

Gesture Recognition for Human-Robot Interaction: An approach based on skeletal points tracking using depth camera

Masterarbeit

am Fachgebiet Agententechnologien in betrieblichen Anwendungen und der

Telekommunikation (AOT)

Prof. Dr.-Ing. habil. Sahin Albayrak

Fakultät IV Elektrotechnik und Informatik

Technische Universität Berlin

vorgelegt von

Sivalingam Panchadcharam Aravindh

Betreuer: Prof. Dr.-Ing. habil. Sahin Albayrak,

Dr.-Ing. Yuan Xu

Sivalingam Panchadcharam Aravindh

Matrikelnummer: 342899

Sparrstr. 9

13353 Berlin

Statement of Authorship

I declare that I have used no other sources and aids other than those indicated. All passages quoted from publications or paraphrased from these sources are indicated as such, i.e. cited and/or attributed. This thesis was not submitted in any form for another degree or diploma at any university or other institution of tertiary education

Place, Date

Signature

Abstract

Human-robot interaction (HRI) has been a topic of both science fiction and academic speculation even before any robots existed [?]. HRI research is focusing to build an intuitive and easy communication with the robot through speech, gestures, and facial expressions. The use of hand gestures provides an attractive alternative to complex interfaced devices for HRI. In particular, visual interpretation of hand gestures can help in achieving the ease and naturalness desired for HRI. This has motivated a very active research concerned with computer vision-based analysis and interpretation of hand gestures. Important differences in the gesture interpretation approaches arise depending on whether 3D based model or appearance based model of the gesture is used [?].

In this thesis, we attempt to implement the hand gesture recognition for robots with modeling, training, analyzing and recognizing gestures based on computer vision and machine learning techniques. Additionally, 3D based gesture modeling with skeletal points tracking will be used. As a result, on the one side, gestures will be used command the robot to execute certain actions and on the other side, gestures will be translated and spoken out by the robot.

We further hope to provide a platform to integrate Sign Language Translation to assist people with hearing and speech disabilities. However, further implementations and training data are needed to use this platform as a full fledged Sign Language Translator.

Keywords

Human-Robot Interaction (HRI), Computer Vision, Depth Camera, Hand Gesture, 3D hand based model, Skeleton tracking, Gesture Recognition, Sign Language Translation, Hidden Markov Model, NAO

Acknowledgements

Der Punkt Acknowledgements erlaubt es, persönliche Worte festzuhalten, wie etwa:

- Für die immer freundliche Unterstützung bei der Anfertigung dieser Arbeit danke ich insbesondere...
- Hiermit danke ich den Verfassern dieser Vorlage, für Ihre unendlichen Bemühungen, mich und meine Arbeit zu foerdern.
- Ich widme diese Arbeit

Die Acknowledgements sollte stets mit großer Sorgfalt formuliert werden. Sehr leicht kann hier viel Porzellan zerschlagen werden. Wichtige Punkte sind die vollständige Erwähnung aller wichtigen Helfer sowie das Einhalten der Reihenfolge Ihrer Wichtigkeit. Das Fehlen bzw. die Hintanstellung von Personen drückt einen scharfen Tadel aus (und sollte vermieden werden).

Contents

Statement of Authorship	II
Abstract	III
Contents	V
1 Introduction	1
2 Background	2
2.1 Computer Vision in Robotics	2
2.2 NAO - The Humanoid Robot	2
2.2.1 Construction	3
2.2.2 Motion	3
2.2.3 Audio	5
2.2.4 Vision	7
2.2.5 Computing	8
2.2.6 Networking	8
2.3 Gesture Recognition	8
2.3.1 Gesture Modeling	9
2.3.2 Gestural Taxonomy	10
2.3.3 Feature Extraction	11
2.3.3.1 OpenNI 2	11
2.3.3.2 NiTE 2	13
2.3.3.3 Algorithm	14
2.3.4 Gesture Classification and Prediction	19
2.3.4.1 Adaptive Naive Bayes Classifier (ANBC)	19
2.3.4.2 Gesture Recognition Toolkit (GRT)	21
Bibliography	27

List of Figures	28
List of Tables	29
Abbreviations	30
Attachments	31
A HRI Module	31
.1 HRI Module	31
.1.1 Gesture Tracker	31

Chapter 1

Introduction

Huge influence of computers in society has made smart devices, an important part of our lives. Availability and affordability of such devices motivated us to use them in our day-to-day living. The list of smart devices includes personal automatic and semi-automatic robots which are also playing a major role in our household. For an instance, Roomba [?] is an autonomous robotic vacuum cleaners that automatically cleans the floor and goes to its charging station without human interaction.

Interaction with smart devices has still been mostly through displays, keyboards, mouse and touch interfaces. These devices have grown to be familiar but inherently limit the speed and naturalness with which we can interact with the computer. Usage of robots for domestic and industrial purposes has been continuously increasing. Thus in recent years, there has been a tremendous push in research toward an intuitive and easy communication with the robot through speech, gestures and facial expressions.

Tremendous progress had been made in speech recognition and several commercially successful speech interfaces are available. However, speech recognition systems have certain limitations such as misinterpretation due to various accents and background noise interference. It may not be able to differentiate between your speech, other people talking and other ambient noise, leading to transcription mix-ups and errors.

Furthermore, there has been an increased interest in recent years in trying to introduce other human-to-human communication modalities into HRI. This includes a class of techniques based on the movement of the human arm and hand, or hand gestures. The use of hand gestures provides an attractive alternative for Human-robot interaction than the conventional cumbersome devices.

Chapter 2

Background

2.1 Computer Vision in Robotics

Computer vision is a broad field that includes methods for acquiring, processing, analyzing and understanding images and, in general, high-dimensional data from the real world in order to produce numerical or symbolic information, e.g., in the forms of decisions [?].

Proper vision is the utmost importance for the function of any vision based autonomous robot. Areas of artificial intelligence deal with autonomous planning or deliberation for robotic systems to navigate through an environment. A detailed understanding of these environments is required to navigate through them. High-level information about the environment could be provided by a computer vision system that is acting as a vision sensor.

In this thesis, we will focus on the hand gesture recognition using computer vision techniques for a humanoid robot named as NAO, as shown in the figure ???. NAO is an autonomous, programmable humanoid robot developed by Aldebaran Robotics. The NAO Academics Edition was developed for universities and laboratories for research and education purposes. Table ?? shows the specification of NAO according to Aldebaran Robotics.

2.2 NAO - The Humanoid Robot

NAO is an autonomous programmable humanoid robot invented by Aldebaran Robotics. NAO Academics Edition was developed for universities and laboratories for research

and educational purposes. Follow subsections discuss briefly about the specifications of NAO as described by Aldebaran Robotics.

Table 2.1: NAO V5 specification

Height	58 centimetres (23 in)
Weight	4.3 kilograms (9.5 lb)
Battery autonomy	60 minutes (active use), 90 minutes (normal use)
Degrees of freedom	21 to 25
CPU	Intel Atom @ 1.6 GHz
Built-in OS	Linux
SDK compatibility	Windows, Mac OS, Linux
Programming languages	C++, Python, Java, MATLAB, Urbi, C, .Net
Vision	2 x HD 1280x960 cameras
Connectivity	Ethernet, Wi-Fi
Sensors	4 x directional microphones 1 x sonar rangefinder 2 x IR emitters and receivers 1 x inertial board 9 x tactile sensors 8 x pressure sensors

2.2.1 Construction

NAO has a body with 25 degrees of freedom (DOF) whose key elements are electric motors and actuators as show in the figure 2.1. It has 48.6-watt-hour battery that provides NAO with 1.5 or more hours of autonomy, depending on usage. Additional specifications of NAO are shown in the table 2.1. —— Add more info ——

2.2.2 Motion

NAOs walking algorithm uses a simple dynamic model (linear inverse pendulum) and quadratic programming. It is stabilized using feedback from joint sensors. This makes the walking robust and resistant to small disturbances, and torso oscillations in the frontal and lateral planes are absorbed. It can walk on a variety of floor surfaces, such as carpeted, tiled, and wooden floors.

NAOs motion module is based on generalized inverse kinematics, which handles Cartesian coordinates, joint control, balance, redundancy, and task priority. This means that when asking it to extend its arm, it bends over because its arms and leg joints are taken into account.

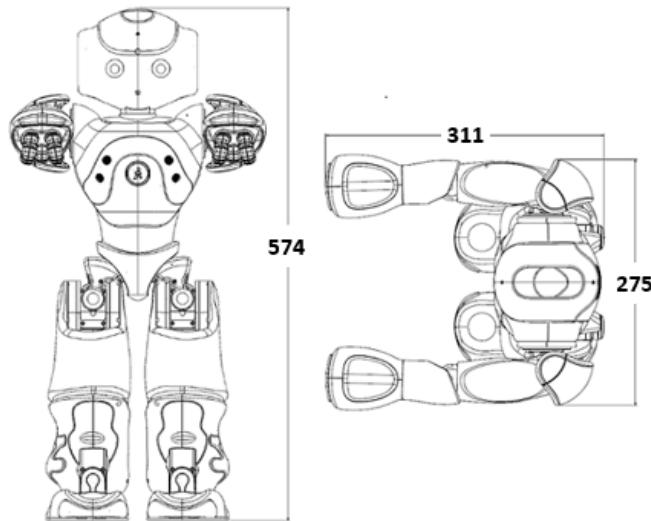


Figure 2.1: Construction of NAO

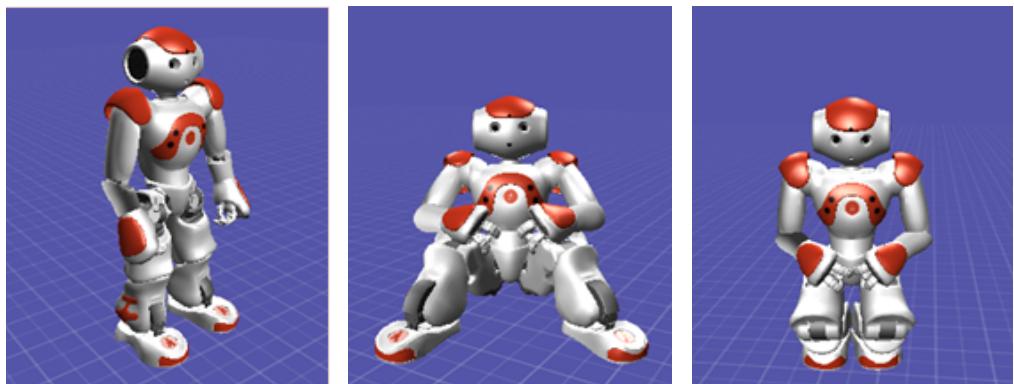


Figure 2.2: Standing, Sitting and Crouching positions of NAO using ALRobotPosture proxy

In this thesis, we used locomotion and stiffness control of Motion API to move NAO to a position in the two dimensional space. Robot Posture API was also used to make the robot go to the predefined posture such as Stand, Sit and Crouch as shown in the figure 2.2. Python code 2.2.2 shows how the robot can be moved to another position at the given normalized velocity using Motion API.

```
# To move the robot at the given normalized velocity using
Aldebaran Motion and Posture API

import time
from naoqi import ALProxy

robotIP = "nao.local"
```

```
PORT = 9559
motionProxy = ALProxy("ALMotion", robotIP, PORT)
postureProxy = ALProxy("ALRobotPosture", robotIP, PORT)

# Wake up robot
motionProxy.wakeUp()

# Send robot to Pose Init
postureProxy.goToPosture("StandInit", 0.5)

# x - normalized, unitless, velocity along X-axis. +1 and -1
# correspond to the maximum velocity in the forward and
# backward directions, respectively.
# y - normalized, unitless, velocity along Y-axis. +1 and -1
# correspond to the maximum velocity in the left and right
# directions, respectively.
# theta - normalized, unitless, velocity around Z-axis. +1 and
# -1 correspond to the maximum velocity in the
# counterclockwise and clockwise directions, respectively.

x = 1.0
y = 0.0
theta = 0.0
frequency = 1.0
motionProxy.moveToward(x, y, theta, [["Frequency", frequency]])

# Lets make him stop after 3 seconds
time.sleep(3)
motionProxy.stopMove()

# Go to rest position
motionProxy.rest()
```

2.2.3 Audio

NAO uses four directional microphones to detect sounds and equipped with a stereo broadcast system made up of 2 loudspeakers in its ears as shown in the figure 2.3.

NAOs voice recognition and text-to-speech capabilities allow it to communicate in 19 languages.

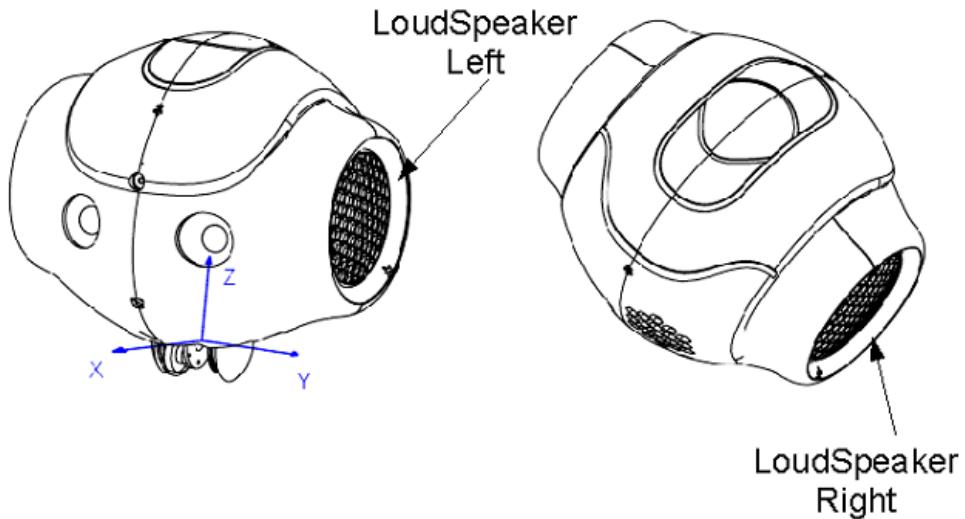


Figure 2.3: NAO Audio

In this thesis, we used Text-To-Speech API of NAO to say some words loud to communicate with the user. Python code 2.2.3 shows how NAO can say words given as strings.

```
# To say the specified string of characters in the specified
language.

from naoqi import ALProxy

robotIP = "nao.local"
PORT = 9559
tts = ALProxy("ALTextToSpeech", robotIP, PORT)

# Set the language to english
tts.setLanguage("English")

# Say the given word
tts.say("Hello World")
```

2.2.4 Vision

Two identical video RGB cameras are located in the forehead of NAO as shown in the figure 2.4. They provide up to 1280x960 resolution at 30 frames per second. NAO contains a set of algorithms for detecting and recognizing faces and shapes.

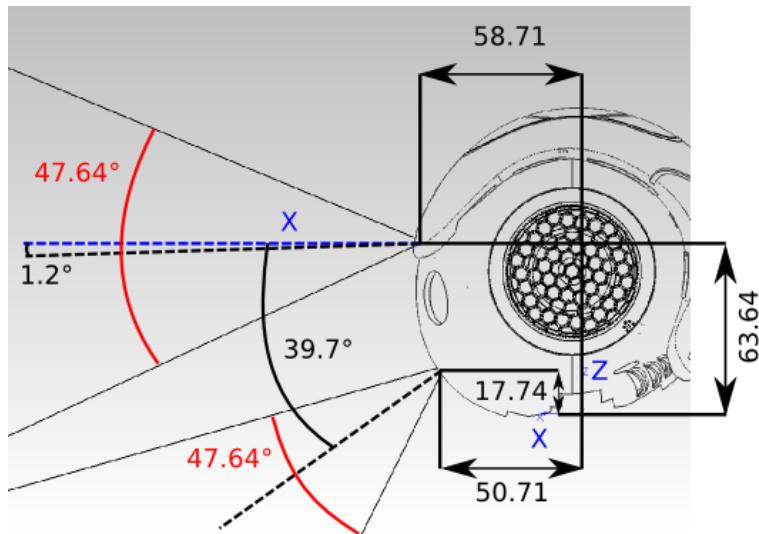


Figure 2.4: NAO Vision

Skeletal points based gesture recognition needs three dimensional data of the human bone joints. However, sensors integrated with NAO could not provide precise three dimensional data for processing heavy algorithms to track human skeletal joints. 3D cameras such as Microsoft Kinect and Asus Xtion are used not only for gaming but also for analyzing 3D data, including algorithms for feature selection, scene analysis, motion tracking, skeletal tracking and gesture recognition [?].

Asus Xtion PRO LIVE uses infrared sensors, adaptive depth detection technology, color image sensing and audio stream to capture a user's real-time image, movement, and voice, making user tracking more precise.



Figure 2.5: Asus Xtion Pro Live

Therefore in this thesis, we attempted to use Asus Xtion PRO LIVE as an external camera that was mounted on the head of NAO as shown in the figure 2.5.



Figure 2.6: Depth Image recorded by depth camera Asus Xtion Pro Live

— this section needs revision —

2.2.5 Computing

NAO is equipped with Intel ATOM 1.6 GHz CPU in the head that runs a 32 bit Gentoo Linux to support Aldebarans proprietary middleware (NAOqi). NAOqi Framework is the programming framework used to program Aldebaran robots. This framework allows homogeneous communication between different modules such as motion, audio, video. NAOqi executable which runs on the robot is a broker. The broker provides lookup services so that any module in the tree or across the network can find any method that has been advertised.

Computational limitations of NAOs CPU hinders us to build a real time gesture recognition based on human skeletal joints. Therefore, we used an off-board computer to execute the gesture recognition program and communicated with NAO via NAOqi proxies.

2.2.6 Networking

Boost ASIO udp unicast socket wifi lan

2.3 Gesture Recognition

Human hand gestures are a means of nonverbal interaction among people. They range from simple actions of using our hand to point at, to the more complex ones that express

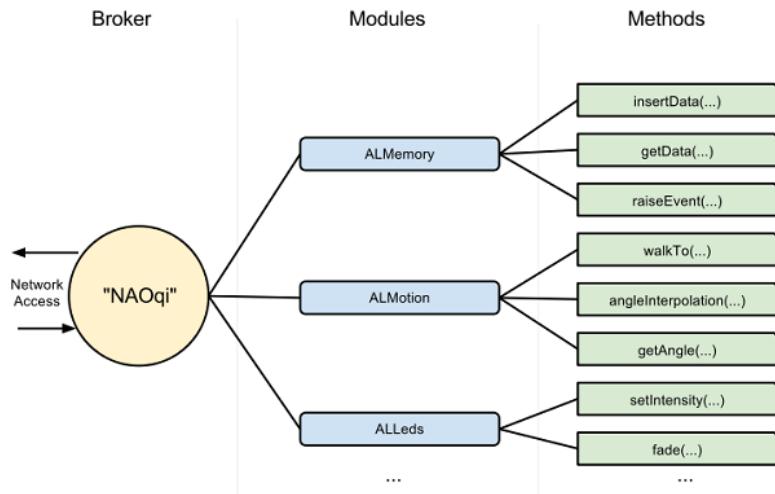


Figure 2.7: NAOqi Proxy

our feelings and allow us to communicate with others. To exploit the use of gestures in HRI, it is necessary to provide the means by which they can be interpreted by robots. The HRI interpretation of gestures requires that dynamic and/or static configurations of the human hand, arm and even other parts of the human body, be measurable by the machine [?].

Initial attempts to recognize hand gestures resulted in electro-mechanical devices that directly measure hand and/or arm joint angles and spatial position using sensors [?]. Glove-based gestural interfaces require the user to wear such a complex device that hinders the ease and naturalness with which the user can interact with the computer controlled environment.

Even though such hand gloves are used in highly specialized domain such as simulation of medical surgery or even in the real surgery, the everyday user will be certainly deterred by such sophisticated interfacing devices. As an active result of the motivated research in HRI, computer vision based techniques were innovated to augment the naturalness of interaction.

2.3.1 Gesture Modeling

Figure 2.8 shows various types of modeling techniques used for Gesture modeling [?]. Selection of an appropriate gesture modeling depends primarily on the intended application. For an application that needs just hand gesture to go up and down or left and right, a very simple model may be sufficient. However, if the purpose is a natural-like interaction, a model has to be sophisticated enough to interpret all the possible gesture.

The following section discusses various gesture modeling techniques which are being used by the existing hand gesture recognition applications.

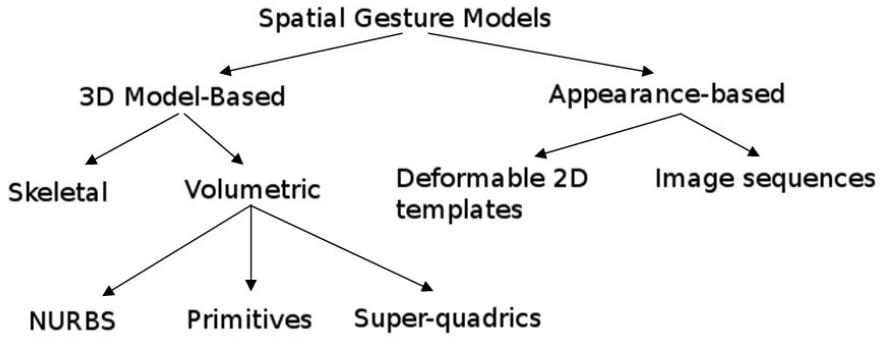


Figure 2.8: Gesture Modeling

Appearance based models don't use the spatial representation of the body, because they derive the parameters directly from the images or videos using a template database. Volumetric approaches have been heavily used in computer animation industry and for computer vision purposes. The models are generally created of complicated 3D surfaces. The drawback of this method is that it is very computational intensive.

Instead of using intensive processing of 3D hand models and dealing with a lot of parameters, one can just use a simplified version that analyses the joint angle parameters along with segment length. This is known as a skeletal representation of the body, where a virtual skeleton of the person is computed and parts of the body are mapped to certain segments [?]. The analysis here is done using the position and orientation of these segments or the relation between each one of them.

In this thesis, we focus on skeletal based modeling, that is faster because the recognition program has to focus only on the significant parts of the body.

2.3.2 Gestural Taxonomy

Several alternative taxonomies have been suggested that deal with psychological aspects of gestures [?]. All hand/arm movements are first classified into two major classes as shown in the figure 2.9.

Manipulative gestures are the ones used to act on objects. For example, moving a chair from one location to another. Manipulative gestures in the context of HRI are mainly developed for medical surgery. Communicative gestures, on the other hand, have purely communicational purpose. In a natural environment they are usually accompanied by speech or spoken as a sign language. In HRI context these gestures are one of

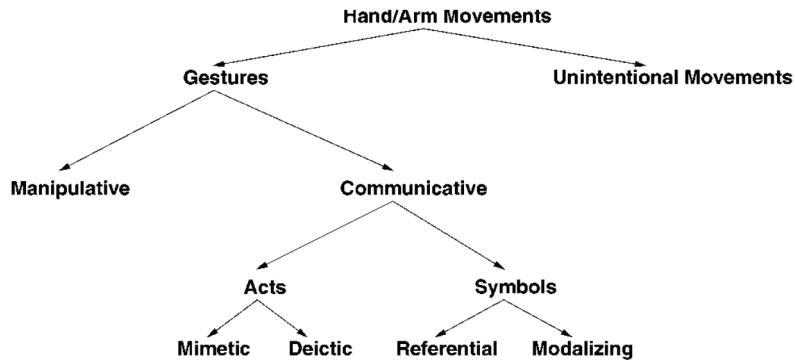


Figure 2.9: Gesture Taxonomy

the commonly used gestures, since they can often be represented by static as well as dynamic hand postures.

In this thesis, we focus on communicative gestures in the form of symbols. They symbolize some referential action. For instance, circular motion of hand may be referred as an alphabet "O" or as an object such as wheel or as a command to turn in a circular motion.

2.3.3 Feature Extraction

Feature extraction stage is concerned with the detection of features which are used for the estimation of parameters of the chosen gestural model. In the detection process it is first necessary to localize the user.

2.3.3.1 OpenNI 2

OpenNI or Open Natural Interaction is a framework by the company PrimeSense and open source software project focused on certifying and improving interoperability of natural user interfaces and organic user interfaces for Natural Interaction (NI) devices, applications that use those devices and middleware that facilitates access and use of such devices. Microsoft Kinect and Asus Xtion are commercially available depth cameras which are compatible with OpenNI.

The OpenNI 2.0 API provides access to PrimeSense compatible depth sensors. It allows an application to initialize a sensor and receive depth, RGB, and IR video streams from the device. OpenNI also provides a uniform interface that third party middleware developers can use to interact with depth sensors. Applications are then able to make use of both the third party middleware, as well as underlying basic depth and video data provided directly by OpenNI.

C++ code 2.3.3.1 shows how a depth camera such as Asus Xtion Pro can be used to retrieve depth information using OpenNI 2 framework.

```
/**  
 * Basic OpenNI setup to read data from depth camera  
 */  
  
  
#include <stdio.h>  
#include <OpenNI.h>  
  
  
#define SAMPLE_READ_WAIT_TIMEOUT 2000  
using namespace openni;  
  
  
int main()  
{  
    OpenNI::initialize();  
    Device depth_camera;  
    depth_camera.open(ANY_DEVICE);  
  
    VideoStream depth_stream;  
    depth_stream.create(depth_camera, SENSOR_DEPTH);  
    depth_stream.start();  
  
    VideoFrameRef depth_frame;  
  
    while (true)  
    {  
        int changedStreamDummy;  
        VideoStream* pStream = &depth_camera;  
        OpenNI::waitForAnyStream(&pStream, 1, &changedStreamDummy,  
            SAMPLE_READ_WAIT_TIMEOUT);  
  
        depth_camera.readFrame(&depth_frame);  
        DepthPixel* pDepth = (DepthPixel*)depth_frame.getData();  
        int middleIndex =  
            (depth_frame.getHeight() + 1) * depth_frame.getWidth() / 2;  
  
        printf("%8d\n", pDepth[middleIndex]);
```

```

    }

depth_camera.stop();
depth_camera.destroy();
depth_camera.close();
OpenNI::shutdown();

return 0;
}

```

2.3.3.2 NiTE 2

PrimeSense's Natural Interaction Technology for End-user is the middleware that perceives the world in 3D, based on the PrimeSensor depth images, and translates these perceptions into meaningful data in the same way that people do. NITE middleware includes computer vision algorithms that enable identifying users and tracking their movements. Figure shows the architecture of NiTE, how it is working together with OpenNI, depth sensors and applications.

Figure 2.10 displays a layered view of producing, acquiring and processing depth data, up to the level of the application that utilizes it to form a natural- interaction based module.

- The lower layer is the PrimeSensor device, which is the physical acquisition layer, resulting in raw sensory data from a stream of depth images.
- The next Cshaped layer is executed on the host PC represents OpenNI. OpenNI provides communication interfaces that interact with both the sensor's driver and the middleware components, which analyze the data from the sensor.
- The sensor data acquisition is a simple acquisition API, enabling the host to operate the sensor. This module is OpenNI compliant interfaces that conforms to OpenNI API standard.
- The NITE Algorithms layer is the computer vision middleware and is also plugged into OpenNI. It processes the depth images produced by the PrimeSensor.
- The NITE Controls layer is an applicative layer that provides application framework for gesture identification and gesture-based UI controls, on top of the data that was processed by NITE Algorithms.

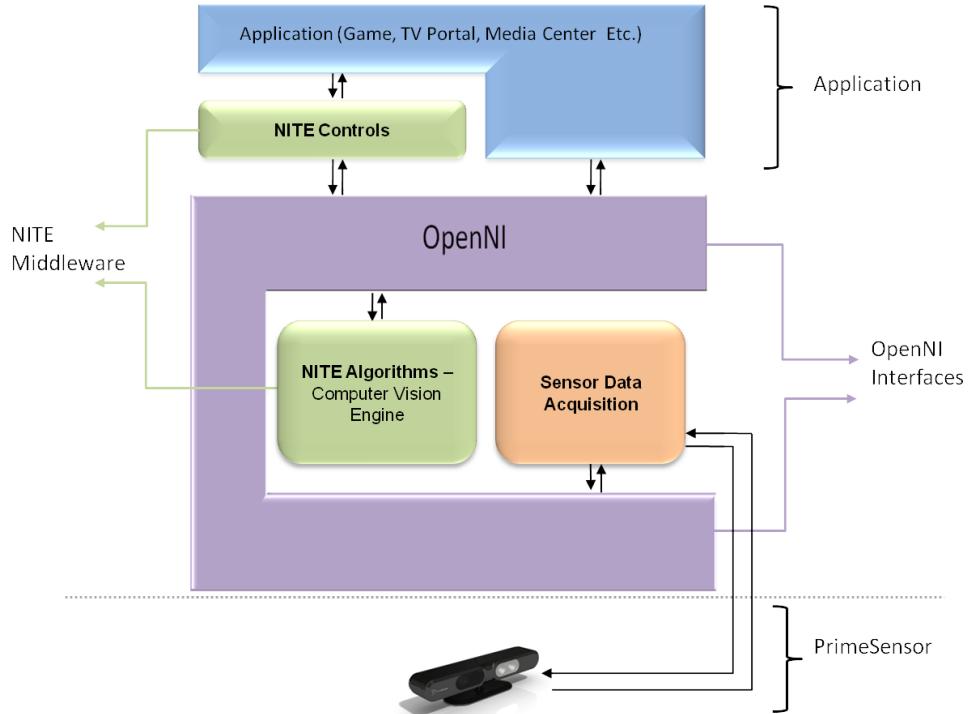


Figure 2.10: NiTE Architecture

2.3.3.3 Algorithm

The lower layer of NiTE middleware that performs the groundwork of processing the stream of raw depth images. This layer utilizes computer vision algorithms to perform the following:

- Scene segmentation is a process in which individual users and objects are separated from the background and tagged accordingly.
- Hand point detection and tracking
- Full body tracking based on the scene segmentation output. Users bodies are tracked to output the current user pose a set of locations of body joints.

NiTE uses machine learning algorithms to recognize anatomical landmarks and pose of human body []. Figure 2.11 shows how NiTE tracks human skeleton from a single input depth image and a per-pixel body part distribution is inferred. Colors indicate the most likely part labels at each pixel and correspond to the joint proposals. Local modes of this signal are estimated to give high-quality proposals for the 3D locations of body joints, even for multiple users.

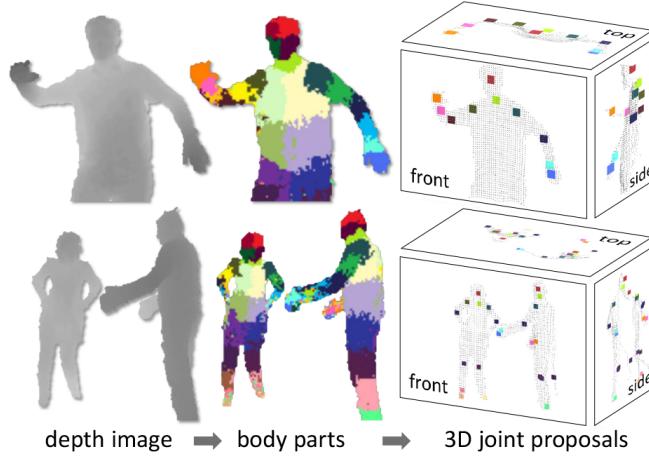


Figure 2.11: NiTE Algorithm to do real-time human pose recognition using depth images

In order to train the system, large collection of synthetic and real representations of human body were recorded and labeled. Each body representations was covered with several localized body part labels as show in the figure 2.12. Some of these parts are defined to directly localize particular skeletal joints of interest, while others fill the gaps or could be used in combination to predict other joints.



Figure 2.12: NiTE Synthetic and real training

Features are located in depth image as shown in the figure 2.13 and labeled. The yellow crosses indicates the pixel x being classified. The red circles indicate the offset pixel. In (a), the two example features give a large depth difference response. In (b), the same two features at new image locations give a much smaller response.

Randomized decision forest is the classification algorithm used by NiTE to predict the probability of a pixel belonging to a body part. Randomized decision trees and forests have proven fast and effective multi-class classifiers for many tasks. Figure 2.14 shows Randomized Decision Forests. A forest is an ensemble of trees. Each tree consists of split nodes (blue) and leaf nodes (green). The red arrows indicate the different paths that might be taken by different trees for a particular input.

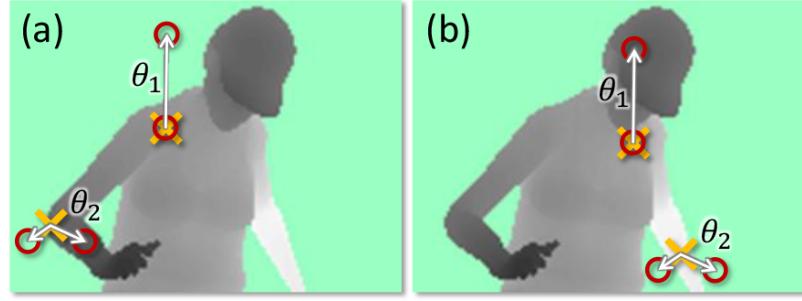


Figure 2.13: Pixels are labeled

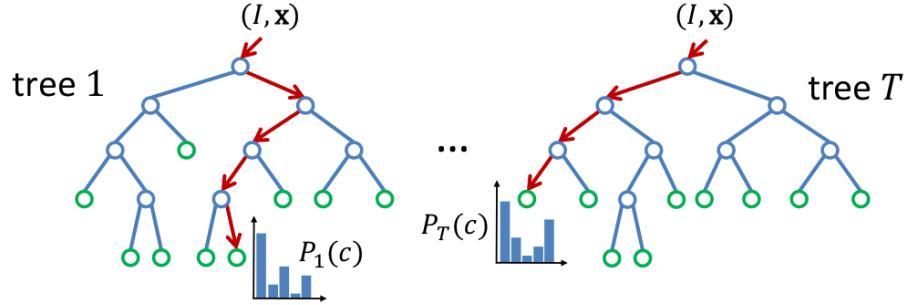


Figure 2.14: Randomized decision forest

To classify pixel x in image I using Randomized decision tree, one starts at the root and repeatedly evaluates equation 2.1, branching left or right according to the comparison to threshold τ . At the leaf node reached in tree t , a learned distribution $P_t(c|I, x)$ over body part labels c is stored. The distributions are averaged together for all trees in the forest to give the final classification.

$$P_t(c|I, x) = \frac{1}{T} \sum_{t=1}^T P_t(c|I, x) \quad (2.1)$$

Each tree is trained on a different set of randomly synthesized images. A random subset of 2000 example pixels from each image is chosen to ensure a roughly even distribution across body parts. Training phase was conducted in distributed manner by training 3 trees from 1 million images on 1000 core cluster.

After predicted the probability of a pixel belonging to a body part, the body parts are recognized and reliable proposals for the positions of 3D skeletal joints are generated. These proposals are the final output of the algorithm and used by a tracking algorithm to self initialize and recover from failure.

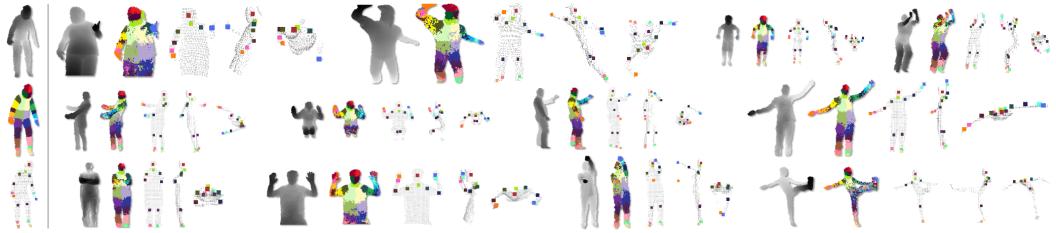


Figure 2.15: Pose Joint Proposal

Figure 2.15 shows example inferences. Synthetic (top row); real (middle); failure modes (bottom). Left column: ground truth for a neutral pose as a reference. In each example we see the depth image, the inferred most likely body part labels, and the joint proposals show as front, right, and top views (overlaid on a depth point cloud). Only the most confident proposal for each joint above a fixed, shared threshold is shown.

Finally NiTE API returns positions and orientations of the skeleton joints as shown in the figure 2.16. As well as it returns the lengths of the body segments such as the distance between returned elbow and shoulder. Joint positions and orientations are given in the real world coordinate system. The origin of the system is at the sensor. +X points to the right of the, +Y points up, and +Z points in the direction of increasing depth.

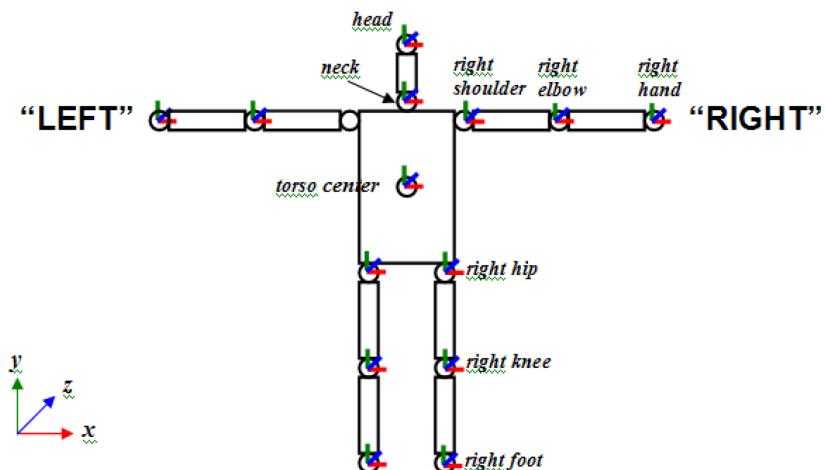


Figure 2.16: NiTE Skeletal Points

Even though NiTE framework can recognize full human body, in this thesis we have used only hand recognition and tracking due to computational limitation of NAO. NiTE provides an interface track only a hand in real time. In order to start tracking a hand, a focus gesture must be gesticulated. There are two supported focus gestures: click and

wave. In the click gesture, you should hold your hand up, push your hand towards the sensor, and then immediately pull your hand back towards you. In the wave gesture, you should hold your hand up and move it several times from left to right and back. Once hand is been found and it will be tracked till the hand leaves the field of view of the camera or hand point is lost due to various factors such as hand was touching another object or closer to another body part. Figure 2.17 shows how hand points are tracked using NiTE and trail of the hand positions in real world coordinates are mapped on to the depth image.



Figure 2.17: NiTE Hand Tracking

Focus gestures of NiTE is can be detected even after initiating the hand tracking. NITE gestures are derived from a stream of hand points which record how a hand moves through space over time. Each hand point is the real-world 3D coordinate of the center of the hand, measured in millimeters. Gesture detectors are sometimes called point listeners (or point controls) since they analyze the points stream looking for a gesture.

NiTE recommends user to follow these suggestions to gain maximum efficiency from its API.

- Try to keep the hand that performs the gesture at a distance from your body.
- Your palm should be open, fingers pointing up, and face the sensor.
- The movement should not be too slow or too fast.
- WAVE movement should consist of at least 5 horizontal movements (left-right or right-left)
- CLICK movement should be long enough (at least 20 cm).

- Make sure CLICK gesture is performed towards the sensor.
- If you have difficulty to gain focus, try to stand closer to the sensor (around 2m), and make sure that your hand is inside the field of view.

Finally, below given example shows how hand tracking can be initiated using a focus gesture.

2.3.4 Gesture Classification and Prediction

Like most other recognitions such as speech recognition and biometrics, the tasks of gesture recognition involve modeling, feature extraction, training, classification and prediction as shown in the figure 2.18. Though the alternatives such as Dynamic Programming (DP) matching algorithms have been attempted, the most successful solutions involves feature-based statistical learning algorithm. Previous sections explained how gesture is modelled and feature is extracted from raw depth images, and the following sections discuss how extracted feature inputs are trained, classified and predicted.

In this thesis, we have chosen a machine learning technique based on Adaptive Naive Bayes Classifier (ANBC) with the help of Gesture Recognition Toolkit. ANBC is an extension to the well-known Naive Bayes, one of the most commonly used supervised learning algorithms that works very well on both basic and more complex recognition problems.

2.3.4.1 Adaptive Naive Bayes Classifier (ANBC)

ANBC is a supervised learning algorithm that can be used to classify any type of N-dimensional signal. It is based on simple probabilistic classifier called Naive Bayes classifier. It fundamentally works by fitting an N-dimensional Gaussian distribution to each class during the training phase. New gestures can then be recognized in the prediction phase by finding the gesture that results in the maximum likelihood value that is calculated by calculating Gaussian distribution for each sample.

ANBC like Naive Bayes classifier makes a number of basic assumptions with input data that all the variables in the data are independent. However, despite these naive assumptions, Naive Bayes Classifiers have proved successful in many real-world classification problems. It has also been shown in a study that the Naive Bayes Classifier not only performs well with completely independent features, but also with functionally dependent features.

ANBC algorithm is based on Bayes' theory and gives the likelihood of event A occurring given the observation of event B. In the equation 2.2, $P(A)$ represents the prior probability of event A occurring and $P(B)$ is a normalizing factor to ensure that all the posterior probabilities sum to 1.

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (2.2)$$

The weighting coefficient adds an important feature for the ANBC algorithm as it enables one general classifier to be trained with multi-dimensional inputs, even if a number of inputs are only relevant for one particular gesture. For example, if it is used to recognize hand gestures, the weighting coefficients would enable the classifier to recognize both left and right hand gestures independently, without the position of the left hand affecting the classification of a right handed gesture. In this case left hand gestures will have weights 1,1,1,0,0,0, right hand gestures will have weights 0,0,0,1,1,1 and both hand gestures will have weights 1,1,1,1,1,1.

Using the weighted Gaussian model, the ANBC algorithm requires $G(3N)$ parameters, assuming that each of the G gestures require specific values for the N -dimensional μ_k , σ_k^2 and ϕ_k vectors. Assuming that ϕ_k is set by the user, μ_k and σ_k^2 values can easily be calculated in a supervised learning scenario by grouping the input training data X into a matrix containing M training examples each with N dimensions, into their corresponding classes. The values for μ and σ^2 of each dimension (n) for each class (k) can then be estimated by computing the mean and variance of the grouped training data for each of the respective classes.

After the Gaussian models have been trained for each of the G classes, an unknown N -dimensional vector x can be classified as one of the G classes using the maximum a posterior probability estimate (MAP). The MAP estimate classifies x as the k th class that results in the maximum a posterior probability given by:

— add equation 14 and 15th —

Using equation (15), an unknown N -dimensional vector x can be classified as one of the G classes from a trained ANBC model. If x actually comes from an unknown distribution that has not been modeled by one of the trained classes (i.e. if it is not any of the gestures in the model) then, unfortunately, it will be incorrectly classified against the k th gesture that gives the maximum log-likelihood value. A rejection threshold, k , must therefore be calculated for each of the G gestures to enable the algorithm to classify any of the G gestures from a continuous stream of data that also contains non-gestural data.

One key element of the Naive Bayes Classifier, is that it can easily be made adaptive. Adding an adaptive online training phase to the common two-phase (training and prediction) provides some significant advantages for the recognition gestures. During the adaptive online training phase the algorithm will not only perform real-time predictions on the continuous stream of input data but also it will also continue to train and refine the models for each gesture. This enables the user to initially train the algorithm with a low number of training examples after and during the adaptive online training phase, the algorithm can continue to train and refine the initial models, creating a more robust model as the number of training examples increases.

ANBC algorithm works well for the classification of static postures and non-temporal pattern recognition. However, The main limitation of the ANBC algorithm is that, because it uses a Gaussian distribution to represent each class, it does not work well when the data you want to classify is not linearly separable. Also when ANBC is working on online training, a small number of incorrectly labelled training examples could create a loose model that becomes less effective at each update step and ultimately lead to a poor performance and accuracy.

2.3.4.2 Gesture Recognition Toolkit (GRT)

Gesture Recognition Toolkit is a cross-platform open-source C++ library designed and developed mainly by Nicholas Gillian at MIT Media Lab to make real-time machine learning and gesture recognitions. Emphasis is placed on ease of use, with a consistent, minimalist design that promotes accessibility while supporting flexibility and customization for advanced users. The toolkit features a broad range of classification and regression algorithms and has extensive support for building real-time systems. This includes algorithms for signal processing, feature extraction and automatic gesture spotting.

In this thesis, we have chose GRT as framework to execute most of the tasks involved in a gesture recognition problem. Figure 2.18 shows that GRT provides the full fledge pipeline to achieve a real-time gesture recognition.

GRT provides an API to reduce the need for repetitive boilerplate code to perform common functionality, such as passing data between algorithms or to preprocess data sets. GRT uses an object-oriented modular architecture and it is built around a set of core modules and a gesture-recognition pipeline. The input to both the modules and pipeline consists of an N-dimensional double-precision vector, making the toolkit flexible to the type of input signal. The algorithms in each module can be used as stand-alone classes; alternatively a gesture-recognition pipeline can be used to chain modules

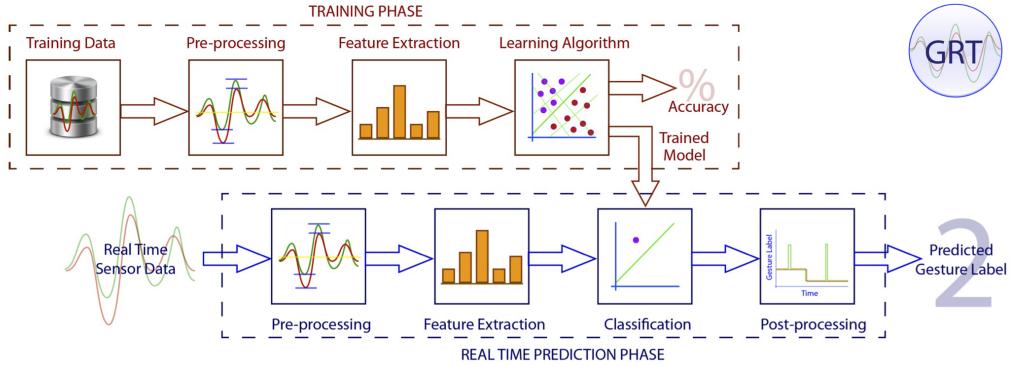


Figure 2.18: Gesture Recognition Pipeline

together to create a more sophisticated gesture-recognition system. Modularity of GRT pipeline offers developers opportunities to work on each stages of gesture recognition independently. Additionally pipeline can be stored and loaded dynamically so that an compiled application can work in many different configurations. Code below shows the basic lines of code needed to build a gesture recognition application.

Accurate labeling of data sets is also critical for building robust machine-learning based systems. The toolkit therefore contains extensive support for recording, labeling and managing supervised and unsupervised data sets for classification, regression and time series analysis. ClassificationData is the data structure used for supervised learning problems and for most of the non-temporal classification algorithms. Code below shows the features of ClassificationData.

GRT allows us to store and load the training data in GRT format or Comma Separated Values (CSV). Since the training data are stored in human readable format, it enables us to add more samples collected separately or remove false samples from the training dataset. Figure shows saved training data in GRT format.

Important part of training is recording positive samples of gestures. Therefore, GRT provides a feature called TrainingDataRecordingTimer that takes recording time and preparation time in milliseconds. Once it is started by calling startRecording(preparationTime, recordTime) method, it waits for given preparation time before it actually starts to store the data. This feature helps the trainer get into the right pose before samples are added to the training data and as well as train all the gestures for the same time duration.

GRT features a broad range of machine-learning algorithms such as AdaBoost, Decision Trees, Dynamic Time Warping (DTW), Hidden Markov Models (HMM), K-Nearest Neighbor (KNN), Linear and Logistic Regression, Adaptive Naive Bayes

(ANBC), Multilayer Perceptrons (MLP), Random Forests and Support Vector Machines (SVM).

Another important feature of GRT is that many of the algorithms are implemented with Null Rejections Thresholds. It means that these algorithms can automatically spot the difference between trained gestures and unintended gestures that can happen when the gesticulator moves the hand in freely. It can be enabled by the method `enableNullRejection(true)` and the range of the null rejection region can be set by this method `setNullRejectionCoeff(double nullRejectionCoeff)` of the classifier. Algorithm such as the Adaptive Naive Bayes Classifier and N-Dimensional Dynamic Time Warping, learn rejection thresholds from the training data, which are then used to automatically recognize valid gestures from a continuous stream of real-time data.

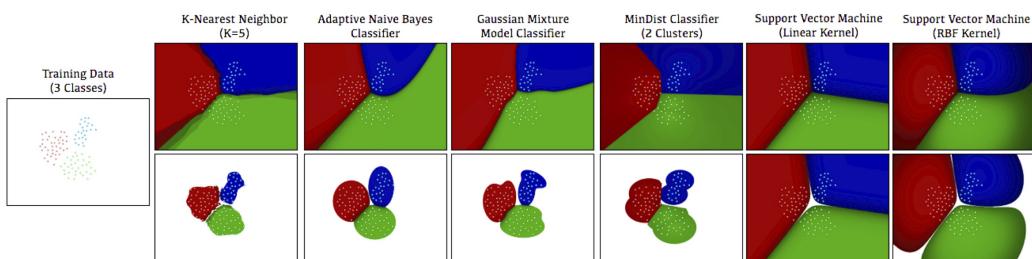


Figure 2.19: GRT Null Rejection

Figure 2.19 above shows that the decision boundaries computed by training six of classification algorithms on an example dataset with 3 classes. After training each classifier, each point in the two-dimensional feature space was colored by the likelihood of the predicted class label (red for class 1, green for class 2, blue for class 3). The top row shows the predictions of each classifier with null rejection disabled. The bottom row shows the predictions of each classifier with null rejection enabled and a null rejection coefficient of 3.0. Rejected points are colored white. Note that both the decision boundaries and null-rejection regions are different for each of the classifiers. This results from the different learning and prediction algorithms used by each classifier.

Real-time classification faces normalization problems when the range of training data differ from prediction input. To solve this problems, there are few solutions such as Z-score Standardization and Feature Scaling. GRT presents a simple solution called as Minimum-Maximum scaling.

Min-Max scaling rescales the range in [0, 1] or [-1, 1]. Selecting the target range depends on the nature of the data. Classifier's `enableScaling(true)` method scales input vector between the default min-max range that is from 0 to 1. The cost of having this

bounded range is that model will end up with smaller standard deviations, which can suppress the effect of outliers. Equation below shows how Min-Max scaling is done.

In many real-world scenarios, the input to a classification algorithm must be pre-processed and have salient features extracted. GRT therefore supports a wide range of pre/post-processing modules such as Moving Average Filter, Class Label Filter and Class Label Change Filter, embedded feature extraction algorithms such as AdaBoost, dimensionality reduction techniques such as Principal Component Analysis (PCA) and unsupervised quantizers such as K-Means Quantizer, Self-Organizing Map Quantizer.

There will not be any need of preprocessing modules in this project since raw data received from depth sensor is processed by NiTE framework. However, post processing modules such as Class Label Filter and Class Label Change Filter may be needed for a reasons that depth sensor samples 30 frames per second, therefore 30 input samples per second are supplied to the classifier for prediction and the output must be triggered once for every gesture.

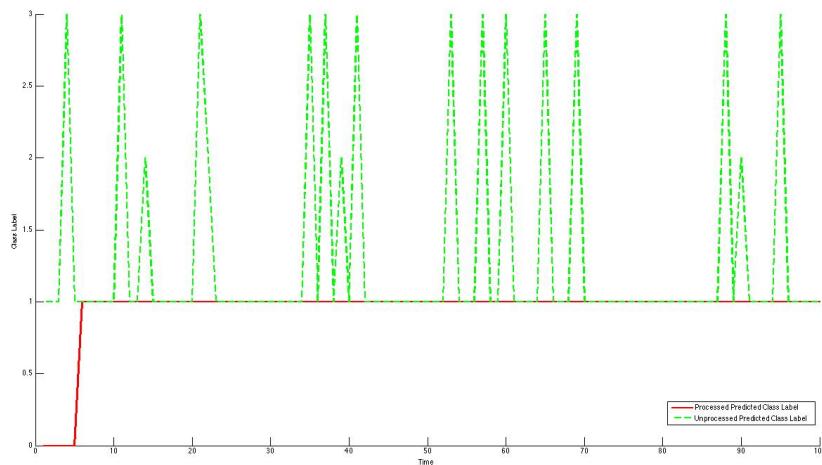


Figure 2.20: GRT Label Filter

Class Label Filter is a useful post-processing module which can remove erroneous or sporadic prediction spikes that may be made by a classifier on a continuous input stream of data. Figure 2.20 that the classifier correctly outputs the predicted class label of 1 for a large majority of the time that a user is performing gesture 1. However, may be due to sensor noise or false samples in the training data, the classifier outputs the class label of 2. In this instance the class label filter can be used to remove these sporadic prediction values, with the output of the class label filter in this instance being 1.

Class Label Filter module is controlled through two parameters: the minimum count value and buffer size value. The minimum count sets the minimum number of label values that must be present in the buffer to be output by the Class Label Filter. The size of the class labels buffer is set by the buffer size parameter. If there is more than one type of class label in the buffer then the class label with the maximum number of instances will be output. If the maximum number of instances for any class label in the buffer is less than the minimum count parameter then the Class Label Filter will output the default null rejection class label of 0.

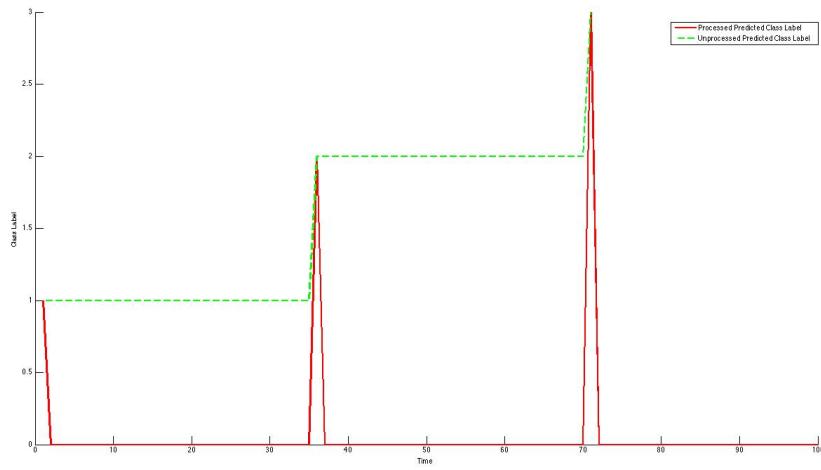


Figure 2.21: GRT Label Change Filter

Class Label Change Filter is one of the useful postprocessing module that triggers when the predicted output of a classifier changes. Figure 2.21 shows that, if the output stream of a classifier was 1,1,1,1,2,2,2,2,3,3, then the output of the filter would be 1,0,0,0,2,0,0,0,3,0. This module is useful to trigger a gesture only once if the user is gesticulating the same gesture for longer time duration.

Figure 2.22 shows GRT-GUI which is an application that provides an easy-to-use graphical interface developed in C++ to setup and configure a gesture recognition pipeline that can be used for classification, regression, or timeseries analysis. Data and control commands are streamed in and out of this application as Open Sound Control (OSC) packets via UDP . Therefore, it acts as a standalone application to record, label, save, load and test the training data and performs a real-time prediction for the incoming data, send output to another application.

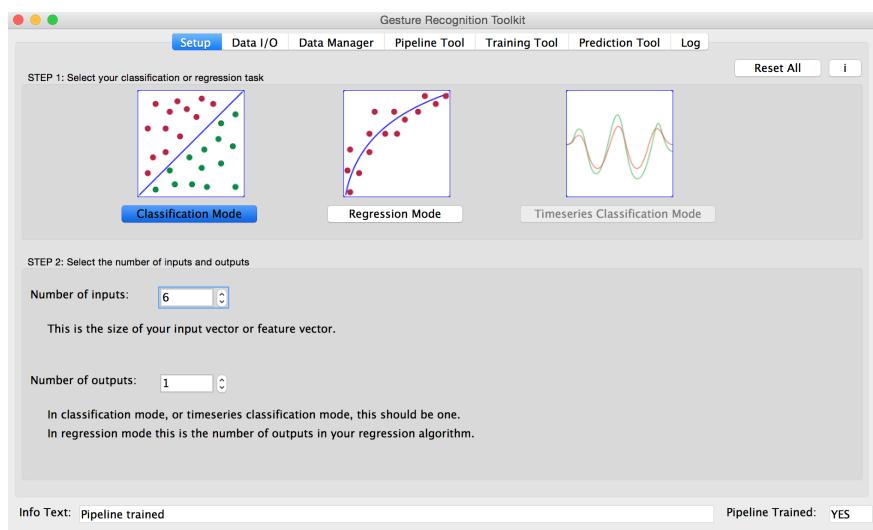


Figure 2.22: GRT GUI

Bibliography

List of Figures

2.1	Construction of NAO	4
2.2	Standing, Sitting and Crouching positions of NAO using ALRobotPosition proxy	4
2.3	NAO Audio	6
2.4	NAO Vision	7
2.5	Asus Xtion Pro Live	7
2.6	Depth Image recorded by depth camera Asus Xtion Pro Live	8
2.7	NAOqi Proxy	9
2.8	Gesture Modeling	10
2.9	Gesture Taxonomy	11
2.10	NiTE Architecture	14
2.11	NiTE Algorithm to do real-time human pose recognition using depth images	15
2.12	NiTE Synthetic and real training	15
2.13	Pixels are labeled	16
2.14	Randomized decision forest	16
2.15	Pose Joint Proposal	17
2.16	NiTE Skeletal Points	17
2.17	NiTE Hand Tracking	18
2.18	Gesture Recognition Pipeline	22
2.19	GRT Null Rejection	23
2.20	GRT Label Filter	24
2.21	GRT Label Change Filter	25
2.22	GRT GUI	26

List of Tables

2.1 NAO V5 specification	3
------------------------------------	---

Abbreviations

AES	Advanced Encryption Standard (Symmetrisches Verschlüsselungsverfahren)
ASCII	American Standard Code for Information Interchange (Computer-Textstandard)
BMP	Windows Bitmap (Grafikformat)
dpi	dots per inch (Punkte pro Zoll; Maß für Auflösung von Bilddateien)
GIF	Graphics Interchange Format (Grafikformat)
HTML	Hypertext Markup Language (Textbasierte Webbeschreibungssprache)
JAP	Java Anon Proxy
JPEG	Joint Photographic Experts Group (Grafikformat)
JPG	Joint Photographic Experts Group (Grafikformat; Kurzform)
LED	Light Emitting Diode (lichtemittierende Diode)
LSB	Least Significant Bit
MD5	Message Digest (Kryptographisches Fingerabdruckverfahren)
MPEG	Moving Picture Experts Group (Video- einschließlich Audiokompression)
MP3	MPEG-1 Audio Layer 3 (Audiokompressionformat)
PACS	Picture Archiving and Communication Systems
PNG	Portable Network Graphics (Grafikformat)
RGB	Rot, Grün, Blau (Farbmodell)
RSA	Rivest, Shamir, Adleman (asymmetrisches Verschlüsselungsverfahren)
SHA1	Security Hash Algorithm (Kryptographisches Fingerabdruckverfahren)
WAV	Waveform Audio Format (Audiokompressionsformat von Microsoft)
YUV	Luminanz Y, Chrominanzwerte U, V (Farbmodell)

Attachments

Hier befinden sich für Interessierte Quelltexte sowie weitere zusätzliche Materialien wie Screenshots oder auch weiterführende Informationen.

.1 HRI Module

.1.1 Gesture Tracker

```
/**  
 * Author: Aravindh Panchadcharam  
 * Email: me@aravindh.info  
 * Date: 28/12/14.  
 * Project: Gesture Recognition for Human-Robot Interaction  
 */  
  
#include <iostream>  
#include <sstream>  
#include <boost/log/trivial.hpp>  
#include <boost/lexical_cast.hpp>  
#include <boost/shared_ptr.hpp>  
#include "NiTE.h"  
#include "gesture_tracker.h"  
  
using boost::property_tree::ptree;  
using boost::property_tree::write_json;  
  
/**  
 * Constructor  
 */
```

```
*/  
gesture_tracker::gesture_tracker(udp_server *server) {  
    server_ = server;  
  
    // Initialize the variables here because compiling on  
    // virtual nao forbides the initialization in the header  
    // file  
    leftHand = 0;  
    rightHand = 0;  
    lastLostHand = 0;  
    handsSize = 0;  
}  
  
/**  
 *  
 * Initializes NiTE and takes available OpenNI device  
 * NiTE::initialize() takes OpenNI deviceId as argument, if  
 * there ar many  
 * Wave and Click Gestures are initiated  
 *  
 */  
void gesture_tracker::init_nite(){  
    niteRc = nite::NiTE::initialize();  
    if (niteRc != nite::STATUS_OK)  
    {  
        BOOST_LOG_TRIVIAL(error) << "NiTE initialization failed"  
        ;  
    }  
    else  
    {  
        BOOST_LOG_TRIVIAL(info) << "NiTE initialized" ;  
    }  
  
    // Check for OpenNI device and start hand tracker.  
    niteRc = handTracker.create();  
    if (niteRc != nite::STATUS_OK)  
    {
```

```
BOOST_LOG_TRIVIAL(error) << "Couldn't create hand
tracker" ;
}

else
{
    BOOST_LOG_TRIVIAL(info) << "Hand tracker created" ;
}

// handTracker.setSmoothingFactor(0.1);
handTracker.startGestureDetection(nite::GESTURE_WAVE);
handTracker.startGestureDetection(nite::GESTURE_CLICK);
BOOST_LOG_TRIVIAL(info) << "Wave your hand to start the
hand tracking" ;
}

/***
*
* Serializes gesture data with position of hand at which
gesture was detected
* Send it to the connected client
*
*/
void gesture_tracker::send_gesture(const nite::GestureData&
gesture){

    std::string gestureJson;
    if(gesture.getType() == 0){
        gestureJson = "{\"GESTURE\":\"WAVE\"}";
    }else if (gesture.getType() == 1){
        gestureJson = "{\"GESTURE\":\"CLICK\"}";
    }

    boost::shared_ptr<std::string> message(new
        std::string(gestureJson));
    server_->send(message);
}
```

```
void gesture_tracker::send_info(std::string info) {  
  
    std::stringstream infoBuffer;  
    infoBuffer << "{\"INFO\": \"" << info << "\"}";  
  
    boost::shared_ptr<std::string> message(new  
        std::string(infoBuffer.str()));  
    server_->send(message);  
}  
  
/**  
 *  
 * Serializes hand data with position and hand id  
 * Send it to the connected client  
 *  
 */  
ptree gesture_tracker::parseToJson(const nite::HandData& hand) {  
  
    // keys and values for json array  
    ptree joint_array, xAxis, yAxis, zAxis;  
  
    xAxis.put("", hand.getPosition().x);  
    yAxis.put("", hand.getPosition().y);  
    zAxis.put("", hand.getPosition().z);  
  
    // Make an array of hand joint for x,y,z axes  
    joint_array.push_back(std::make_pair("", xAxis));  
    joint_array.push_back(std::make_pair("", yAxis));  
    joint_array.push_back(std::make_pair("", zAxis));  
  
    return joint_array;  
}  
  
/**
```

```
*  
* Serializes hand data with position and hand id  
* Send it to the connected client  
*  
*/  
void gesture_tracker::send_hand(const nite::HandData& hand) {  
  
    ptree handJson;  
    std::string handName = getHandName(hand.getId());  
  
    if (!handName.empty())  
    {  
        // Parse it json array and add to object  
        handJson.add_child(handName, parseToJSON(hand));  
  
        // Stringify the ptree  
        std::ostringstream hand_buffer;  
        write_json (hand_buffer, handJson, false);  
        boost::shared_ptr<std::string> message (new std::string(  
            hand_buffer.str()));  
  
        server_->send(message);  
    }  
}  
  
/*  
*  
* Serializes hand data with position and hand id  
* Send it to the connected client  
*  
*/  
void gesture_tracker::send_hand(const nite::HandData& hand1,  
    const nite::HandData& hand2) {  
  
    ptree handJson;  
    std::string handName1 = getHandName(hand1.getId());  
    std::string handName2 = getHandName(hand2.getId());
```

```
if (!handName1.empty() && !handName2.empty())
{
    // Parse it json array and add to object
    handJson.add_child(handName1, parseToJson(hand1));
    handJson.add_child(handName2, parseToJson(hand2));

    // Stringify the ptree
    std::ostringstream hand_buffer;
    write_json (hand_buffer, handJson, false);
    boost::shared_ptr<std::string> message (new std::string(
        hand_buffer.str()));

    server_->send (message);
}

/*
 *
 *
 * Finds the hand name based on given handId
 *
 */
std::string gesture_tracker::getHandName(int handId) {
    std::string handName;

    if (handId == leftHand) {
        handName = "LEFT";
    }
    else if (handId == rightHand) {
        handName = "RIGHT";
    }

    return handName;
}
```

```
/**  
 *  
 * Starts Gesture recognition and Hand tracking based on the  
 * position of Hand found by WAVE gesture  
 * It tracks it till there is a keyboard ESC Hit and stops  
 *  
 */  
void gesture_tracker::track_gestures()  
  
    nite::HandTrackerFrameRef handTrackerFrame;  
  
    while (!utils::wasKeyboardHit())  
    {  
        niteRc = handTracker.readFrame(&handTrackerFrame);  
        if (niteRc != nite::STATUS_OK)  
        {  
            BOOST_LOG_TRIVIAL(error) << "Get next frame failed";  
            continue;  
        }  
  
        const nite::Array<nite::GestureData>& gestures =  
            handTrackerFrame.getGestures();  
        for (int i = 0; i < gestures.getSize(); ++i)  
        {  
            if (gestures[i].isComplete())  
            {  
                send_gesture(gestures[i]);  
                nite::HandId newId;  
  
                // Dont track more than 2 hands.  
                // handTrackerFrame.getHands().getSize() sometime  
                // goes to 3 even though  
                // there are 2 active hands  
                if(gestures[i].getType() == 0){  
                    handTracker.startHandTracking(gestures[i].getCurrentPosition(),  
                        &newId);  
                }  
            }  
        }  
    }  
}
```

```
        }

    }

const nite::Array<nite::HandData>& hands =
    handTrackerFrame.getHands();

// hands.getSize() gives the number of active hands, but
// handId increases
// hands.getSize() = 0 and goes to hands.getSize() = 2,
// when there is a new hand detected
// hands.getSize() = 0, this happens for a fraction of
// second when tracking is troubled
for (int i = 0; i < hands.getSize(); ++i)
{
    // Get Hand data
    const nite::HandData& hand = hands[i];

    // If hand is lost, update the losthandid
    if(!hand.isTracking())
    {
        lastLostHand = hand.getId();
        BOOST_LOG_TRIVIAL(info) <<
            getHandName(hand.getId()) << " Hand with id "
            << hand.getId() << " is Lost";

        // When there is no active hands, reset all the
        // values
        // Last active hand
        if(hands.getSize() == 1){
            leftHand = 0;
            rightHand = 0;
            lastLostHand = 0;
            handsSize = 0;

            send_info("Both hands are lost");
        }
        else
        {
    }
```

```
send_info( getHandName(hand.getId()) + " Hand
           is lost");
}

}

// If new hand is found
if(hand.isNew()) {
    BOOST_LOG_TRIVIAL(info) << "Found new hand with id
        " << hand.getId();

    handsSize++;

    // Check if it is a hand for the first time or
    // second time
    if(handsSize == 1 && lastLostHand == 0) {
        BOOST_LOG_TRIVIAL(debug) << "First hand is
            found";
        rightHand = hand.getId();
    }
    else if (handsSize == 2 && lastLostHand == 0) {
        BOOST_LOG_TRIVIAL(debug) << "Second hand is
            found";
        leftHand = hand.getId();
    }
    // If a hand was lost and a hand is active, then
    // update the appropriate id to left or right hand
    else if(handsSize > 2 && lastLostHand > 0) {
        if(lastLostHand == leftHand) {
            leftHand = hand.getId();
        }else if(lastLostHand == rightHand) {
            rightHand = hand.getId();
        }
    }

    send_info( getHandName(hand.getId()) + " Hand is
        new");
}
```

```
    if(hand.isTouchingFov()) {
        send_info( getHandName(hand.getId()) + " Hand is
                    at FOV");
    }

}

if(hands.getSize() == 2 && hands[0].isTracking() &&
   !hands[0].isNew() && hands[1].isTracking() &&
   !hands[1].isNew()) {
    send_hand(hands[0], hands[1]);
}
else if(hands.getSize() == 1 && hands[0].isTracking() &&
         !hands[0].isNew()) {
    send_hand(hands[0]);
}

}

nite::NITE::shutdown();
}



---



```
 /**
 *
 * Called by the main thread and Boost ioService keeps it in
 * the loop
 *
 */
void gesture_tracker::run() {
 init_nite();
 track_gestures();
}
```


```