

# 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, we programmed a NAO robot from Aldebaran to mimic human movements captured by a webcam using Mediapipe. While this type of human imitation has been done before, we aimed to approach it differently. Unlike most other projects utilizing multiple cameras or depth sensors like the LIDAR sensor, our focus was on a novel methodology. The project is still in progress, with the robot currently capable of moving its arms, albeit not with high precision. Nonetheless, we see this as a promising starting point for future developments in computer vision and robotics.

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

## I. STATE OF THE ART

In this section, we will explore various methods for controlling the movement of a NAO robot using the Mediapipe Framework, as well as discuss related work involving human imitation with the NAO robot.

According to the NAO robot documentation [1], there are two primary methods for posing the NAO robot: using joint angles and employing inverse kinematics. While other methods exist for walking or performing predefined poses, these two are fundamental to our project's scope.

Regardless of whether inverse kinematics is employed, there are two main ways to control the effectors to achieve a pose:

- Animation methods, which operate within fixed time frames and involve blocking functions.
- Reactive methods, suitable for real-time control with non-blocking functions.

Within the *ALMotionProxy* module, numerous functions pertain to joint movement. Some notable animation methods include:

- *ALMotionProxy.angleInterpolation*, enabling movement from one specific position to another within a specified time.
- *ALMotionProxy.angleInterpolationWithSpeed*, facilitating movement between specific positions within defined time and speed parameters.

An essential method is:

- *ALMotionProxy.setAngles*, enabling the manipulation of single joint angles or sets of joint angles on the robot.

Furthermore, the robot's control extends to inverse kinematics. Here, classical IK solvers and generalized IK solvers

(Whole Body control) are employed. Some pertinent animation and reactive methods related to inverse kinematics include:

- *ALMotionProxy.positionInterpolation*, facilitating movement between specific positions within a specified time frame.
- *ALMotionProxy.setPosition*, enabling precise positioning of the robot in 3D space.

BlazePose [2] is an architecture for single-person 3D pose detection based on convolutional neural networks. It is integrated into the Mediapipe framework, enabling the detection of poses from videos or images. The model is capable of detecting 33 landmarks (see Figure 1), which can be interpreted as points in a 3D space.

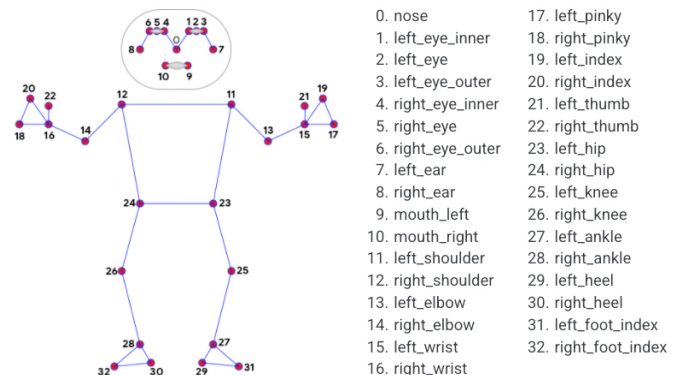


Fig. 1. Landmarks from the Mediapipe library.

The Landmarker offers various configuration options, including:

- *model\_complexity*, This parameter allows selection among three different models - lite, full, and heavy. Depending on the chosen model, both accuracy and speed can vary.
- *min\_detection\_confidence*, It enables setting the minimum confidence level required for detecting a pose.
- *min\_tracking\_confidence*, This parameter allows setting the minimum confidence level required for tracking a pose.

In a thesis titled "Mimetización de movimiento por robot antropomorfo basado en imágenes de sensor de profundidad" [3], the authors utilize the NAO robot to mimic human

movements using a Kinect sensor. They adopt a client-server approach to transmit data from the Kinect sensor to the robot via the NAOqi API. In the client-side implementation, they employ the *pykinect2* library to retrieve data from the Kinect sensor, which consists of 25 landmarks representing the human body. This data is then relayed to the server-side using the *socket* library. On the server-side, the angles between the joints are computed and transmitted to the robot using the *ALMotionProxy.AngleInterpolation* method.

## II. APPROACH

### A. Objectives

During this project we 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, we 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 we had to add as an objective to make the two libraries work together.

### B. Methodology

The methodology that we used to solve the objectives above was to first divide the problem into modules and try 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, such as the video "AI Pose Estimation with Python and MediaPipe" [4], that helped us a lot on getting the landmarks and calculate some angles. On the other hand working with the NAO robot was a much more difficult, because we did not know how to start.

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

### D. Robot movement

Then we started working with the robot movement part, to do this we used the NAOqi API. We 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. We decided to use the first one because we thought that with the landmarks, obtained with Mediapipe, we could easily get the angles of each join of the robot, and also the project that we 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 II-D)

Right Arm joints

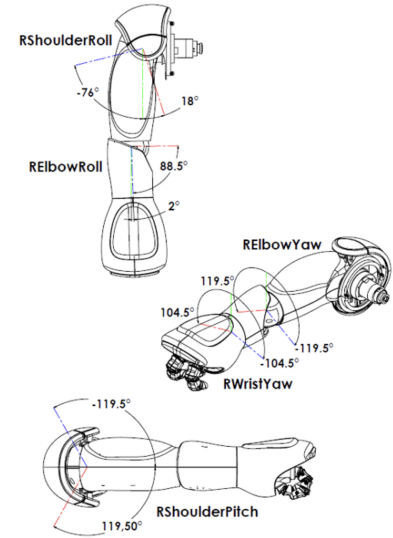


Fig. 2. Joints in the arm of the NAO robot.

### E. Connection between the Computer vision and the Robot movement

After doing some testing on both parts, we had to somehow connect them. In order to do the connection between the robot and the computer vision part, we 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 we had some problems sending the data from the computer vision part to the part in which we move the robot, so we decided to try a different approach.

We decided to use sockets to send the data, as we got inspiration from a project that we 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 us 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. We just had to do some encoding to send the data.

### F. Testing the implemented code and solving problems

After changing the code that we found to work with Mediapipe and OpenCV instead of the kinect sensor. We did some testing on the functions that were already implemented in the server side, but they performed very bad, but at least, we 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. We 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 we said before after changing a lot the parameters and definition of the functions that were implemented before we realized that we was not sending the right data to the robot from the server, so we decided to change the approach and instead of sending the coordinates of the landmarks to the server, we decided to send the angles between them and send them directly to the server. This would help us to see what angles we was sending to the robot by showing them on the displayed image.

After we was able to see the angles that were being sent to the robot, we 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 we had to change the defintion 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 we decided to change the function. Instead of using geometry we tried with vectors and the angles between them.

#### G. Vector approach

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

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

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

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

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

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

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

We implemented a function in python to do the calculations mentioned above and this function did not work well neither when we tested it, as we could see the angles on the displayed

image, we realized that the Mediapipe library did not detect the z axis very well, even with a more complex model. Then we 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 we did not continue with this approach and as we did not have more time to finish the project, we decided to leave the functions that were already implemented.

### III. CONTRIBUTION

The main difference between our 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. There is also a contribution to the work done before in "Mimetización de movimiento por robot antropomorfo basado en imágenes de sensor de profundidad" [3], because we can see the angles that are being sent to the robot, which can be useful for further research.

### IV. EXPERIMENTS

To do the experiments and testing we used Choreographe, which has a simulator mode, in which you can see a NAO robot and move it using the localhost as the robot IP (See figure 3). We used this to test the robot movements, but we also used a real robot to see the movements in a real robot (See figure 4). We had a script to try the robot movments and another script just to try the computer vision part. The union of both scripts, as we said before, was done using sockets. So we had to run two scripts in the terminal, one for the server and one for the client.

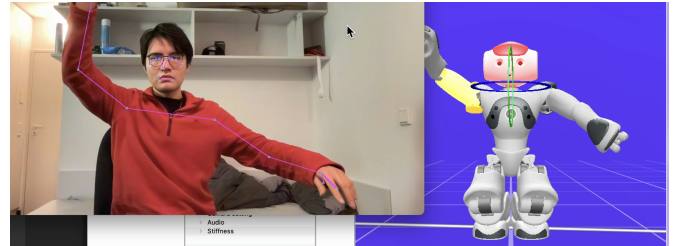


Fig. 3. NAO robot in the Choreographe simulator.

We used a webcam and videos to see the angles that are being sent to the robot. To first initaitate the application you have to run the server part, using a python 2.7 interpreter. This can be done using the following command, in which you have to specify the ip address of the robot as an argument:

```
python2 server.py "ip_address"
```

Then you have to run the client part, using a python 3.10 interpreter. If you want to send the landmarks using a webcam you can use the following command:

```
python3 client.py
```



Fig. 4. NAO robot moving.

On the other hand, if you want to send the landmarks using a video you can use the following command:

```
python3 client.py "path_to_video"
```

Sometimes the robot falls and if this occurs the application might not work or would stop working, so we had to be careful about this, because sudden movements of the robot could make it fall. If the application stops working, there is no other option than to restart the application repeating the steps above.

## V. CONCLUSION AND FUTURE WORK

Even though the project is not finished, we think is a good start for just using computer vision to apply it in robotics. There is a lot of work to do, such as improving the calculation of the angles or use both inverse kinematics and programming the angles to move the robot.

In order to use the legs of the robot is necessary to do further research, because we have to consider many other things such as the mass center of the robot, and take care that the movements are not too abrupt to make the robot fall.

## ACKNOWLEDGMENT

We want to thank our professor, Dr. Christian Groth for helping us during the project and thank also to the Hof

University of Applied Science for giving us the opportunity to work with the NAO robot, which was a great experience.

## REFERENCES

- [1] NAOqi API, <http://doc.aldebaran.com/2-5/naoqi/motion/index.html>
- [2] Mediapipe Pose Landmarking [https://developers.google.com/mediapipe/solutions/vision/pose\\_landmarker](https://developers.google.com/mediapipe/solutions/vision/pose_landmarker)
- [3] I. Irigoyen, "Mimetización de movimiento por robot antropomorfo basado en imágenes de sensor de profundidad," Trabajo Fin de Grado, Facultad de Informática, Grado de Ingeniería Informática, Universidad del País Vasco, 2020.
- [4] Nicholas Renotte - Robotics and AI, "AI Pose Estimation with Python and MediaPipe," YouTube, 2021. [Online]. Available: [https://www.youtube.com/watch?v=06TE\\_U21FK4t=2109sab\\_channel=NicholasRenotte](https://www.youtube.com/watch?v=06TE_U21FK4t=2109sab_channel=NicholasRenotte).