

NAO robot that imitates human movements using Mediapipe

Costilla Caballero, Salvador
Applied Robotics
Hof University of Applied Science
CDMX, Mexico
salvadorcoscab@gmail.com

Abstract—In this project I made a NAO robot, from Aldebaran, move imitating human movements. This movements were captured with a webcam using Mediapipe. Though this human imitation has already been done, I wanted to make it in a different way. Because most of the other projects use multiple cameras or depth sensors, such as the LIDAR sensor. The project is way far from being finished, because the robot only moves his hands and not very precisely. But I think is a good start for future work on computer vision and robotics.

Index Terms—NAO, Mediapipe, Human imitation, aldebaran, robotics

I. STATE OF THE ART

In this section I will describe different ways to move a NAO robot the Mediapipe Framework and a work related to human imitation using the NAO robot.

The NAO robot documentation [7] states that there is two ways to move the NAO robot to do poses. The first one is using the angles from the joints and the second one is using inverse kinematics. Though there is more ways to move the robot, those are more for walking or to perform a specific pose.

There is two way to control the effectors to perform a pose, regardless we use inverse kinematics or not, which are:

- Animation methods, in which the time is fixed and there are blocking functions.
- Reactive methods, that can be used to control the robot in real time, and there are non-blocking functions.

We have various functions contained in the *ALMotionProxy* module, related to the movement by controlling the joints some of the animation methods are:

- *ALMotionProxy.angleInterpolation*, that allows you to move the robot from a specific position to another in a specific time.
- *ALMotionProxy.angleInterpolationWithSpeed*, that allows you to move the robot from a specific position to another in a specific time and speed.

One important method is:

- *ALMotionProxy.setAngles*, that allows you to move set the angles of a single joint or a set of joints of the robot.

As it was stated before the robot can be controlled using inverse kinematics as well, in this case a classical IK solver which uses only the joints of the effector chain to reach a

target and a generalized IK solver (also called Whole Body control) which uses all the robot joints to reach a target. Some of the animation and reactive methods related to the inverse kinematics that I consider important are:

- *ALMotionProxy.positionInterpolation* that allows you to move the robot from a specific position to another in a specific time
- *ALMotionProxy.setPosition*, that allows you to move the robot to a specific position in a 3d space.

BlazePose [7], is an architecture for single person 3D pose detection based in convolutional neural networks. It is available in the Mediapipe framework. And allows you to detect poses from a video or an image. There are 33 landmarks that the BlazePose model can detect (See figure 1). These

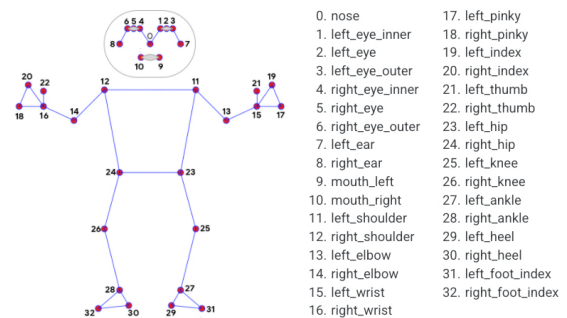


Fig. 1. Landmarks from the Mediapipe library.

landmarks can be read as a set of points in a 3D space. The Landmarker has various configuration options, such as:

- *model_complexity*, that allows you to choose between three different models, the lite, the full and the heavy. Depending on the model chosen the accuracy and the speed can increase.
- *min_detection_confidence*, that allows you to set the minimum confidence that a pose is detected.
- *min_tracking_confidence*, that allows you to set the minimum confidence that a pose is tracked.

In [7] the authors use the NAO robot to imitate the movements of a human using a kinect sensor. They use a client-server approach to send the data from the kinect sensor to the robot. In the client part they use the *pykinect2* library to get the data from the kinect sensor, this data are 25 landmarks that

represent the human body. These data is then sent to the server part using the *socket* library. In the server part they calculate the angles between the joints and then send them to the robot using the *ALMotionProxy.AngleInterpolation* method.

II. APPROACH

A. Objectives

During this project I had some objectives which were:

- Solve the computer vision part, that in this case was done with Mediapipe.
- Make the robot move using the NAOqi API.

The problems were that, in order to work with the NAOqi API, I had to download and setup the SDK which is written (by the moment of this project) in Python 2.7 and the Mediapipe is only available in Python 3.7 or higher. That is why I had to add as an objective to make the two libraries work together.

B. Methodology

The methodology that I used to solve the objectives above was to first divide the problem into modules and tried to solve them separately, and then join them together and solve the problems that arose from the connection between the two modules. The main two modules were the computer vision and the robot movement.

C. Computer vision

In the beginning, the computer vision part was done very easily, because Mediapipe has a lot of usage examples and good documentation. On the other hand working with the NAO robot was a much more difficult.

For the computer vision part I used the opencv library to get images from both camera and videos and then I used the Mediapipe library to get the poses from the images. I used also the opencv library to plot the landmarks that the Mediapipe library detected, so I could see how the landmarks were being detected. This landmarks conform an object with 33 points, which represents landmarks of the human body. (See figure 1)

D. Robot movement

Then I started working with the robot movement part, to do this I used the NAOqi API. I realized, after reading the NAOqi documentation that there is two ways to control the robot for doing poses. The first one was using angles and the second one was using inverse kinematics. I decided to use the first one because I thought that with the landmarks, obtained with Mediapipe, I could easily get the angles of each join of the robot, and also the project that I found on internet used this approach. All of this joint-angle approach is written in the *ALMotion* module of the NAOqi API. In summary the robot has various joints, and each joint can be moved using a specific angle. These joints work in a similar way as the human body joints, allowing some joints to pitch, roll, and yaw. (See figure 2)

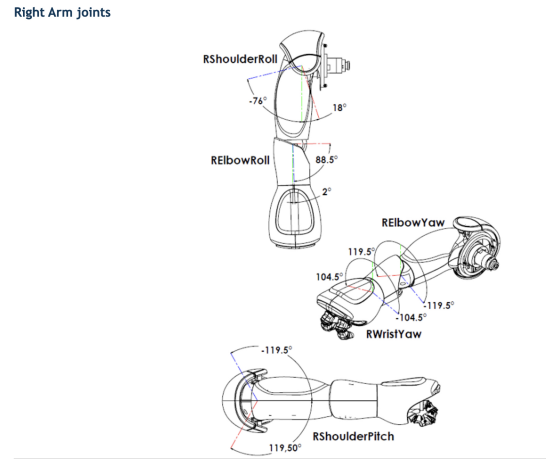


Fig. 2. Landmarks from the Mediapipe library.

E. Connection between the Computer vision and the Robot movement

After doing some testing on both parts, I had to somehow connect them. In order to do the connection between the robot and the computer vision part, I tried with *execnet*, which is a library that allows you to trigger python scripts from another python script, even with different versions of python. But I had some problems sending the data from the computer vision part to the part in which I move the robot, so I decided to try a different approach.

I decided to use sockets to send the data, as I got inspiration from a project that I found with a similar approach, in which the author used sockets to send the data from a kinect sensor to the robot. This solved everything, about the connection between the scripts, because it uses the client-server approach, allowing me to use python 3.10 for the computer vision part that is the client and python 2.7 for the Robot movement part that is the server. I just had to do some encoding to send the data.

F. Testing the implemented code and solving problems

After changing the code that I found to work with Mediapipe and OpenCV instead of the kinect sensor. I did some testing on the functions that were already implemented in the server side, but they performed very bad, but at least, I saw that the robot moved.

Also the display of the images was very slow since the code had a delay to send the data, in order for the server to process it. I solved this using a thread to get the data from the displayed image and send it to the server with the thread without messing up with the display.

As I said before after changing a lot the parameters and definition of the functions that were implemented before I realized that I was not sending the right data to the robot from the server, so I decided to change the approach and instead of sending the coordenates of the landmarks to the server, I decided to send the angles between them and send them directly to the server. This would help me to see what angles

I was sending to the robot by showing them on the displayed image.

After I was able to see the angles that were being sent to the robot, I realized that the angles sent were not very good in the function definition, this is because the kinect sensor and Mediapipe have different coordinate systems, so I had to change the definition of the functions.

This did not work very well, though the movements were more accurate, the robot was not imitating the movement of the hands perfectly. So I decided to change the function. Instead of using geometry I tried with vectors and the angles between them.

G. Vector approach

I saw that as I had the landmarks, that represents a point in the space, I could calculate a vector between two points (See equation 1), and I could use the angles between the vectors to move the robot (See equation 2). But then I realized that I could also calculate the angles between the projection of the vectors (See equation 3) in a generated plane from a normal vector (See equations 4a,4b,4c). I tried this for one joint, that was the right shoulder pitch. I got the normal vector from the left shoulder to the right shoulder to generate the plane, and then I got the vectors from the right shoulder to the right elbow and from the shoulder to the right hip. I projected these vectors in the generated plane and then I calculated the angles between them.

$$\vec{v} = \vec{p}_2 - \vec{p}_1 \quad (1)$$

$$\cos(\theta) = \frac{\vec{v} \cdot \vec{w}}{\|\vec{v}\| \|\vec{w}\|} \quad (2)$$

$$proj_H \vec{v} = (\vec{v} \cdot \vec{u}_1) \vec{u}_1 + (\vec{v} \cdot \vec{u}_2) \vec{u}_2 \quad (3)$$

$$\vec{n} \cdot (\vec{p} - \vec{p}_0) = 0 \quad (4a)$$

$$\vec{p} = \begin{bmatrix} x \\ y \\ z \end{bmatrix} \quad (4b)$$

$$\vec{p}_0 = \begin{bmatrix} x_0 \\ y_0 \\ z_0 \end{bmatrix} \quad (4c)$$

I implemented a function in python to do the calculations mentioned above and this function did not work well neither when I tested it, and as could see the angles on the displayed image, I realized that the Mediapipe library did not detect the z axis very well. Then I remembered that the kinect sensor has a LIDAR sensor, that allows him to detect, not only position in x and y, but also depth in a better way than the Mediapipe library. That is why I did not continue with this approach and as I did not have more time to finish the project, I decided to leave the functions that were already implemented.

III. CONTRIBUTION

The main difference between my project and the others is that it has a different approach to solve the problem, relying only on one camera and machine learning to detect and recreate the human poses, instead of using a lot of sensors. This is a contribution to the state of the art, because it implies that if there were better machine learning models, the robot could imitate the human movements in a better way.

IV. EXPERIMENTS

V. CONCLUSION AND FUTURE WORK

ACKNOWLEDGMENT

The authors would like to thank the support of the Hof University of Applied Science and the support of the Applied Robotics department.

REFERENCES

- [1] G. Eason, B. Noble, and I. N. Sneddon, "On certain integrals of Lipschitz-Hankel type involving products of Bessel functions," *Phil. Trans. Roy. Soc. London*, vol. A247, pp. 529–551, April 1955.
- [2] J. Clerk Maxwell, *A Treatise on Electricity and Magnetism*, 3rd ed., vol. 2. Oxford: Clarendon, 1892, pp.68–73.
- [3] I. S. Jacobs and C. P. Bean, "Fine particles, thin films and exchange anisotropy," in *Magnetism*, vol. III, G. T. Rado and H. Suhl, Eds. New York: Academic, 1963, pp. 271–350.
- [4] K. Elissa, "Title of paper if known," unpublished.
- [5] R. Nicole, "Title of paper with only first word capitalized," *J. Name Stand. Abbrev.*, in press.
- [6] Y. Yorozu, M. Hirano, K. Oka, and Y. Tagawa, "Electron spectroscopy studies on magneto-optical media and plastic substrate interface," *IEEE Transl. J. Magn. Japan*, vol. 2, pp. 740–741, August 1987 [Digests 9th Annual Conf. Magnetism Japan, p. 301, 1982].
- [7] M. Young, *The Technical Writer's Handbook*. Mill Valley, CA: University Science, 1989.