

Moving 3D Pose Estimation with ESP32 Wi-Fi

Name1 Surname^{1,2}[✉], Name2 Surname²[✉], Name3 Surname^{2,3}[✉], Name4 Surname²,
Name5 Surname²[‡], Name6 Surname²[‡], Name7 Surname^{1,2,3}^{*}, with the Lorem Ipsum
Consortium[¶]

1 Affiliation Dept/Program/Center, Institution Name, City, State, Country

2 Affiliation Dept/Program/Center, Institution Name, City, State, Country

3 Affiliation Dept/Program/Center, Institution Name, City, State, Country

[✉]These authors contributed equally to this work.

[‡]These authors also contributed equally to this work.

[✉]Current Address: Dept/Program/Center, Institution Name, City, State, Country

[¶]Deceased

[¶]Membership list can be found in the Acknowledgments section.

^{*} correspondingauthor@institute.edu

Abstract

Doraemon.

Author summary

Doraemon.

Introduction

Wifi is one of common network mediums nowadays. Normally, It is used for establishing a wireless network to connect to the internet. But there are still many more functions Wifi is good at. Wifi can also be applied in fields beside connecting to the internet according to its stability being upgraded continuously. Decent Wifi connectivity can extract more data other than the data to be transmitted like concentration, speed, obstacle between the transmission. Those can be composed to be many useful application on their own like localization, activity prediction and etc.

Camera is a very good tool for monitoring things and is being used as a very effective data collector for mapping to the ground truth to create many popular machine-learning-based usable model like pose estimation, text segmentation, object detection and many more. However, camera may be unavoidably judged as a serious privacy infringement since the data obtained like a photo or video are too clear and possess too much information that might be used in a bad way.

There are many works tried to extract those extractable features like camera's videos does from Wifi. But, they are mostly working with very specific tools and Network Interface Card (NIC) connected to a laptop running Linux that is currently one of the ways allowing to obtain fine-grain Channel State Information (CSI), the descriptive data of the Wifi propagating in that environment. Those limitation significantly decrease steamline of implementation. It is hard for public demonstration and intregation with many updated tools in operating system like Windows or OSX.

Actually, there are other existing ways for obtaining the CSI. One is from a ubiquitously used microprocessor, ESP32. which is still not much be explored in Wifi exploiting field. It is simple to implement and can be easily integrated with others tools in many platforms due to its massively produced external tools.

Human pose estimation is one of the most popular topic in machine-learning-based field. It can be used as visualization directly and also applied for further applications like Game-controlling, Activity classification, Violence detection and many more. Most of the model extract human pose from videos which are taken from camera which conduct privacy issue like statement above. We are expecting to extract human pose from Wifi CSI to solve the issue. So, this paper proposes a machine-learning-based model to create a mapping rule from Wifi CSI to 3D moving human pose estimation by using ESP32.

Background

Pose Estimation

There many machine-learning-based human body pose estimation tools had been proposed and available online. Those can de found both 2D and 3D. In our paper, we chose a light weight 3D pose estimation as an annotation due to its simplicity and the hypothesis that 3D should suit the most in our work. The project can be found at [github-lw3d] and is based on [paper-Lightweight OpenPose] and [paper- Single-Shot Multi-Person 3D Pose Estimation From Monocular RGB]. Its job is to simply create 3D human pose annotation from an image. Then, feed to our works training process as annoations.

Wifi

Wifi is a well-known connectivity with no wire needed (wireless). It has been used as a medium for connecting to the internet for over 10 years. However, the Wifi is the name covering IEEE 802.11 n/g/ac protocols. It mostly deliver data through 2.4/5GHz frequency with multiple channels. The bandwidth in each channel is 22MHz. the data are to be transmitted pararely with multiplexing technique named orthogonal frequency division multiplexing (OFDM). Each carrier may propagate to a reciever with encountering many obstacles. The effect of that situation is the Doppler Effect. So, Channel State Information (CSI) is represented as physical layer indicator that can be used to investigate how each channel propagate to the reciever or back to the transmitter.

If a sender sends data to a reciever through Wifi, the data will be surely not transmitted without any loss.

CSI data

As mentioned in that data propagating to the reciever while touching surrounding environment, the CSI is a variation of the data. The CSI can be found at both sender and reciever since reciever may transmit data back. In this paper, We consider to mainly use CSI at the transmitter. Let the sender use the modulation method of 16-quadrature amplitude modulation (16-QAM) which one carrier can carry 4 bits. When the sender needs to send a '1111', the modulation returns $x = 1 + 1i$ then, transmit to the reciever. At the reciever, let the obtained data is $y = 0.8 + 0.9i$. So, the CSI can be computed by the variation $h = y/x = 0.2 + 3.4i$.

Human body is literally water which reflect radio wave like Wifi. [] and [] have proven that human body can affect the CSI.

ESP32

ESP32 is a very popular single-board computer (SBC). With its affordable price and many available additional tools, ESP32 is commonly used in Internet of Things field. Moreover, it can be applied in research field. Quantitative CSI can be obtained from Wifi in ESP according to [Wi-ESP]. The number of available subcarriers in ESP32 is 64.

According to the detail about Wifi mentioned in , the Wifi in ESP32 has some limitation. It supports only 2.4GHz frequency and can be set only one channel over a connection. The bandwidth of each channel is 22MHz. The CSI can be both obtain from Access point (AP) and station (STA) as shown in Fig 1. The frequency of each channel is as 802.11 standard.

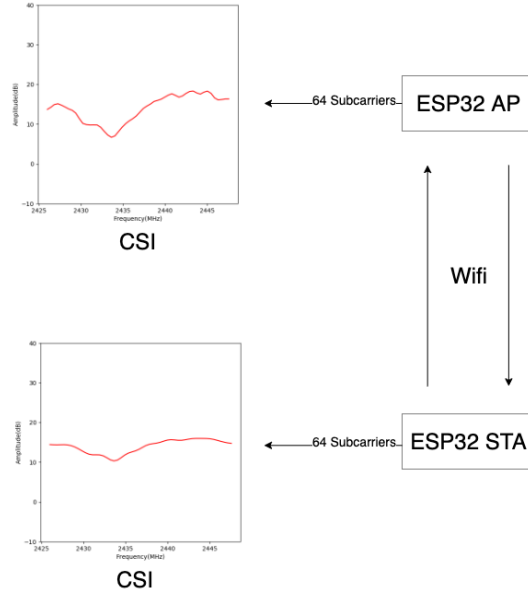


Fig 1. CSI from ESP32s with channel 6.

Materials and methods

Concept

Other famous proposed works like [1] and [2] focus on line-of-sight (LOS) in between AP and STA while our work uses 2 directional Wifi antennae and focus on reflection from human body as shown in Fig 2 on the left.

The reason we name “Moving Pose Estimation” instead of “Pose Estimation” is that CSI is not only affected by human body but mainly by overall environment. This means that 2 corresponding human poses can result obviously different CSIs if the environment around are not exactly the same as shown in Fig 2.

So, the detection of human standing still in every environment is nearly impossible since the CSI of that situation may be found exactly matched to a CSI of the environment that a big bottle of water placed in front of ESP32. In short, if it does not move, we do not if it is human.

Meanwhile, the moving pose is totally different because we focus on its change instead. The example of mapping CSI’s change to Activity Classification can be found

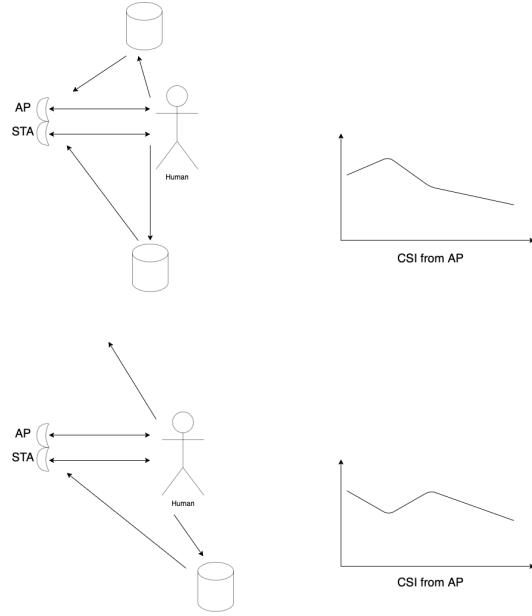


Fig 2. 2 different CSIs resulted from corresponding human poses.

in [1] and [2]. Our work does likewise but focusing on Pose Estimation instead of Activity Classification.

In different environment, the CSIs are different. But, the corresponding moving pose may affect to the same changing pattern of CSI. This hypothesis is investigated in the upcoming parts.

All steps of training method is shown in Fig 3.

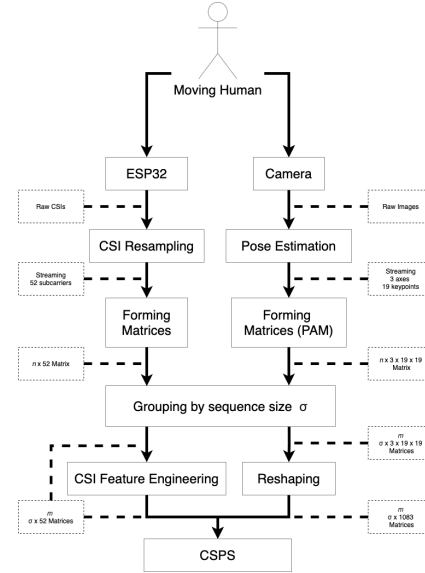


Fig 3. All steps of training method.

Pre-processing

CSI Resampling

As mentioned in , There are 64 subcarriers in CSI data from ESP32 but there are only 52 those are usable while the rest are null. So, we can construst a tensor of 1×52 to represent each CSI. We are to map CSI from the ESP32 to 3D human pose annotation from a camera. The sampling rate of the camera are set to 30Hz. So, we have 30 human pose annotations for one second. For the ESP32, the sampling rate is originally unpredictable and not constant but it is running around 120Hz. So we do a process called “Resampling” to obtain CSI at rate 30Hz in order to map to each human pose annotation.

An example of CSI Resampling is shown in Fig 4. The top graph shows that the the original CSI is logged unstably. The bottom one is to pick a timestamp at rate 30Hz and calculated each with the closet data from the original with a simple mathamatical weight equation as in Eq. 1 in order to predict CSI at timestamp corresponding to each human pose annotation.

$$CSI_{now} = CSI_{before} + \left(\frac{ts_{now} - ts_{before}}{ts_{after} - ts_{before}} \times (CSI_{after} - CSI_{before}) \right) \quad (1)$$

, where

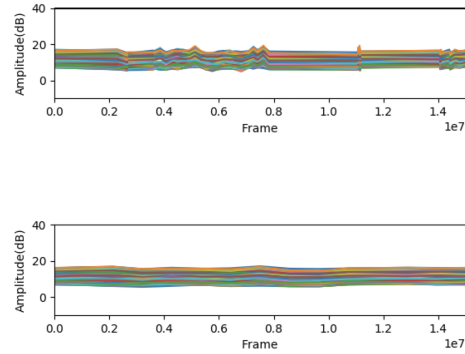


Fig 4. An example of CSI resampling.

Now we can map each CSI sample to each 3D human pose annotation with the corresponding timestamp.

Pose Estimation and PAM formation

We use [github-lw3d] to estimate 3D human pose from image as stated above. the estimation gives us a 19×3 matrix for each image. The 19×3 matrix are used as an annotation where 19 is for keypoints in human body and 3 is for 3 axes coordination as shown in Fig 5. But, the annotation can still possess too much independency. Some alignment of keypoints may lead to some impossible pose e.g. body keypoint found very far from neck keypoint or head keypoint attached to hip keypoint. We assume those poses are quite not possible for normal human pose. To preserve those constraints, we form a pose adjacent matrix (PAM) from an original 19×3 matrix. the PAM is applied for all x, y and z axes. Each are to be form their 19×19 matrix by the following equations.



Fig 5. 19 keypoints of human pose.

$$x'_{i,j} = \begin{cases} x_i - x_j & i \neq j \\ x_i & i = j \end{cases} \quad (2)$$

$$y'_{i,j} = \begin{cases} y_i - y_j & i \neq j \\ y_i & i = j \end{cases} \quad (3)$$

and

$$z'_{i,j} = \begin{cases} z_i - x_j & i \neq j \\ z_i & i = j \end{cases} \quad (4)$$

The PAM is finally a $3 \times 19 \times 19$ matrix created from 3 matrices of x', y' and z' stacked. Apparently, one PAM represent one human pose.

Conclusively, we are making a model by mapping a sequence of 1×52 matrix from CSI to each sequence of $19 \times 3 \times 3$ PAM from moving human pose annotation with the corresponding timestamp as shown in Fig 6.

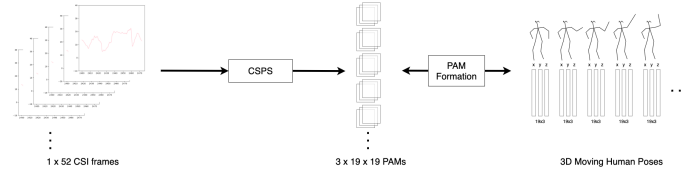


Fig 6. A mapping rule from CSI frames to 3D human poses.

Processing

Mapping CSI and PAM

Let D be a set of synchronized pose and CSI data package. Each pair has corresponding timestamp.

$$D = (C_t, P_t), t \in [1, n], \quad (5)$$

where n is a number of pairs, C is for CSI data from ESP32, P is for PAM annotation as the ground truth, t is the timestamp index when those 2 were collected.

Grouping By Sequence Size

Let σ be an adjustable window size. σ size of D are fed to the model solely. So, $\sigma \leq n$. And, Let m be the number of feeding iteration, $m = \lfloor \frac{n}{\sigma} \rfloor$.

Forming Network Layer

The network swallow m training data as an input, where each is a sequential set of (P_t, C_t) at a corresponding timestamp with size σ . Let Γ be a representation of each sequential set.

$$\Gamma = (C_u, P_u), u \in [1, \sigma], \quad (6)$$

where u is the timestamp index when those 2 were collected. Apparently, the number of Γ is m .

As the Γ is a sequential set with size σ , Long-short term memory (LSTM) [] is suitable for this type of data.

(Optional) CSI Feature Extractor We may extract feature from the CSI before feed it to the model. Or ignore it and feed original CSI to the model directly.

(Optional) CSI Feature Engineering As mentioned in Concept, CSI (C) value implies environment where signal propagating and we need to ignore it in order to universalize the circumstance. It means variousity on CSI data can confuse the model significantly. To solve it, we simply sequentially subtract C_u in each Γ backward to preserve only how CSI changes over each sequence by the following equation,

$$C_u = \begin{cases} 0 & u = 1 \\ C_u - C_{u-1} & u > 1 \end{cases}, C_u \in \Gamma. \quad (7)$$

CSI Sequence to Pose Sequence (CSPS) The summation of layers is shown in Table 1. It is designed to shape a CSI Sequence ($\sigma \times 52$ tensor) to a predicted Pose Sequence ($\sigma \times 3 \times 19 \times 19$ tensor).

Table 1. Layers in CSPS.

Layer (type)	Output Shape	Param #
Bidirectional LSTM	(None, 200)	122400
RepeatVector	(None, σ , 200)	0
Bidirectional LSTM	(None, σ , 200)	240800
TimeDistributed over Dense	(None, σ , 1083)	217683

Table notes Phasellus venenatis, tortor nec vestibulum mattis, massa tortor interdum felis, nec pellentesque metus tortor nec nisl. Ut ornare mauris tellus, vel dapibus arcu suscipit sed.

We first use Bidirectional LSTM as an encoder layer with the input size of $(\sigma, 52)$ to satisfy dimension of sequential C_u in Γ . Then, the repeat vector is added with time σ to make the model treat with correct time-step. Afterward, we placed a decoder layer with another Bidirectional LSTM.

Next, we dense the decoder to be 1083 where can be reshaped into $3 \times 19 \times 19$ later.

Lastly, the TimeDistributed layer is used in order to make the model treat the output for each time-step individually.

Results

Data Collection

dance in ma room ((We collected data under an approval of Carnegie Mellon University IRB 4 . We recruited 8 volunteers, and asked them to do casual daily actions in two

rooms of the campus, one laboratory room and one class room. Floor plans and data collection positions are illustrated in Fig. 8. During the actions, we run the system in Fig. 3 to record CSI samples and videos, simultaneously. For each volunteer, data of his first 80test the networks. The data size of training and testing are 79496 and 19931, respectively))

Experimental Result

All steps of testing method is shown in Fig 7.

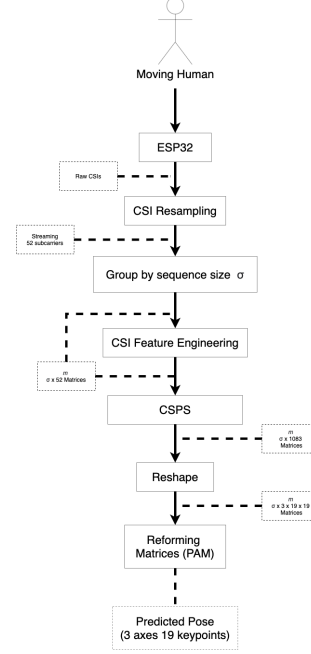


Fig 7. All steps of testing method.

Percentage of Correct Keypoint (PCK) is widely used to evaluate the performance of human pose estimation according to [1].

$$PCK_i@a = \frac{1}{N} \sum_{i=1}^N I\left(\frac{\|pd_i - gt_i\|_2^2}{\sqrt{rh^2 + lh^2}} \leq a\right), \quad (8)$$

where $I(\cdot)$ is a binary indicator that outputs 1 while true and 0 while false, N is the number of frames, i is the index of keypoints that $i \in 1, 2, \dots, 19$. The rh and lh are for the positions of the right shoulder and the left hip from the ground truth, respectively. So, the $\sqrt{rh^2 + lh^2}$ is considered as the length of the upper limb from the ground truth, which is used to normalize the prediction error length $\|pd_i - gt_i\|_2^2$, and pd_i and gt_i are coordinates of prediction and ground-truth at the keypoint i respectively.

Table 2 shows the estimation performance of 19 body keypoint in PCK@5, PCK@10, PCK@20, PCK@30, PCK@40, and PCK@50 of the σ 15.

Github¹.

Discussion

Doraemon.

¹<https://github.com/rtmtree/>

Table 2. Table of PCK.

Order	Keypoint	PCK@5	PCK@10	PCK@20	PCK@30	PCK@40	PCK@50
1							
2							
3							
4							
5							
6							
7							
8							
9							
10							
11							
12							
13							
14							
15							
16							
17							
18							
19							

Table notes Phasellus venenatis, tortor nec vestibulum mattis, massa tortor interdum felis, nec pellentesque metus tortor nec nisl. Ut ornare mauris tellus, vel dapibus arcu suscipit sed.

Table 3. Table caption Nulla mi mi, venenatis sed ipsum varius, volutpat euismod diam.

Heading1				Heading2			
cell1row1	cell2 row 1	cell3 row 1	cell4 row 1	cell5 row 1	cell6 row 1	cell7 row 1	cell8 row 1
cell1row2	cell2 row 2	cell3 row 2	cell4 row 2	cell5 row 2	cell6 row 2	cell7 row 2	cell8 row 2
cell1row3	cell2 row 3	cell3 row 3	cell4 row 3	cell5 row 3	cell6 row 3	cell7 row 3	cell8 row 3

Table notes Phasellus venenatis, tortor nec vestibulum mattis, massa tortor interdum felis, nec pellentesque metus tortor nec nisl. Ut ornare mauris tellus, vel dapibus arcu suscipit sed.

Conclusion

Doraemon.

Supporting information

S1 Fig. Bold the title sentence. Add descriptive text after the title of the item (optional).

S2 Fig. Lorem ipsum. Maecenas convallis mauris sit amet sem ultrices gravida. Etiam eget sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec euismod ligula. Curabitur fringilla pulvinar lectus consectetur pellentesque.

S1 File. Lorem ipsum. Maecenas convallis mauris sit amet sem ultrices gravida. Etiam eget sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec euismod ligula. Curabitur fringilla pulvinar lectus consectetur pellentesque.

S1 Video. Lorem ipsum. Maecenas convallis mauris sit amet sem ultrices gravida. 205
Etiam eget sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec euismod ligula. 206
Curabitur fringilla pulvinar lectus consectetur pellentesque. 207

S1 Appendix. Lorem ipsum. Maecenas convallis mauris sit amet sem ultrices 208
gravida. Etiam eget sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec 209
euismod ligula. Curabitur fringilla pulvinar lectus consectetur pellentesque. 210

S1 Table. Lorem ipsum. Maecenas convallis mauris sit amet sem ultrices gravida. 211
Etiam eget sapien nibh. Sed ac ipsum eget enim egestas ullamcorper nec euismod ligula. 212
Curabitur fringilla pulvinar lectus consectetur pellentesque. 213

Acknowledgments 214

We thank Sirindhorn International Institute of Technology for providing technical 215
environment and supportive information. 216

References

1. Conant GC, Wolfe KH. Turning a hobby into a job: how duplicated genes find new functions. Nat Rev Genet. 2008 Dec;9(12):938–950.
2. Ohno S. Evolution by gene duplication. London: George Alien & Unwin Ltd. Berlin, Heidelberg and New York: Springer-Verlag.; 1970.
3. Magwire MM, Bayer F, Webster CL, Cao C, Jiggins FM. Successive increases in the resistance of Drosophila to viral infection through a transposon insertion followed by a Duplication. PLoS Genet. 2011 Oct;7(10):e1002337.