

Estimating Human pose from depth images using Convolutional Neural Networks

*When all you got is a powerful GPU, everything
looks like a CNN problem*

Bård-Kristian Krohg



Thesis submitted for the degree of
Master in Informatics: Robotics and Intelligent Systems
60 credits

Institute for informatics
Faculty of mathematics and natural sciences

UNIVERSITY OF OSLO

Spring 2020

Estimating Human pose from depth images using Convolutional Neural Networks

*When all you got is a powerful GPU, everything
looks like a CNN problem*

Bård-Kristian Krohg

© 2020 Bård-Kristian Krohg

Estimating Human pose from depth images using Convolutional Neural Networks

<http://www.duo.uio.no/>

Printed: X-press printing house

Estimating Human pose from depth images using Convolutional Neural Networks

Bård-Kristian Krohg

2nd March 2020

Abstract

This work is part of a larger project where we explore bringing robotics into geriatric care. The goal of this project is to create a robotic system that can assist in optimizing the use of physical personell, so they are used where they are needed.

This work will focus on capturing information about the user, anonymization of the data, what data is neccessary or ethical to capture, limitations for on-location data processing and what data can be sent for further processing in the cloud, or for human analysis.

We will also implement an ethical data-collection suite for the open-source Robotic Operating System, which can be implemented on a wide variety of robots.

Convolutional Neural Networks has been used for solving object recognition in 2D images with great success. This work aims to use the same techniques to extract 3D human pose from depth images in real time. We will use two multi-staged CNNs, one to encode the location of each joint, and another to encode the association between the joints to do this.

Preface

Contents

1	Introduction	1
1.1	Multimodal Elderly Care System	1
1.2	Hardware	2
1.3	Privacy	3
2	Background	4
2.1	Convolutional Neural Networks	4
2.2	Pose Estimation	4
3	Human 3D pose from depth images	6
3.1	Architecture	6
3.1.1	Depth feature extraction	7
3.1.2	3D object detection	8
3.1.3	Articulation network	8
3.2	Training data preparation	8
4	Experiments	10
5	Future Work	11
6	Conclusions	12
Appendices		13
.1	Depth sensors	14
.1.1	Stereo vision	14
.1.2	Structured light	14
.1.3	Time-of-Flight	15
.2	Robotic Operating System	15
.3	Classifiers	15
Glossary		18

List of Figures

2.1	OpenPose pipeline	5
3.1	Main architecture	7
3.2	Numbering for keypoint markers	9

List of Tables

3.1 Names/coordinates for detected landmarks	9
--	---

Chapter 1

Introduction

Finding human pose can have many applications inside the field of Human Robot Interaction (HRI), to understand the intent of a person, if they are aware of the robot, their position in relation to the robot, help understand their reaction to an interaction via body-language, or the action the person is doing at any one time. Finding human pose can also be used to create cheap motion capture technology for the ever popular Virtual Reality (VR) applications, which ranges from games to virtual training or simulations.

Human pose can also help professionals understand a patient better. If tracked over time, one can observe changes in the persons gait, their posture, or track their daily activities. All this can be used to assess if treatments or prescriptions are working, or if a persons activity schedule needs to be adjusted.

The method described in this work seeks to solve the human pose estimation problem for the Multimodal Elderly Care System (MECS) project. This means finding accurate human pose in 3D, based only on a single depth image from *one* viewpoint as input.

An algorithm using sequential refinement-steps with a combination of bottom-up evidence gathering, and a top-down neural network for articulation estimation is explored. The work also explores how shallow the network can be, as well as other optimization techniques for running on mid- to low-end hardware.

1.1 Multimodal Elderly Care System

As life expectancy increases in Norway, so does the population who needs geriatric care either at home, or in a geriatric facility. According to [4], it is projected that especially the elderly population in Oslo will increase in the coming years, and that this will lead to increased pressure on the healthcare services. This means that it will be even more important to encourage health-promoting, and take preventative

measures.

As part of the effort to let people live independently at home for as long as possible, the MECS was proposed. One of the goals for the MECS is to function as an autonomous safety alarm, a device that lets elderly living at home call the emergency services at the click of a button. However, if the emergency is an accident which renders the user incapacitated, or otherwise unconscious, an alarm that requires interaction will not be of much help. In contrast, the MECS can monitor a user, and warn healthcare personnel in case of an occurring, or predicted, emergency. The MECS will even be able to send contextual information to first-responders, which lets them better prepare for the situation.

The system is specified to be non-invasive, as this will increase the convenience for the user, in that the MECS will not require any interaction from the user to function. Taking the user out of the operational loop has other advantages as well. For example, a monitoring system in the form of a smartwatch, will be inconvenient and ineffective the user forgets to put it on, whereas a passive system will not need the user to perform any action to be effective.

Further, it is hypothesized that the information gathered by the MECS can help doctors or physical therapists prescribe or recommend health-promoting activities for the user. This could help prevent accidents or lifestyle diseases – which again will help relieve pressure on the communal healthcare services.

1.2 Hardware

The MECS is envisioned to be a small, mobile unit as this can be introduced to any home without extensive alterations to the environment, thus lowering the cost of the system and reusability of the units. As the MECS would need a charging station anyway, we propose a master/slave configuration between a stationary and a mobile unit, which should communicate through a secure wireless connection, for example a WLAN. We let the stationary unit take care of the power consuming complex processing, extending the operational time for the mobile unit between battery charges. We will therefore assume the system has access to mid- to high-end personal GPU/processing hardware, when we evaluate the real world practicality/runtime of the algorithm.

To provide as good a service for the user as possible, we believe that gathering many channels of information will be helpful. We wish to learn the users daily activity patterns or vital signs so when unhealthy or risk-filled patterns emerge, preventative actions can be implemented. Human Activity Recognition (HAR), gate/mood recognition, or detection of vital signs all require us to know where the user is in the scene.

This also places some requirements on our system. If the system solely relies on

this work to find humans in the scene, we set our lower framerate limit to 8 fps to be able to preform human heart rate aquisition [8]¹. Further, the MECS needs to be able to recognize humans in unstructured environments, in a variety of poses and to diffrentiate between multiple people.

In order to log the users activity patterns, and detect anomalies or deviations in this, we envision the MECS doing HAR.

1.3 Privacy

The information gathered by this system should only be available to the users designated doctors/physisians and to the user themselves, and should be treated at the same classification level as any persons medical journal. This means that the dataprocessing from the MECS must happen on-site, or that any cloud processing happens under a user agreement that protects this data from those owning the servers. The same agreement should count for anyone making the hardware/software that is used in the MECS.

The MECS is also capable of gathering additional data that is not relevant for the healthcare personell. For example, by using depth sensors we are able to create 3d maps of the users home or environment. this is helpful for the MECS, however it is not relevant for the healthcare personel. (with the exception of first responders, which could get the information either through a descriptive message – “the patient is in the bathroom on the second floor, no elevator, steep stairs” or via the actual internal map the MECS has created for its own internal navigation.)

¹If the maximum human heart rate is 220 bpm, and we want to measure it accuratley using video sequences produced by the MECS, our sampling frequency needs to be higher than 7. $\bar{3}$ Hz in order to satisfy the Nyquist rate.

Chapter 2

Background

2.1 Convolutional Neural Networks

CNNs are widely used in image classification tasks, because they look at a collection of spatially connected pixels, and is therefore robust to objects being framed differently.

2.2 Pose Estimation

In this work we are heavily inspired by [2] that uses two convolutional neural networks to find human pose. One of the networks finds the probability for the 2D location of a set of joints, we get N confidence maps for the locations of each N number of joints. The other creates M Part Affinity Field (PAF) the probability maps for M number of limbs¹. We then use the results from this PAF to find out which joints from the first result should be connected.

Using the method described in [2] we however don't get any information if some pairs of joints are missing, for example due to occlusion or failure to detect the body landmark, leading to an incomplete skeleton.

A lot of research² has been going into extracting human pose either from RGB images, or using depth images. Although many methods exist such as *Histogram of Gradients (HoG)* classifiers,

A lot of research has been done in estimating human pose in two dimensions, as this is where we have quite large datasets, such as the MPII, or the Human 3.6M datasets [1, 3].

¹In this work we'll stick to the convention of using the term limb to describe any connection between any pair of body landmarks, which we call joints.

²TODO: Cite Research

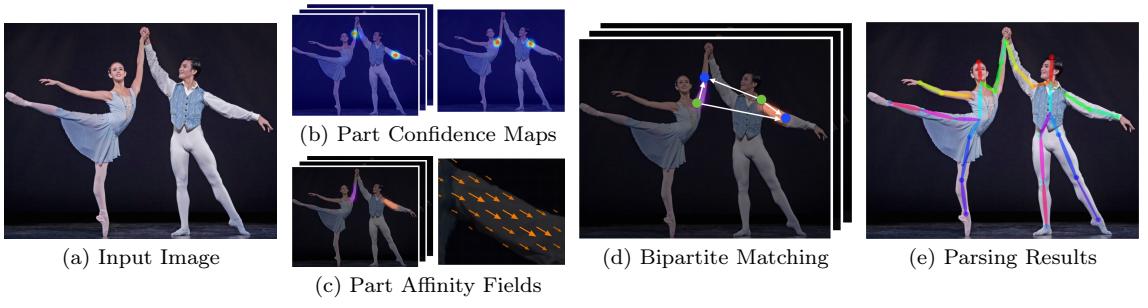


Figure 2.1: The pipeline described in [2]. The input image 2.1a is fed into the two networks which produce joint detections in confidence maps 2.1b and PAFs 2.1c. They then perform bipartite matching in 2.1d to determine which detected joints should be connected by a limb. 2.1e shows the finished results.

Chapter 3

Human 3D pose from depth images

To train any network using supervised learning, we need large amounts of training data. One of the goals for the MECS project is to do *Human Activity Recognition (HAR)*, so one can track the user from day to day and look for patterns that could lead to worsening living conditions. We also want to be able to recognize the activity from any viewpoint, and this is where a 2D approach will lack robustness. This is because any HAR model trained solely on 2D data, will only be able to recognize the activity from the views it has seen the activity being preformed. A 3D approach will provide us with robustness in respect to view-independentness.

We implemented an algorithm for extraction of human pose in 3D. Applying methods used on 3-channel (RGB) images to depth images, we show that the same methods can be used to extract objects in 2d images, can be used to extract objects in depth images as well, when it comes to human pose.

As in [2] we will use two networks to create the PAFs and the confidence maps for the joints. However, instead of training on 3-channel RGB images, we will use a single channel depth image to discover the body landmarks/joints. However, since depth images are single channel, and thus have less information than the RGB images, we propose using a shallower network. This also means we have to do the first step of feature extraction which was already done in a However since the depth images are less detailed than normal RGB images some landmarks might be harder to detect, such as eyes or nose or placing the joint on outstretched limbs.

This was considered when preparing the training data.

3.1 Architecture

When we are creating a neural network, it is often helpful to have in mind *what* we want to detect in each layer. It is therefore segregated into a couple of different steps to make it easier to follow along.

The architecture of this project is *recurrent* in that it repeats itself for a number of iterations. As with the [2] architecture, we have to have a first step which produces the first outputs we can use in later steps. However, this step is not illustrated in Fig.3.1, since the first step will be identical to the next steps except it won't have the additional inputs produced by the outputs of the previous step.

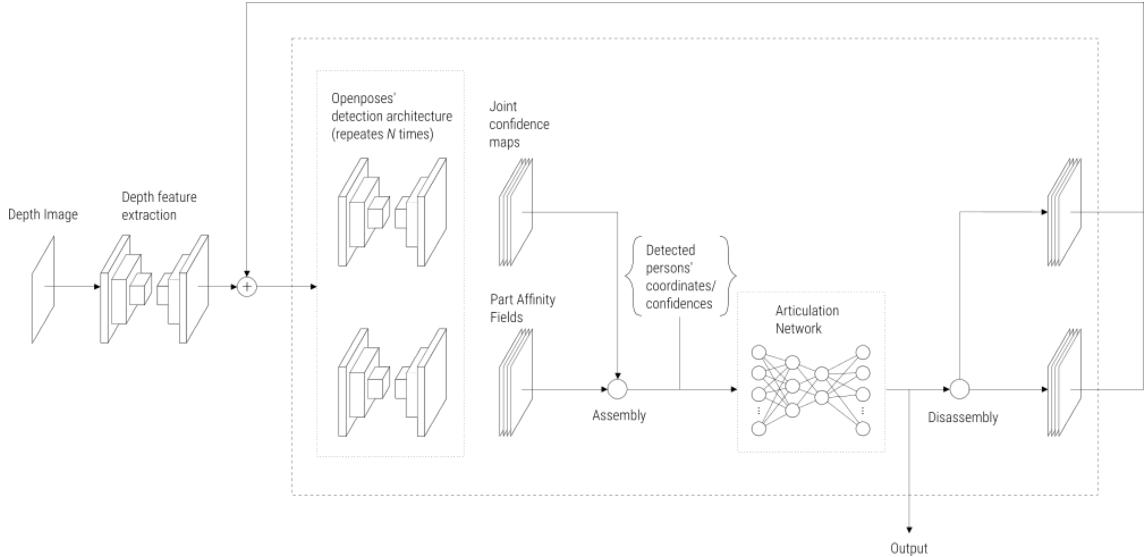


Figure 3.1: Main architecture. CNNs are illustrated as hourglass networks, but this is not necessarily representative of the final architecture. A depth image is fed to a depth-feature extraction network (3.1.1) which in consecutive steps are combined with output PAFs and joint confidence maps, before this is fed to the 3D object detection networks (3.1.2). These are run for a small number of iterations to get some estimate for the PAFs and confidence maps, before we move on to the assembly step. For each of the skeletons detected (with a ceiling for maximum detected skeletons) the skeletons with the strongest evidence are assembled. In the case that a body landmark is not found, we use preconceived coordinates from a standard defined skeleton. The coordinate frame of the skeleton in the articulation step till be the root joint (middle hip) with the z-axis pointing up, and the y-axis orthogonal to the line through the two side hip joints. The articulation network (3.1.3) tries to scale and articulate the skeleton based on detected evidences. In the disassembly step, we place all the skeletons back into the camera coordinate frames, and update the PAFs and confidence maps with confidences and placements from the articulation network, and use them in the following steps. This is repeated a small number of times (3-4).

3.1.1 Depth feature extraction

This network is built from the ground up. Therefore we want some layers to for example detect edges, and one for detecting slanting gradients, or connected surfaces.

For limbs, we might want to find surfaces that are shaped like tubes or oblong spheroids.

3.1.2 3D object detection

In this part of the network we borrow some of the architecture described in [2]. The purpose is to find joints and limbs and put them into PAFs or joint confidence maps.

3.1.3 Articulation network

Here, we get some coordinates for different joints, along with the confidences for those coordinates. The network will try to find out what it thinks the joints with low confidences, or no detections, should be. It is hypothesized that this network will learn things like symmetry (left and right limbs should have the same length), proportionality (limbs should be proportional to each other), possible articulations and natural poses.

For joints which we haven't detected or have very low confidence for, the network will input some standard coordinates for that joint, scaled by the limbs we already have the strongest confidences for. The standard scaling/coordinates is hard-coded, based on the human anatomy surveys in [6]. The exception is eyes and ears, which is set to a best guess. In addition, the depth coordinate for each joint is set to 0. Numbering and a visualization of the skeleton can be seen in Figure 3.2.

The architecture is visualized as a simple fully connected neural network. Though, it might be enough to connect the neurons responsible for connected limbs. A neuron in the second layer connected to the foot, knee and one of the hip-joints does not need to be connected to the inputs from a hand or shoulder. Subsequent hidden layers can however be fully connected.

If we had some input describing the direction of the camera, and the visual hull containing the undetected points, this network may perform better. However, that would require a fundamental change to the network, which is not done in this work.

This network will also be trained separately, as all it needs is poses and random confidences as inputs.

3.2 Training data preparation

Skeleton models

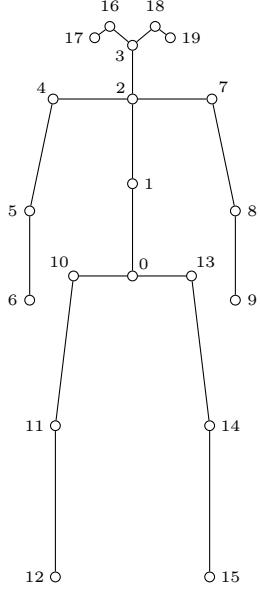


Figure 3.2: Numbering for detected landmarks/keypoint markers.

ID	Description	Std.Coord.
0	Middle hip	(0.00, 0.00)
1	Middle back	(0.00, 1.22)
2	Neck	(0.00, 2.34)
3	Nose	(0.00, 3.05)
4	Right shoulder	(-1.05, 2.34)
5	Right elbow	(-1.36, 0.86)
6	Right wrist	(-1.36, -0.32)
7	Left shoulder	(1.05, 2.34)
8	Left elbow	(1.36, 0.86)
9	Left wrist	(1.36, -0.32)
10	Right hip	(-0.78, 0.00)
11	Right knee	(-1.02, -1.98)
12	Right ankle	(-1.02, -3.98)
13	Left hip	(0.78, 0.00)
14	Left knee	(1.02, -1.98)
15	Left ankle	(1.02, -3.98)
16	Right eye	(-0.30, 3.30)
17	Right ear	(-0.50, 3.15)
18	Left eye	(0.30, 3.30)
19	Left ear	(0.50, 3.15)

Table 3.1: Numberings, names/descriptions and standard coordinates for recognized landmarks

Chapter 4

Experiments

Chapter 5

Future Work

Finding more accurate pose from a temporal algorithm which takes input from a series of depth images, from a single viewpoint.

Finding more accurate pose from multiple viewpoints.

Heart/respiration rate monitoring using frequency search in changing rgb and depth-pixel values for automatically selected RoIs.

Mood detection on facial expressions.

Human activity recognition using 3D pose provided by the method proposed in this paper.

Train the network over a larger dataset in unstructured environments and with multiple people present.

Train an accompanying network that takes a sequence of estimated limb positions and their probability as input, and trying to refine the estimation based on earlier detection. This could also be done through a Kalman filter.

This should all accumulate in an LSTM network for predicting diseases. – requires dataset acquired over possibly years, dispersed over many users, and their daily activities, as possible. Other factors that should be taken into consideration is environmental factors such as humidity, temperature and weather. (As they may be risk factors for certain conditions such as heatstroke or depression.) With such a diverse dataset we could possibly do PCA to determine certain risk factors for different diseases.

Train and test network on the Human 3.6M dataset using TOF data

Chapter 6

Conclusions

Appendices

.1 Depth sensors

Since depth sensors are widely used in different robotics applications for tasks such as SLAM, odometry and object detection, we selected this as our main source of information for monitoring the user. There are mainly three different technologies to choose from: *Structured light*, *Time-of-Flight (ToF)* and *Stereo vision*.

.1.1 Stereo vision

uses two cameras that are observing parts of the same scene. In commercial packages the cameras are usually calibrated, so we have measurements to put into the camera matrix as well as the rotation and translation between the two camera matrices.

However, to get a 3D structure, we need to find common feature points between the two cameras. To do this, we can use various feature descriptors such as ORB, SWIFT and SURF. When good matches has been found between the images, we measure the disparity between the points, and triangulate the distance. The depth measurements for the rest of the image are then calculated by matching pixels close to the found featurepoints.

Since this is an optically based technology, it will work well in well-lit scenes that contain many unique featurepoints. If we operate in an homogenous environment with few, or similar textures it will be difficult to find featurepoints to map the environment. An example of this could be on the seabed or inside buildings with limited light conditions, for example during a blackout.

.1.2 Structured light

uses a projected pattern of light points onto the scene which is registered by a calibrated camera. Usually, the projected light pattern and camera operate in the infrared part of the electromagnetic spectrum¹. This means that in locations where one can expect a lot of IR radiation, this technology will not work very well. Since the IR radiation from the sun usually is much stronger than the one emitted from the projector on the sensor, this technology will not work well outside in well-lit conditions. It will however work inside and in conditions where no external light source are provided.

In addition, since the light is structured and the sensor is calibrated, we can skip the step where we find common featurepoints to triangulate the distance which we have to do in the stereo vision case.

¹The Microsoft Kinect V2 sensor uses a wavelength around 827-850nm stated by a developer in [this forum](#). **TODO : fix source**

.1.3 Time-of-Flight

cameras uses the known constant of c to calculate distances in the image, by measuring the time a light-pulse emitted from the camera uses to be reflected onto the camera sensor. For example, Microsofts Kinect v2 uses a specialized ToF-pixel array in conjunction with a timing generator and modulated laser diodes to obtain per-pixel depth images [7].

As with structured light sensors, this is susceptible to interference from external light sources, or specular surfaces, and has limited range because of light fall-off. However, since the distance calculations are timing based, we can obtain framerates up to 30 fps in the Kinect v2 sensor [5].

.2 Robotic Operating System

In order to make the system easier to use and available to as many platforms as possible, it was decided to create it for the Robotic Operating System (ROS). ROS is a collection of libraries and a runtime environment making communication with different modules and programs on the robot possible.

.3 Classifiers

write a bit about what classifiers are, and how we use them to find the different keypoints in the image.

Glossary

CNN A type of neural network that uses convolved filters to create spacial recognition more robust.

Abbreviations

HAR Human Activity Recognition.

HRI Human Robot Interaction.

MECS Multimodal Elderly Care System.

PAF Part Affinity Field.

ToF Time-of-Flight.

Symbols

c Speed of Light.

Bibliography

- [1] Mykhaylo Andriluka et al. ‘2D Human Pose Estimation: New Benchmark and State of the Art Analysis’. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2014.
- [2] Zhe Cao et al. ‘Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields’. In: *CVPR*. 2017.
- [3] Catalin Ionescu et al. ‘Human3.6M: Large Scale Datasets and Predictive Methods for 3D Human Sensing in Natural Environments’. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 36.7 (July 2014), pp. 1325–1339.
- [4] Helseetaten Oslo Kommune. *Oslohelsa – Kortversjonen, Oversikt over helsetilstand og påvirkningsfaktorer*. June 2016. URL: https://www.oslo.kommune.no/getfile.php/13139280/Innhold/Politikk%20og%20administrasjon/Statistikk/Oslohelsa_kortversjon.pdf.
- [5] Elise Lachat et al. ‘Assessment and Calibration of a RGB-D Camera (Kinect v2 Sensor) Towards a Potential Use for Close-Range 3D Modeling’. In: *Remote Sensing* 7.10 (Oct. 2015), pp. 13070–13097. ISSN: 2072-4292. DOI: [10.3390/rs71013070](https://doi.org/10.3390/rs71013070). URL: <http://dx.doi.org/10.3390/rs71013070>.
- [6] Drillis R. and Contini R. *Body Segment Parameters*. 1966. URL: <http://edge.rit.edu/edge/P13032/public/FinalDocuments/Detailed%20Analysis/Anthropometric%20Data/Drillis%20%26%20Contini.pdf>.
- [7] John Sell and Pat O’Connor. ‘XBOX One Silicon’. Hot Chips 25. Aug. 2013. URL: http://www.hotchips.org/wp-content/uploads/hc_archives/hc25/HC25.10-SoC1-epub/HC25.26.121-fixed-%20XB1%2020130826g.pdf.
- [8] Hao-Yu Wu et al. ‘Eulerian Video Magnification for Revealing Subtle Changes in the World’. In: *ACM Transactions on Graphics (Proc. SIGGRAPH 2012)* 31.4 (2012).