and hospitals. Several studies have utilized unimodal and multimodal methods to estimate in-bed human poses. The unimodal studies generally employ RGB images, whereas the multimodal studies use modalities including RGB, long-wavelength infrared, pressure map, and depth map. Multimodal studies have the advantage of using modalities in addition to RGB that might capture information useful to cope with occlusions. Moreover, some multimodal studies exclude RGB and, this way, better suit privacy preservation. To expedite advancements in this domain, we conduct a review of existing datasets and approaches. Our objectives are to show the limitations of the previous studies, current challenges, and provide insights for future works on the in-bed human pose estimation field.

*Index Terms*—In-Bed Human Pose Estimation, Review

## I. INTRODUCTION

Human pose estimation involves predicting the joints, e.g., head, elbow, and knee, of a human body from an image or video. Since human pose serves as a key technology for human perception, it has been utilized in a wide range of applications [1]. The medical field has also benefited from human pose analysis by continuously monitoring individuals in home and hospital environments. Examples of practical applications include assistive systems for elderly individuals [2], estimating infant poses [3], automatic rehabilitation systems [4], and tracking the poses of surgeons and clinicians [5]. In addition to the given topics, an emerging research area, in-bed human pose estimation has also received attention recently. This topic involves tracking individuals' poses while lying or sleeping in a bed. These systems can provide intuitive insights into diagnosing disorders by observing the symptoms during sleep, evaluating sleep quality, and monitoring the sleeping position changes to improve post-surgery healing. However, occlusion brings a challenge to the task when the person is covered with a blanket. Hence, the available data modalities should be utilized to see under the occlusion.

To accelerate the studies on in-bed pose estimation and showcase the current limitations and potentials, our study will be the first review paper that covers the available datasets, used metrics, and the previous studies specific to the in-bed pose
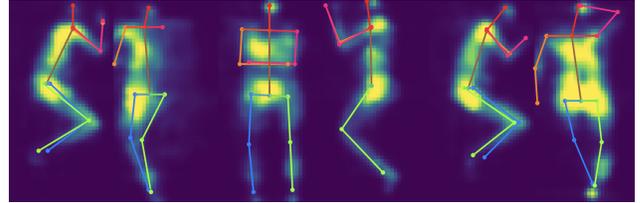


Fig. 1. Sample pressure maps and estimated poses from the Pressure-Sensing Mat Dataset [6].

estimation topic. First, we will review the available datasets, the number of samples they include, the modalities they are composed of, and the evaluation metrics used to evaluate the joint predictions. Second, we will categorize the in-bed pose estimation works into two groups: unimodal and multimodal methods. Unimodal studies use only a single modality to estimate human pose, while multimodal studies use multiple modalities to benefit from color, heat, pressure, and depth maps of the scene. Additionally, since the main application areas of such systems are the hospital rooms and intensive care units, one concern of the approaches is the preservation of the patient's privacy. Thus, besides RGB, the use of long-wavelength infrared (LWIR), depth maps, and pressure maps have been mainly used in the literature. Finally, the current limitations and directions of in-bed pose estimation will be given.

## II. DATASETS

In this section, we review the publicly available datasets collected for in-bed pose estimation tasks. These encompass various patient poses, lighting conditions, and occlusion scenarios. A summary of the datasets is also given in Table I.

*1) Pressure-Sensing Mat Dataset [6]:* This is a public unimodal pose estimation dataset with a pressure map of 17 human participants in a hospital bed. The dataset includes over 28,000 pressure maps with 3D human pose annotations in a 17-joint skeleton model. The dataset includes two configurations, e.g., supine (0°) and seated (60°), and the participants took various poses during the data collection. The goal of the dataset is to enable patient-assisting robots to estimate the 3D position of the patient while minimizing the impact of lighting and occlusion on the bed. Sample images and annotations from the dataset can be seen in Figure 1.

| Dataset | Modality | # Participants | # Samples | Labels | Type |
|---|---|---|---|---|---|
| Pressure-Sensing Mat Dataset [6] | Pressure Map | 17 | 28,000 images | 3D Pose | Image |
| Mannequin In-Bed Dataset [7] | RGB, LWIR | 2 | 419 images | 2D Pose | Image |
| BlanketSet [8] | RGB, LWIR, Depth Map | 14 | 303,965 frames | Action Class | Video |
| SLP Dataset [9] | RGB, LWIR, Depth Map, Pressure Map | 109 | 14,715 frames | 3D Pose | Video |
| Patient MoCap Dataset [10] | RGB, Depth Map | 10 | 180,000 frames | 3D Pose | Video |



Fig. 2. Samples images from the Mannequin In-Bed Dataset [7] in two modalities: LWIR and RGB modalities from left to right.

*2) Mannequin In-Bed Dataset [7]:* This is the first public dataset composed for in-bed pose estimation studies, collected by utilizing realistic male and female mannequins in different poses, e.g., supine, left-side lying, and right-side lying, in a hospital room setting. Images are taken with an infrared selective (IRS) image acquisition system, and two sets of normal and infrared-illuminated images are included in the dataset. The poses of the images were annotated in a 14-joint format, and a total of 419 non-occluded pose images with 2D ground-truth skeletons are available in the dataset. Sample images and annotations from the dataset can be seen in Figure 2.

*3) BlanketSet Dataset [8]:* This is a public multimodal video-based action recognition dataset with RGB, LWIR, and depth map modalities recorded from 14 participants in eight movement sequences: Foot-to-knee, knee bend, swinging legs, hands to shoulders, belly-down spread, torso lean, stretched arms and feet stretched. The videos were recorded with three different blanket positions and blanket types, e.g., thickness, color, and weight. In total, 405 videos were collected.
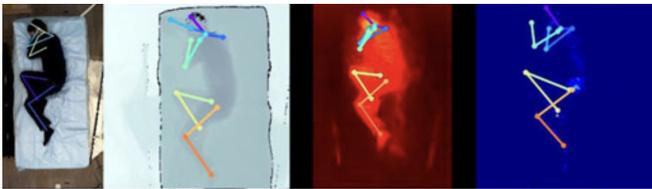


Fig. 3. Samples images from the SLP dataset [9] in four modalities: RGB, depth map, LWIR, and pressure map from left to right.

*4) Simultaneously-Collected Multimodal Lying Pose (SLP) Dataset [9]:* This is a large-scale public multimodal pose estimation dataset with 109 subjects and 14,715 samples in RGB, LWIR, depth, and pressure map modalities. Occlusion cases with different blanket thicknesses were also included. The dataset is divided into two parts: In-home

(102 participants) and hospital settings (7 participants). From the participants, 15 poses were collected in three different positions; supine, left, and right-side sleeping were collected, and annotated in 14 joints. Sample images and annotations from the dataset can be seen in Figure 3.
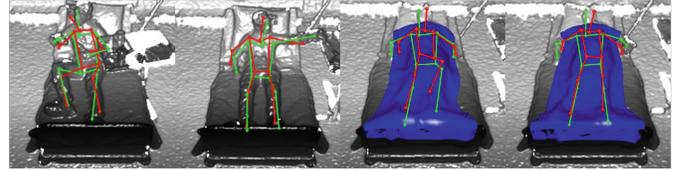


Fig. 4. Non-occluded and occluded sample depth maps from the Patient MoCap Dataset [10].

*5) Patient MoCap Dataset [10]:* This dataset comprises sequences in RGB and depth map, each with ground truth pose information obtained from five calibrated motion capture cameras. These cameras track 14 rigid targets on each subject, allowing for the inference of 14 body joint locations. The subjects, consisting of five females and five males, performed ten sequences involving various activities such as getting in/out of bed, sleeping on horizontal/elevated beds, eating with/without clutter, using objects, reading, clonic movements simulating epileptic seizures, and a calibration sequence. The dataset includes 180,000 video frames captured by a calibrated and synchronized RGB-D sensor. Sample images and annotations from the dataset can be seen in Figure 4.

## III. METHODS

This section includes the unimodal and multimodal approaches used for the in-bed pose estimation task. Each study will be investigated by the utilized modalities, selected models, and the proposed approach. A summary of the approaches is given in Table II.

### A. Unimodal In-bed Pose Estimation

Several studies utilize a single source of information for in-bed pose estimation. In [11], as one of the earliest approaches to in-bed pose estimation topic, the authors introduce a vision-based tracking system designed for long-term monitoring of in-bed postures in various environments. The system uses a latent variable-based hierarchical inference model to generate in-bed posture tracking history reports from top-view videos captured by regular off-the-shelf cameras. Despite being supervised, the model is person-independent and can be applied to new users without retraining. Experiments are performed

TABLE II
A SUMMARY OF THE PROPOSED APPROACHES ON IN-BED POSE ESTIMATION

| Studies | Highlight | Dataset | Metric | Modality |
|---|---|---|---|---|
| Liu and Ostadabbas [11] | A person-independent, latent variable-based hierarchical inference model was designed. | Mannequin [7] | Accuracy | RGB |
| Liu *et al.* [7] | A CNN-based pose estimation with HOG rectification in non-occluded scenarios was performed. | Mannequin [7] | PCK | LWIR |
| Davoodnia *et al.* [12] | An autoencoder model was used to change the pressure map representation for pose estimation. | Pressure [6] | MPJPE | Pressure |
| Bigalke *et al.* [13] | Anatomical loss-based unsupervised and source-free domain adaptation methods were implemented. | SLP [9] | MPJPE | Depth |
| Afham *et al.* [14] | Two GAN-based augmentation methods were used to increase the labeled images with occlusion. | SLP [9] | PCK | LWIR |
| Casas *et al.* [15] | Hash and CNN-based two approaches were implemented for pose estimation on pressure maps. | *Private* | MAE | Pressure |
| Obeidavi *et al.* [16] | A multi-scale CNN model was proposed for thermal image-based pose estimation. | SLP [9] | PCK | LWIR |
| Cao *et al.* [17] | A variational autoencoder is designed to retrieve the missing RGB modality during inference time. | SLP [9] | PCK | RGB, LWIR |
| Dayarathna *et al.* [18] | Different modality fusion approaches and best performing modality couple were investigated. | SLP [9] | PCK | RGB, LWIR, Depth, Pressure |
| Yin *et al.* [19] | A modality fusion pyramid with an attention-based reconstruction method was designed. | SLP [9] | MPJPE | RGB, LWIR, Depth, Pressure |

on the mannequin dataset and a private human dataset with 358 samples from 12 participants. The results demonstrate the posture detection accuracy as 91.0% on the mannequins dataset [7] and 93.6% on the private dataset with human participants.

In [7], the authors propose a pre-trained Convolutional Neural Network (CNN) model, named Convolutional Pose Machine (CPM) to estimate in-bed poses on their composed Mannequin In-Bed Dataset. While collecting the images, the challenges unique to in-bed poses, such as lighting variations and unconventional perspectives, are addressed through an infrared selective (IRS) image acquisition technique. The collected images are processed with a Histogram of Oriented Gradient (HOG) rectification method to handle unusual positions and angles of the subject. They fine-tune the intermediate layers of the model stages with limited data and achieve an 82.6% probability of correct keypoint (PCK) with a threshold of 20%, a metric for joint correctness based on the distance between the predicted and ground truth joints compared to a percentage of a person's torso length.

In [12], the authors explore the application of in-bed pose estimation using pressure data. The study evaluates various strategies for detecting body pose from ambiguous pressure data by leveraging pre-existing pose estimators. The approaches include using pose estimators pre-trained on the RGB domain, re-training them on pressure datasets, and employing a learnable pre-processing domain adaptation step to transform pressure maps into a new representation. The authors introduce a domain adaptation method with a fully convolutional network, PolishNetU, to generate robust representations for common pre-trained pose estimation models. PolishNetU incorporates an objective function addressing pose identification loss, reconstructing the lost body parts, and penalizing large deviations from the original pressure maps. On the Pressure-Sensing Mat Dataset [6], the authors re-train pre-existing image-based pose estimators using the new representations of the pressure maps, which significantly increases the pose estimation accuracy compared to using the pressure maps directly.

In [13], the authors propose a novel domain adaptation method that adapts a labeled source model to an unlabeled target domain, focusing on in-bed pose estimation. Two complementary adaptation strategies based on prior knowledge about human anatomy are introduced. The first strategy –

Unsupervised Domain Adaptation (UDA) – guides the learning process in a supervised manner for the source domain and an unsupervised manner for the target domain, achieved through embedding anatomical constraints into an anatomical loss function. The second strategy – Source-Free Domain Adaptation (SFDA) – involves filtering pseudo labels for self-training based on their anatomical plausibility. The evaluation is conducted under two adaptation scenarios using the depth maps from the SLP [9] dataset and a newly created dataset from 13 participants. The proposed method reduces the domain gap and surpasses the baseline model by a 95.5 mean per joint position error (MPJPE), calculated as the average L2 Norm of differences between predicted and ground truth joint coordinates.

In [14], the researchers introduce a methodology for cross-domain in-bed pose estimation, particularly in the LWIR modality. The formulation of the problem involves standard 2D human pose estimation tasks with labeled source domain data and unlabeled target domain data. To mitigate the domain gap, based on the work of Zhu *et al.* [20], the authors propose a two-fold data augmentation pipeline, incorporating CycAug for unpaired image-to-image translation and ExtremeAug for introducing more covering artifacts and occlusions. The augmented images are then utilized for pose estimation, extending the input space for optimization. Furthermore, knowledge distillation is employed to transfer knowledge from a teacher to a student model, enhancing student performance. On the SLP dataset, by training only with the uncovered images and performing domain transfer into covered images, the approach achieves a 76.13% PCK0.5 score.

In [15], the objective is to explore in-bed motion monitoring by utilizing pressure sensors as a privacy-preserving alternative to video-based methods. Two approaches are presented: a hashing content-retrieval approach and a deep learning method for human pose estimation from pressure sensor data. The hashing-based approach is used to retrieve the closest body poses by searching the nearest neighbors, and the ground truth poses are averaged to predict the body pose in 3D. On the other hand, the deep learning approach uses a CNN model to estimate the human pose from the pressure maps. The experiments on a private dataset show a 12.20 mean absolute error (MAE) between the ground truth and the predicted joints for the hashing approach, while the CNN model achieves 8.00 MAE.

In [16], the study focuses on the use of a thermal camera image-based pose estimation model, both sustaining the privacy of the patient and robustness in occluded cases. The proposed Multi-Scale Stacked Hourglass (MSSHg) network is applied to enhance the processing of the images in different scales for pose estimation. The results demonstrate the effectiveness of the MSSHg network, achieving an accuracy of 96.8% in the PCK0.2 on the SLP dataset [9].

### B. Multimodal In-bed Pose Estimation

Multimodal in-bed pose estimation studies employ multiple data sources, such as RGB, LWIR, pressure maps, and depth maps. In [17], the authors design a multimodal pipeline to extract more informative features for training. However, they remove the non-privacy preserving RGB modality during inference time. To maintain the performance of the model, the authors propose a solution using a multimodal conditional variational autoencoder (MCVAE) in conjunction with the HRNet [21] model. This approach enables the reconstruction of features from missing modalities during test time, providing a self-supervised mechanism for learning. The results show that the proposed framework effectively predicts the joints from the available modality by achieving a 95.8 PCK0.5 score, compared to the method fusing all the modalities which achieved 95.6 PCK0.5 on the SLP dataset [9].

In [18], the study utilizes multiple modalities such as RGB, depth maps, LWIR, and pressure maps. There are two primary objectives of this study: (1) The effective fusion of information from different modalities to enhance pose estimation and (2) the development of a framework capable of estimating in-bed poses with only the privacy-preserving modalities. Various fusion techniques are explored, including addition, concatenation, fusion via learnable parameters, and an end-to-end fully trainable approach coupled with a state-of-the-art pose estimation model. The study also introduces visible image reconstruction from privacy-preserving LWIR images using a conditional generative adversarial network (GAN). Experiments on the SLP dataset [9] demonstrate the performance increase of the proposed fusion model compared to the recent unimodal approaches in the literature.

In [19], the authors propose a pyramid scheme to effectively fuse different modalities, leveraging the knowledge captured by multimodal sensors. The depth and infrared images are initially fused to generate robust pose and shape estimations. Subsequently, pressure maps and RGB images are fused to refine the results, providing occlusion-invariant information for covered parts and accurate shape details for uncovered parts. An attention-based reconstruction module is proposed to generate uncovered modalities to mitigate occlusion effects. The pyramid fusion scheme and attention-based reconstruction module prove effective in complementing each modality, showcasing the contributions of depth, infrared, pressure map, and RGB data in the pose estimation process. According to the results, the fusion of modalities in multiple levels results in an MPJPE of 80.21 for thick-covered samples on the SLP dataset [9].

## IV. LIMITATIONS AND POTENTIALS

The first limitation of the field can be given as the comparability of the current studies. Since different dataset portions are used during model evaluations in the papers, **a fair comparison between the approaches is not possible**. Furthermore, there is a lack of diverse datasets in terms of age and inclusion of people with disorders. In scenarios with a person having anatomical abnormalities or people from different ages, the proposed algorithms may not perform well due to limb sizes and joint locations from out-of-distribution. Therefore, the future **datasets should be sampled from a larger and diverse population**.

Regarding the algorithm design, a potential that may improve the current state-of-the-art is to **include more modalities during model training**. The previous approaches generally utilize one or two modalities to estimate human poses, and only a few studies are exploring the use of more than two modalities. Besides, RGB images were mainly eliminated in the studies due to the visual attributes they include, which might breach individuals' privacy. However, the color and shape information obtained from the RGB images can be beneficial for the task. Moreover, some of the current approaches utilize RGB images during training time but eliminate them during inference time. To maintain the performance of the model, GAN models are used to generate the missing modalities during test time. However, these approaches can slow down the inference time by increasing the number of floating operations, and may not be suitable for real-time scenarios. Therefore, in the future, generating a new latent representation for the RGB images with lightweight models to benefit from the modality while preserving privacy, and the fusion of more than two modalities would be more appropriate in terms of speed, privacy, and accuracy.

## V. CONCLUSION

In this review, we summarized the publicly available datasets and the in-bed human pose estimation approaches in the literature. It was shown that while the uncovered poses are relatively easier to estimate than the covered scenes, the fusion of different modalities both preserves the privacy of the patient and gives a comprehensive understanding of the person under the blanket, thus improving the accuracy. To eliminate the visual attributes of the patients, non-RGB modalities have been focused on more during the model development due to privacy concerns. It was observed that a common benchmark is necessary to evaluate proposed models since the current evaluation of the approaches is dependent on the selection of the training and test subsets. Finally, to evaluate the generalizability of the models, individuals with different anatomical characteristics and from a wider age range should be included in the datasets for future studies.

REFERENCES

[1] J. Wang, S. Tan, X. Zhen, S. Xu, F. Zheng, Z. He, and L. Shao, "Deep 3d human pose estimation: A review," *Computer Vision and Image Understanding*, vol. 210, p. 103225, 2021.

[2] J. Jang, D. Kim, C. Park, M. Jang, J. Lee, and J. Kim, "Etri-activity3d: A large-scale rgb-d dataset for robots to recognize daily activities of the elderly," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 10990–10997, IEEE, 2020.

[3] D. Groos, L. Adde, R. Støen, H. Ramampiaro, and E. A. Ihlen, "Towards human-level performance on automatic pose estimation of infant spontaneous movements," *Computerized Medical Imaging and Graphics*, vol. 95, p. 102012, 2022.

[4] Y. Qiu, J. Wang, Z. Jin, H. Chen, M. Zhang, and L. Guo, "Pose-guided matching based on deep learning for assessing quality of action on rehabilitation training," *Biomedical Signal Processing and Control*, vol. 72, p. 103323, 2022.

[5] V. Srivastav, A. Gangi, and N. Padoy, "Unsupervised domain adaptation for clinician pose estimation and instance segmentation in the operating room," *Medical Image Analysis*, vol. 80, p. 102525, 2022.

[6] H. M. Clever, A. Kapusta, D. Park, Z. Erickson, Y. Chitalia, and C. C. Kemp, "3D Human Pose Estimation on a Configurable Bed from a Pressure Image," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 54–61, IEEE, 2018.

[7] S. Liu, Y. Yin, and S. Ostadabbas, "In-Bed Pose Estimation: Deep Learning With Shallow Dataset," *IEEE Journal of Translational Engineering in Health and Medicine*, vol. 7, pp. 1–12, 2019.

[8] J. Carmona, T. Karácsony, and J. P. S. Cunha, "BlanketSet–A Clinical Real Word Action Recognition and Qualitative Semi-synchronised MoCap Dataset," *arXiv preprint arXiv:2210.03600*, 2022.

[9] S. Liu, X. Huang, N. Fu, C. Li, Z. Su, and S. Ostadabbas, "Simultaneously-collected multimodal lying pose dataset: Enabling in-bed human pose monitoring," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 1, pp. 1106–1118, 2022.

[10] F. Achilles, A.-E. Ichim, H. Coskun, F. Tombari, S. Noachtar, and N. Navab, " Patient MoCap: Human Pose Estimation Under Blanket Occlusion for Hospital Monitoring Applications," in *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2016: 19th International Conference, Athens, Greece, October 17-21, 2016, Proceedings, Part I 19*, pp. 491–499, Springer, 2016.

[11] S. Liu and S. Ostadabbas, "A Vision-Based System for In-Bed Posture Tracking," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pp. 1373–1382, 2017.

[12] V. Davoodnia, S. Ghorbani, and A. Etemad, "Estimating pose from pressure data for smart beds with deep image-based pose estimators," *Applied Intelligence*, vol. 52, no. 2, pp. 2119–2133, 2022.

[13] A. Bigalke, L. Hansen, J. Diesel, C. Hennigs, P. Rostalski, and M. P. Heinrich, "Anatomy-guided domain adaptation for 3d in-bed human pose estimation," *Medical Image Analysis*, vol. 89, p. 102887, 2023.

[14] M. Afham, U. Haputhanthri, J. Pradeepkumar, M. Anandakumar, A. De Silva, and C. U. Edussooriya, "Towards Accurate Cross-Domain in-Bed Human Pose Estimation," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 2664–2668, IEEE, 2022.

[15] L. Casas, N. Navab, and S. Demirci, "Patient 3d body pose estimation from pressure imaging," *International Journal of Computer Assisted Radiology and Surgery*, vol. 14, pp. 517–524, 2019.

[16] S. Obeidavi, M. Gandomkar, and G. Hirtz, "In-Pose Estimation of Covered and Uncovered Human Body from Thermal Camera Images Using Multi-Scale Stacked Hourglass (MSSHg) Network," in *2022 16th International Conference on Signal-Image Technology & Internet-Based Systems (SITIS)*, pp. 84–90, IEEE, 2022.

[17] T. Cao, M. A. Armin, S. Denman, L. Petersson, and D. Ahmedt-Aristizabal, "In-Bed Human Pose Estimation from Unseen and Privacy-Preserving Image Domains," in *2022 IEEE 19th International Symposium on Biomedical Imaging (ISBI)*, pp. 1–5, IEEE, 2022.

[18] T. Dayarathna, T. Muthukumarana, Y. Rathnayaka, S. Denman, C. de Silva, A. Pemasiri, and D. Ahmedt-Aristizabal, "Privacy-preserving in-bed pose monitoring: A fusion and reconstruction study," *Expert Systems with Applications*, vol. 213, p. 119139, 2023.

[19] Y. Yin, J. P. Robinson, and Y. Fu, "Multimodal in-bed pose and shape estimation under the blankets," in *Proceedings of the 30th ACM International Conference on Multimedia*, pp. 2411–2419, 2022.

[20] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Towards Accurate Cross-Domain in-Bed Human Pose Estimation," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2223–2232, 2017.

[21] K. Sun, B. Xiao, D. Liu, and J. Wang, "Deep High-Resolution Representation Learning for Human Pose Estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5693–5703, 2019.