

Kinect Wizard Game for Upper Limbs Rehabilitation

Xuzhe Li

Newcastle University

Newcastle upon Tyne

United Kingdom

07413575723

X.Li42@newcastle.ac.uk

ABSTRACT

In this paper, Kinect and gesture recognition will be introduced as well as serious game. Although gesture recognition is simplified by Kinect, that could still be a complex problem when the candidates are complicated gestures, and it also depends on the algorithms that to be implemented. This project aims to create a system to record and recognize upper-limbs gestures and to build a game which can potentially help patients with the rehabilitating process.

Keywords

Human-Computer Interaction, Kinect Sensor, Upper-limb Rehabilitation, Gesture Recognition, Unity, Serious Games

1. INTRODUCTION

The development of computational devices expands the ways of human-computer interaction, and it is not only limited to mice and keyboards nowadays. Vision-based human-computer interaction is one of the more effective ways to get information (body movement, gestures, postures, etc.) from people. Because the technology of motion sensor devices such as Kinect becomes mature and easily accessible, applications using this technology are developed in broad areas such as smart home, simulation, rehabilitation and games. Actually, the research of motion capture and recognition has been doing for decades, and the early work has been done with the help of wearable equipment (i.e. instrumented gloves, body suit, etc.) to gather data. The appearance of motion-sensing input devices allows people to be freed from the equipment, and the research interests have transferred to the area of markerless motion recognition. Motion sensors simplify the process of capturing the motion data and provide the function of tracking the human body movement; however, the data is just the basis of human gesture recognition. The main challenge is analyzing the raw data captured by the sensor, and using proper methods to extract the correct gestures.

Most rehabilitation takes place at hospitals or rehabilitation centers, and that is not always convenient for patients to attend the rehabilitation treatments. In recent years, motion capture technics are used widely in the field of medicine, and it makes home-based rehabilitation possible with the help of proper hardware and well developed software [1]. Traditional rehabilitation is usually supervised by therapists. Traditionally, therapists will teach the patients the exercise they need to do and observe the patients' performance to adjust the progress and methods. With the help of machines, all these work could be handled automatically. Machines can teach the patients by playing exercise videos or using 3D models for demonstration, store the performance data and analyze it to choose the suitable approaches for patients accordingly. Using proper devices and software in rehabilitation could potentially provide a cheaper alternative for people who need treatment.

Kinect (Figure. 1) is used as input device in this project. It is a popular motion input device in the market. A Kinect sensor has four major components:

- A RGB camera which stores three channels of color data in the resolution of 1280 * 960.
- A depth sensor (infrared emitter) which emits infrared light beams and reads the beams that reflected back by objects. It is used for measuring the distance of objects.
- A multi-array microphone which consists of four microphones records the audio and locates the direction of the source.
- A 3-axis accelerometer configured for 2 times gravity to determine the orientation of Kinect.

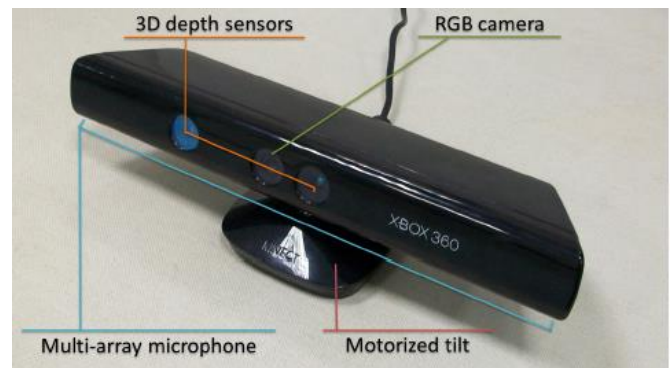


Figure 1. The Microsoft Kinect for Xbox 360 [15]

Except for the data of RGB image and depth image, Microsoft's Kinect Software Development Kits (SDK) also provides Application Programming Interfaces (API) to track the positions of skeleton (figure. 2). The Kinect SDK uses the information from RGB camera and IR sensor to track the 3D coordinates of 20 joints (30 frames per second) [2]. The data of skeleton is not very accurate especially when the joints overlap, but it is enough for the area of rehabilitation which does not have high requirement for precision. With the tracked skeleton joints data provided by Kinect, the remaining task for developers is to use the data to recognize gestures.

Unity3D Personal version is used for the purpose of making a rehabilitation PC game prototype. Unity is a cross-platform engine which supports development for PC, consoles, mobile devices and websites. C# is the programming language that is used for writing the scripts with the help of Unity APIs and some Windows system libraries. Unity does not support .NET framework 4.0, so some of the APIs in the recent Kinect SDK does not work. Therefore, the recognition system is also built in Unity to avoid problems that could possibly occur when

integrating the system into the game engine. It is built on the Microsoft official Unity package [19] since the package provides some examples to show how to access the data from the Kinect sensor in Unity.

The rest of paper is organized as follow: in the section of background, the usage of Kinect, general research on gesture recognition will be reviewed as well as the concept of serious game. Some specific cases of gesture recognition system will be discussed in detail, especially the \$1 recognizer and approaches using DTW algorithm; the design and implementation of the gesture recognition system and the rehabilitation game is proposed and justified in section 3; then, the test results of the recognition system is shown in section 4 along with the analyzing and explanation; at last, the dissertation is summarized in conclusion and discussed how the project can be taken forward.

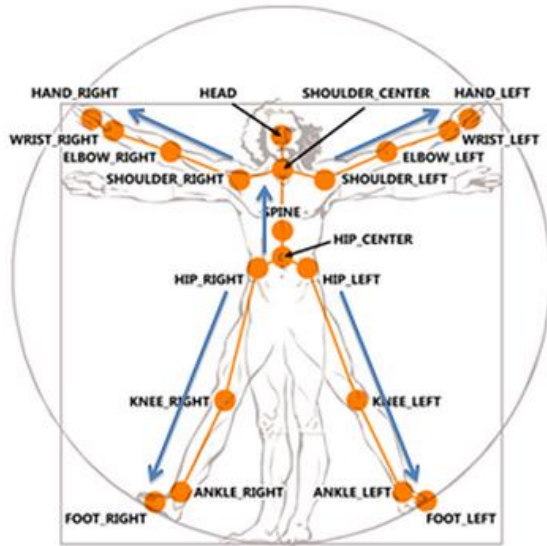


Figure 2. The skeleton joints from the Kinect SDK (Source: <https://msdn.microsoft.com/en-us/library/jj131025.aspx>)

2. BACKGROUND

2.1 The Usage of Kinect

Kinect was originally designed to support the natural interaction in game applications and provide an immersed entertaining environment; however, it is also used in the fields like education, medicine, and modeling.

Technologies are used commonly in classrooms, and it brings a lot of benefits. In the case of Kinect, the kinesthetic features and its gesture-based interaction could encourage students to engage more in the classrooms and improve the interaction of the lessons. Although Kinect cannot stand alone in the classroom without other devices and proper software, it is still a capable tool to enhance teaching and to support learning [23].

Computed Tomography (CT) and Magnetic Resonance Imaging (MRI) are important in clinical practice nowadays. The ways to access CT and MRI data are printing it out slice by slice or using desktop-based UI on computers, but they are not efficient enough for environment like operating rooms (OR) in which surgeons uses the images to navigate the patients' body. To use Kinect in time-crucial situations like OR, it could save time for surgeons by allowing them to access the data with some gestures to simulate

the operations on the images such as selection, rotation, zooming and etc. [24] to improve the efficiency of operation.

Kinect can also be used as depth camera to build 3-D models of indoor environment. It is cheaper and faster comparing with the laser scanning approach, and it is more accurate and robust than vision only approach. Henry et al [25] introduced a framework called RGB-D Mapping to generate dense indoor 3-D maps. Although the capability of Kinect limits the depth precision and the field of view, the framework is still capable of being used in applications like robotic mapping and indoor navigation.

2.2 Gesture Recognition in General

For the perspective of machines, a gesture is a sequence of movement generated by a human and recorded as a series of data. As Turaga et al [3] states, the process of resolving the data can be divided into three steps: 1) extracting the human body information from images; 2) Tracking the related data from each frame; 3) recognizing the gesture. As they point out, most of the raw information from images is irrelevant to identify gestures, such as the background, the color, and etc. How to extract the useful information from images and track it correctly itself is a challenge of recognizing gestures. However, this job has already been done by Kinect as well as tracking the data every frame, so the focus of this project will be how to use that data to recognize gestures.

The representation of gestures can be roughly divided into two categories: the features in 3D space and images in 2D planes. The advantage of using 3D features to represent gestures is that no matter the point of view, 3D features remain the same. However, 3D models with massive amount of parameters usually require a lot of training and computation. Furthermore, matching the gestures with corresponding models is also a difficult problem, and the implementation may require a lot of knowledge not only in the field of computing but also in math. In comparison, the representation of 2D images is relatively simple, but the drawback is that it is only suitable for the situations in which the viewpoints are fixed, so that the sensor can capture the gestures from proper angles to provide correct 2D images [4]. Because of the complex situations in real life, applications using 2D planes to represent gestures will be limited due to the complicated environment. For example, the positions of bodies may change constantly during the process of observation, and images which are captured from the same view point cannot accurately represent the movement.

To recognize the body movement from the gesture data, there are mainly two types of approaches which are used commonly: rule-based approaches and template-matching approaches [5]. The former approach is to use constraints (relationships among joints, absolute positions of joints, etc.) to describe specific gestures. Constraints are stored as data regarding to different gestures, and it is used to check if the movements satisfy the requirement. This approach relies on the parameters and thresholds on joints locations, and users can link gestures to different actions by adjusting the value of thresholds. It is an intuitive way to make the machines recognize human gestures, but in exchange of simplicity, it needs a large amount of assistants from humans by define those gestures. Besides, the process is error-prone, especially when defining those complicated gestures, and people who define the gestures will need relative knowledge and experience in the area.

The template-matching approaches are capable of recognizing gestures more generally, which means humans will not be focused on making rules for specific gestures but on the algorithms of matching the data. It sees the gesture recognition as a

classification problem, and uses the technics of machine learning, data mining, and etc. By using these approaches, machines will be able to record and save the gestures into data sets, and each gesture could be labeled to a classifier to match the movement. When the gestures are presented in front of machines, algorithms evaluate the similarity between the real-time date and each of the gesture patterns in the data set and find out the classifier with the highest similarity. Although this approach offers a better solution for gesture recognition, developers have to implement machine learning algorithms, and maybe to train the machine with a lot of data [6].

2.3 Recognition with Kinect

There is fairly large amount of work has been done in the area of human body movement recognition, and some works are discussed below.

2.2.1 Hands Detection

Kinect does not have much support for hands tracking. The positions of hands cannot always be tracked accurately when the hands overlap with the other parts of human body especially in the situation like folding hands or both hands moving across each other. Because of that, some works are focused on tracking the movements of hands and detection of hands gestures.

Zhou et al [7] use the Kinect SDK to locate the rough position of the hand and add a pair of depth values as thresholds to get the region of the hand. They also use a black wrist belt to help them to get more accurate results by using RANSAC (Random sample consensus) to filter the image data after detecting the black color pixels of the belt. After removing irrelevant information but the shape of hands on the frame images, they represent the data as a curve in time manner which shows the distance between each point on the outline of the hand and the centroid point.

Li [8] uses a similar way to estimate the rough region of the hands by manually setting depth thresholds, but his system can track two hands simultaneously. He uses 2D images on which the hands are projected to do advanced analysis. To be more detailed, he applies the K-means clustering algorithm in order to partition the points into two groups. The distance of centroid points of the two clusters is used to decide if there are two hands or one. If the distance is smaller than a pre-defined minimum value, then there is only one hand, and the two clusters will be merged to represent that hand; otherwise, each cluster will represent one hand. After partitioning the clusters, he uses the Graham Scan Algorithm and a modified Moore-Neighbor Tracing Algorithm to compute and detect the outline of the hands. Above approaches can track the

hands effectively; however, because of the depth thresholds they set in order to detect the hands, the movement users will be limited. If the hands move out of the depth range (like punching, swing hands forward and backward), they cannot be detected.

2.2.2 Rule-based Approach

Flexible Action and Articulated Skeleton Toolkit (FAAST) [11] and Full Body Interaction (FUBI) Framework [12] recognize gestures through detecting sequences of postures as well, but they use simple rule-based approach to recognize postures. They describe gestures using the relations among specific joints as well as the time limit and directions, and then they pre-store them in configuration files. In run-time, the joints' positions are used to compare with the configuration every frame. The speed and directions of joints are also calculated constantly to combine with the constraints of postures to recognize specific gestures. When the joints data meets the requirement of both time constraints and specific postures, the gesture will transit to next phase until a full gesture sequence is completed. As mentioned in the previous section, a lot of effort, experience, and knowledge are required for defining and testing the rules that are set to gestures. Additionally, they can effectively recognize gestures like waving hands, but they are not suitable for the complex gestures which need to be recognized in this project (Figure 6) since many of them are very difficult to be distinguished by only constraints.

2.2.3 Machine Learning and Data Mining

Machine learning and Data mining are proved to be suitable for gesture recognition, and works below using these technics to detect static postures instead of dynamic body movement, and those postures are used as key frames to recognize completed gestures. Patsadu et al [9] proposed an approach using data mining classification method in video stream. They get the data of 3-D vectors for the 20 body joints from Kinect sensor and use various method of data mining classification (BPNN, SVN, decision tree, naive Bayes) to compare the performance among them. An open-source data mining tool named KNIME is used to achieve those classification methods. They tested sit down, stand, and lie down three gestures and found all the methods are effective.

Miranda et al [10] uses a pose descriptor to represent human body poses. The descriptor uses joint angles representation for the nodes of the skeleton. This way of representing the joints is more robust than the original frame skeleton data provided by the Kinect sensor. It provides invariance to sensor orientation and removes the redundant data but the information for the

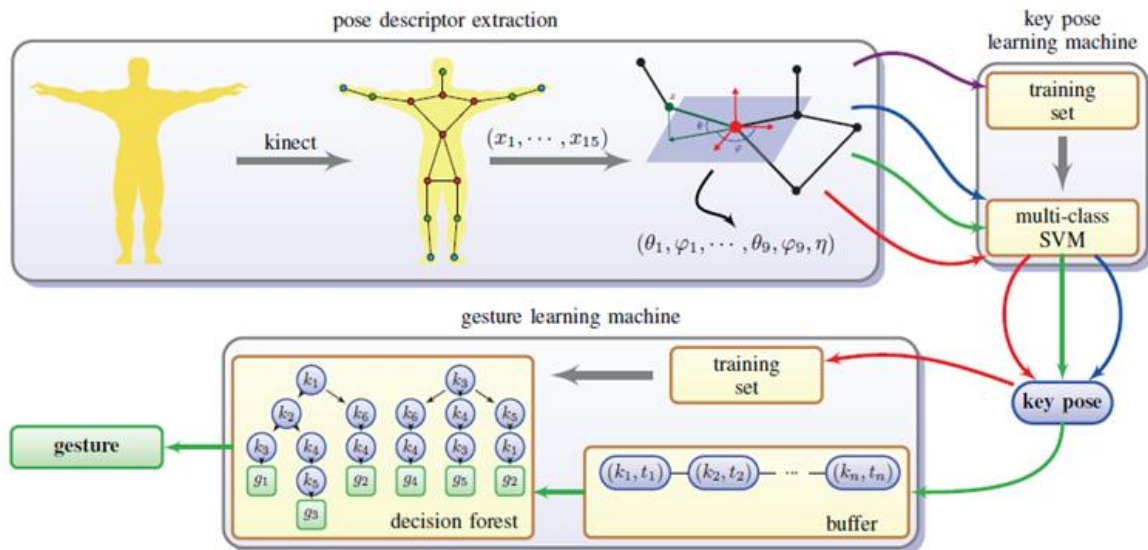


Figure 3. Work flow of the recognition system proposed by Miranda et al

classification of key poses. Training sets are used to build multiple SVM learning machines to identify key poses in real-time. After the learning machine is trained to recognize the key poses, users can define gestures as sequences of key poses, and a decision forest will be built so that the poses are connected to improve the efficiency of searching gestures. It can even take the time/speed constraints into consideration when recognizing gestures.

These approaches are robust and perform very well on recognizing specific gestures, but knowledge of machine learning and data mining is required in order to design and implement a robust system. Besides, the training process may take a long period of time, and some assistants may be required from others.

2.4 Dynamic Time Warping Algorithm

Dynamic time warping (DTW) is a dynamic programming technique, and it was originally applied to the field of speech recognition, and it is one of the more popular algorithms in gesture recognition. The objective of DTW is to compare two sequences, $X := (x_1, x_2, \dots, x_N)$ of length $N \in \mathbb{R}$ and $Y := (y_1, y_2, \dots, y_M)$ of length $M \in \mathbb{R}$ which are not only differ in amplitudes but also in time progression, and to calculate the warping distances between the two.

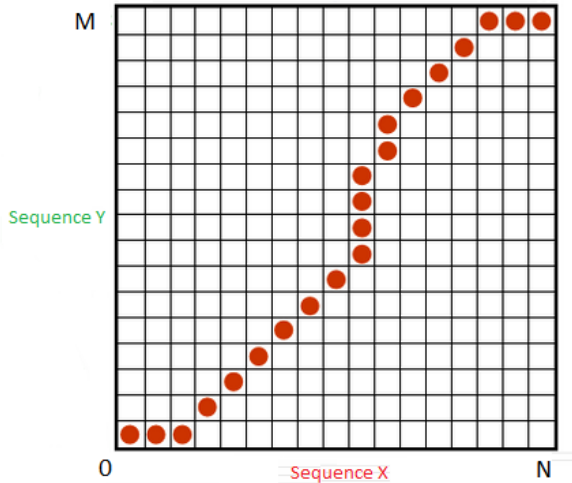


Figure 4. The red points consist of warping path

When the algorithm is applied to gesture recognition, it calculates warping paths and distances between a test sequence and the template sequences. Assuming X is the sequence of the gesture data, and Y is one of the template sequences. In order to get the warping path, a local distance matrix $D \in R_{N \times M}$ which stores the distances of every pairs of points in two sequences will be first constructed. To be more specific, each element of D is the value of distance between $X(n)$ and $Y(m)$. Euclidean distance is typically used (this project also chooses to use Euclidean distance), and Manhattan distance will be used as an alternative in some occasions.

$$D(n, m) = \sqrt{(X(n)_x - Y(m)_x)^2 + (X(n)_y - Y(m)_y)^2}$$

A warping path matrix $W \in R_{N \times M}$ then can be calculated based on the matrix D . The initial condition for constructing the warping matrix is to assign the element $W(0,0)$ with the value of $D(0,0)$. First, the first row and the first column which start from the first element $W(0,0)$ need to be calculated. To calculate the values in

the first row of W , only the horizontal path needs to be considered (using the equation $W(n-1, m) + D(n, m)$). Similarly for the first column, only the vertical extensions will be considered (using the equation $W(n, m-1) + D(n, m)$), and values of the rest elements will be calculated using the following equations:

$$W(n, m) = \min \begin{cases} W(n, m-1) \\ W(n-1, m) \end{cases} + D(n, m)$$

There are four constraints regarding to calculate the warping path:

- **Monotonicity:** the alignment paths always have increasing indices. It means both n and m are either increased or stay the same.
- **Continuity:** the warping paths do not jump in time series, meaning m and n can only be increased by 1 at a time so that important features will all be captured.
- **Boundaries:** the paths start from $(0,0)$ and end at (N, M) , and that could guarantee the whole sequences will be considered.
- **Warping Window:** A good warping path should be close to the diagonal of the matrix.

The length of the optimal warping path equals to the value of $W(N, M)$, and the path can be get by tracking reversely from $W(N, M)$ to $W(0,0)$ and choosing elements in $W(n-1, m-1)$, $W(n, m-1)$ and $W(n-1, m)$ with the smallest value. To get the final result, the warping path length will need to be normalized using a Normalized Factor (NF), and the template has the smallest normalized path length will be the best alignment of the gesture.

To use NF to normalize warping path lengths is necessary because each template will have a different path length, and that would affect the lengths of warping paths and lead to inaccurate results. Normalizing the sequences has been discussed in detail by [21]. They use two different approaches to compare the approaches: reinterpolation all the sequences to have the same length; normalising the returned warping distance. As the result of their experiment shows, the two approaches have similar accuracy. The second approach is used in this project and the NF is applied as the following four conditions:

- No normalization on the distance.
- Normalize by the length of the optimal warping path.
- Normalize by the length of the shorter time series (for each pair of the time series during each DTW computation).
- Normalized by the length of the longer time series.

2.3.1 Approaches using DTW

EasyGR (Easy Gesture Recognition) [6] is a tool using DTW to help to reduce the human effort involved in gesture recognition. In terms of matching templates, it uses a different algorithm but works similarly with the system in this project. Template gestures are first performed by humans and recorded by the machine, and they are sequences of joints data. Then the data will be moved to the center of the detection field and to be normalized. Thresholds then will be generated by comparing the two most different samples of same gestures. Thus, when users perform the same gestures in front of the camera, the system will store the joints' positions in buffer data and use DTW algorithm (or HMM alternatively) to get values for each template, and then the best match will be finally found out.

Another work that uses DTW is done by Su et al [1]. Except for implementing the DTW algorithm, they use fuzzy logic to store the result in a range of 0 to 1 instead of using traditional logic to set the value (true or false) by comparing trajectory and speed. They designed a two-stage fuzzy logic system, and it contains two ANFIS (adaptive neuro-fuzzy inference system) based evaluators, which are implemented by using Matlab, for trajectory and speed, and a Mamdani fuzzy evaluator for the overall performance of the gestures. Calculation of the overall performance takes consideration of the results from all three evaluators. The fuzzy logic system not only recognizes the gesture more precisely but also helps therapists to observe the patients' performance with more detailed information.

Celebi et al [22] proposed an approach called Weighted DTW, which computes weights for every joint in each gesture template. It boosts the efficiency of conventional DTW algorithm since different gestures have different weights on different joints, and weights will be adjusted according to different templates when recognizing the gestures. They use a parametric model which is based on the body joints movement (the placement during the performance of gestures) to compute the weights. After filtering the noise in the gesture data (e.g. shaking, trembling), they apply some equations to the data to calculate the result with the effect of calculated weights. Their work helps to improve the discrimination capability of DTW, and it shows a better performance comparing to the original DTW.

2.5 \$1 Unistroke Recognizer

\$1 Unistroke Recognizer is separated in one sub-section and is discussed in more detail because it is the core algorithm used in this project in terms of template-matching. The \$1 Unistroke Recognizer [13] is a 2-D single-stroke recognizer which can recognize mouse-based gestures using Golden Section Search (originally) and Protractor (a faster solution). In the term of machine learning, \$1 is a geometric template matcher with a Euclidean scoring function. It only needs very few templates comparing with huge datasets in typical machine learning method to perform relative accurate recognition. Besides, the core part of this approach is not complicated; on the contrary, it is easy to implement and deploy.

\$1 uses a simple four-step algorithm:

1. Resample the data of points so that a gesture with a different size of M points is replaced by a gesture in the same shape but with a fixed number of points N .
2. Rotate the gesture to make sure the angle between the line formed by first point and the centroid point and the line $y = 0$ is 0 degree.
3. Scale the gesture to a reference square and translate the gesture so that its centroid point is on $(0, 0)$.
4. Find the optimal angle and get the final score.

The performance of gesture recognition for this approach is fairly well according to the test cases in the article. Comparing with DTW and Rubine classifier, it has the similar accuracy in recognizing gestures with DTW, and has the fastest speed among the three algorithms; however, there are also some limitations. First, the results can only show the closest match in one 2D Euclidean space. It is not enough for recognize complex 3D human gesture. For example, waving the hands for 90 degrees with different directions will all generate a track of quarter-circle; however, to observe the movement from the same perspective will

result different tracks from a line to a quarter-circle. Besides, the original algorithm cannot distinguish gestures which are different in orientations in order to provide tolerance to gesture variation. In human gesture recognition, orientations are the only difference between some gestures and need to be taken into consideration. Especially in rehabilitation treatments, different gestures even with slight difference in those attributes are designed to exercise different parts of the body. At last, there is no time involved in \$1 algorithm; thus, gestures with different speed will be treated as the same.

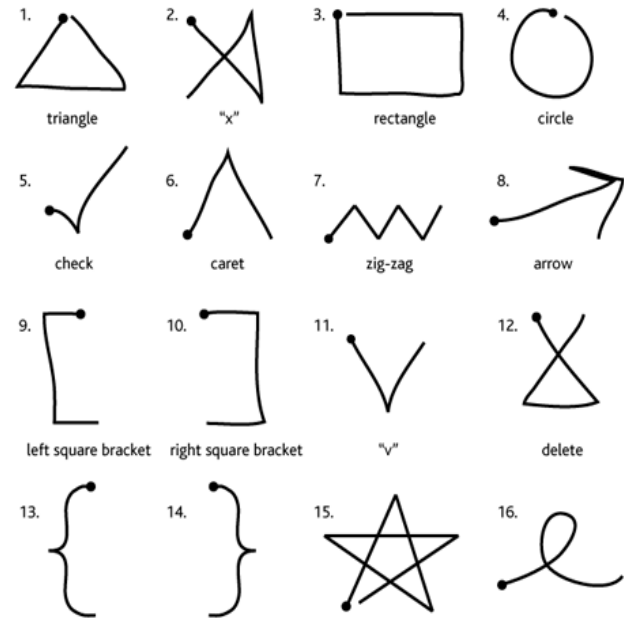


Figure 5. The tested gestures in \$1 [13]

2.6 Serious Games

Using commercial video games as a rehabilitation tool has gained a lot of interest in recent years, and some studies show that games are successful for rehabilitation following stroke and some severe injuries [14]. Serious games are the digital games that are used for not only entertainment, but also for many other purposes such as E-learning, training, rehabilitation and etc. Comparing to the traditional methods in rehabilitation, games have some advantages. For example, games are usually designed highly interactive, engaging and maybe even addictive for players, and this feature could also be applied to patients and help them to involve in the rehabilitation procedure more willingly. Games can create a rich, graphical virtual environment, and that could encourage users to be more motivated, involved and immersed. Virtual environments also provide an approach of safe and cheap training, which can be tailored to particular persons' interests and capabilities. Moreover, the performance and progress can be monitored and analyzed automatically by machines, and the data could be saved and sent back to therapists or be uploaded to remote servers to help therapists observe the treatment and give feedback more accurately.

VR and video-capture systems could make the process even more immersed; however, the high cost of the equipment and difficulties of setting up the system prohibit that from being popular for home rehabilitation. On the other hand, game consoles and user movement capture devices are widely used in serious games. Game consoles and motion capture devices are relatively cheaper and easier to be accessed. Also, the techniques of game consoles are relatively mature, which lower the technical

requirements for the system to be used. There are fairly amount of researchers reported on the usage of different consoles (e.g. Kinect to XBOX, EyeToy to PS2, and Wii). With the proper usage of user movement capture devices, players can interact with the game effectively using body movement [20].

There are some limitations on developing rehabilitation games on consoles. Commercial console games are usually developed by professional game studios whose target users are general players since the most profit of the game industry comes from them. Individual developer may be more willing to develop rehabilitation games; however, they are usually not authorized to do so, and they may also lack of knowledge and experiences on developing console games. Fortunately, companies gradually lower the requirements of developing console games for developers. Microsoft authorizes members of Xbox Creator's Club with only \$99 per year, while Nintendo has a similar platform called WillWare. Besides, the consoles are getting more powerful, and the technical difficulties of developing console games reduce. Moreover, devices like Kinect also works on PCs, and all of these positive conditions open up possibilities for rehabilitation games to be developed and deployed on consoles and PCs [20].

There are some commercial games which are designed for healthy life, such as Wii sports, Tai-Chi and Yoga-based games for Nintendo and Sony consoles. They may not be specially designed for rehabilitation, but in fact they could be suitable and indeed be experimented in some rehabilitation projects. However, as Burk et al [20] pointed out, the paces of those games are still too fast for patients who have just suffered diseases like stroke, and special designed games just for rehabilitation are still necessary. To design games which are suitable for rehabilitation will require a lot of related knowledge as well as experience. They divided the design into two principles: meaningful play and challenges.

Meaningful play is a design pattern about how players interact with the game. The players should be easily noticed on how their actions affect the game play, and feedback is the core part of maintaining meaningful play in games. Feedback helps users to understand how their choices and behaviors affect the game and to measure their progress in achieving goals. The feedback is usually presented to players by visual effect and sound, but it can also exist with other different forms such as numerical scores, health bars, camera movement, and etc. The feedback should occur in all the occasions when human computer interaction happens during the game play, like monsters getting damaged and causing damage, players acquiring new skills and leveling up.

Players usually expect less challenge at the start of a game and desire the challenge increases with the development of the game as they are getting more familiarized with the game and more skilled. It is common for games to have increased difficulties, but achieving that in rehabilitation games would be more difficult. Patients who play the rehabilitation games may vary significantly in experience of playing video games. Some of them maybe play games a lot, while the other may have not played video games. A possible solution for this problem could be adjusting the difficulty of the game according to the collected data from medical organizations and the progress in the game. The pace of rehabilitation games is also an important feature related to the difficulty. Due to the inconvenience of movement for patients, the game speed should be slower, while the movement of player-controlled characters should be amplified.

3. DESIGN AND IMPLEMENTATION

In this article, an easy to implemented approach which is inspired by \$1 Unistroke Recognizer is proposed, and DTW is also implemented to solve the problem of 3D recognition and to be compared with \$1. It combines both rule-based approach and template-matching approach and uses tracked position of hands and other relative joints as data source. The purpose of building this system is to recognize a series of gestures for upper-limb rehabilitation provided by the medical school. It also considers the possibility of the future expansion of the gesture set. Therefore, it will be aimed to develop a system which can generally recognize arbitrary gestures without any particular patterns.

A wizard game which integrates the gesture recognition system is built using Unity3D to help patients who need upper-limbs rehabilitation treatment to recover. By playing a game while doing the rehabilitation exercise, the treatment process could potentially be more effective since patients are more entertained and involved, and that could lead to a better recovery.

3.1 Recognition System

The original idea was using \$1 recognizer as the recognition algorithm since it can recognize many different kinds of shapes which might be generated by the hands movement, see Figure 6. In order to adapt the mouse-based recognition system to a body movement recognition system, the first step is to change the way of input data. The movement of the mouse was replaced by the movement of the hands, and the points that generated by hands movement will be captured by Kinect sensor and be used for recognition (other joints are also stored for checking constraints and future use). The joints data that Kinect output is 3D vector3 in which the value of X and Y are ranged from -1 to 1 and Z is ranged from 0 to 8 (the distances between joints and the sensor in meters).

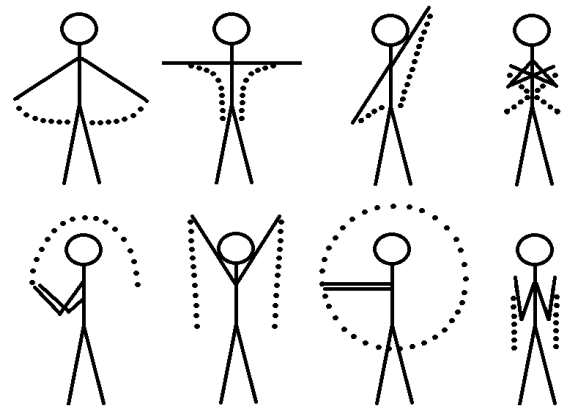


Figure 6. Part of the gestures need to be recognized

The challenge of using hands to generate the shapes is that there is no button that can be clicked or released like a mouse to inform the machine to start and to stop recording data. This is one main difficulty in design of an effective recognition system while another difficulty is that same gestures may be varying between individuals [14]. The second difficulty will not be a challenge in this project since the algorithm will normalize the gesture data, so as long as the gestures performed by different people are the same, the data will be similar enough to be recognized as the same gestures. The first problem still needs to be addressed by finding alternative for the buttons.

Kinect V2 has added two joints (one for the thumb and another one for the rest of the fingers) on the hands so that it is possible to detect grip and release, and that could be a potential replacement for the mouse buttons (click and release). However, the hardware that is used in this project is Kinect V1, and it does not have the extra joints on the hands; therefore, an alternative approach was used. It was to store the hands' positions from the latest M frames in a list with a fixed size in which the oldest position will be removed to free space for the newest position if the list is full. The elements in the list will be used as the gesture data after the hands holding the position (within an acceptable range of movement) for a certain period of time T . The latest N elements in this list are constantly being checked if the size of the list is larger than N to decide if the movement is stopped (the distances between the last point of the list and the rest $N-1$ points are smaller than a predefined value). However, it was found not very efficient because to fill the list with valid data for the candidate gestures is hardly to be achieved since there will be a lot of noise.

The solution has been modified later to avoid the problem of noise as the project goes on. Instead of only giving the system a signal when the gesture finished, the system will also be waiting for a signal to start recording. After the system decides to start recording, the list will be cleared, and the gesture data will be pushed into the list without the existence of noise. The rest of the procedure is the same as it was before with a slight change. When the list is full, the recording will be finished instead of replacing old positions with new positions, and the latest N positions will be removed after finishing recording because they are used for checking the end signal and should be treated as noise.

```
<?xml version="1.0" encoding="utf-8" standalone="yes"?>
<Gesture GesName="hdel2" NumPts="168" Date="Monday, June 08, 2015" TimeOfDay="11:24:56 AM">
  <Constraints Value="1" />
  <Constraints Value="1" />
  <LeftHandPoints X="-0.204528779" Y="-0.524480164" Z="1.79457378" />
  <LeftHandPoints X="-0.1827106" Y="-0.5750704" Z="1.779556" />
  <LeftHandPoints X="-0.168711081" Y="-0.6075141" Z="1.76940835" />
  <LeftHandPoints X="-0.163902923" Y="-0.61880666" Z="1.7660048" />
  <LeftHandPoints X="-0.16085434" Y="-0.625983357" Z="1.76394081" />
  <LeftHandPoints X="-0.159310535" Y="-0.629628539" Z="1.762684" />
  <LeftHandPoints X="-0.1594857" Y="-0.629307" Z="1.76245439" />
  <LeftHandPoints X="-0.160577863" Y="-0.626501441" Z="1.76285732" />
  <LeftHandPoints X="-0.18010664" Y="-0.5601836" Z="1.78225" />
  <LeftHandPoints X="-0.188500792" Y="-0.530782342" Z="1.79133916" />
  <LeftHandPoints X="-0.1926348" Y="-0.516205132" Z="1.79576135" />
  <LeftHandPoints X="-0.195938528" Y="-0.5042188" Z="1.79569137" />
  <LeftHandPoints X="-0.198033571" Y="-0.4965949" Z="1.79530108" />
  <LeftHandPoints X="-0.199361816" Y="-0.4915709" Z="1.79427576" />
  <LeftHandPoints X="-0.200055689" Y="-0.4888119" Z="1.79248464" />
  <LeftHandPoints X="-0.200266153" Y="-0.487824619" Z="1.79031134" />
  <LeftHandPoints X="-0.200881734" Y="-0.487940282" Z="1.78747332" />
  <LeftHandPoints X="-0.201199442" Y="-0.488758177" Z="1.78483415" />
  <LeftHandPoints X="-0.201076478" Y="-0.491325" Z="1.78347874" />
  <LeftHandPoints X="-0.200633064" Y="-0.4933094" Z="1.78205991" />
  <LeftHandPoints X="-0.200272828" Y="-0.495623916" Z="1.78120089" />
  <LeftHandPoints X="-0.1996652" Y="-0.4970827" Z="1.78010664" />
  <LeftHandPoints X="-0.198965738" Y="-0.498746773" Z="1.77904089" />
```

Figure 7. Part of an XML file of template data

There are some drawbacks using this approach. Depending on the lengths of gestures, users may need to adjust the speed of performing a gesture. For example, a gesture which will generate more than $M - N$ positions with the regular speed needs to be performed faster; otherwise, there will not be enough space in the list to store the data. It might be acceptable for the system, but that could lead to a different final trajectory, and the gesture may not be recognized correctly depending on how many positions are missing. Another drawback is that users have to hold their moves both before start to record and after finish performing the gestures for a short while. If extra moves occur before the machine is informed by the signal, that one gesture will be failed to be recognized and needs to be performed again.

The data in the list then will be passed to the recognition system, and the system will process the data so that it can be applied to compare with the templates. The original \$1 algorithm also has a

similar process, and it can be used in gesture recognition with minor changes. First, the position list will be resampled to a fixed size N while keeping the trajectory unchanged. To resample, the total length of the path needs to be calculated first, then the length will be divided by $N-1$ to get the average length I of each increment. To calculate the new list for resampled date, it will start from the first point of the old list and add the distance between the current point and the next point to a variable D . A new position will be interpolated after the value of D exceeds I , and the interpolated point will be pushed to the new list as well as be inserted to the old list to be the start of calculating the next new position.

The second step of processing the date in \$1 is to rotate the data so that the angle between the line formed by first point of the data and the centroid point and the line $y = 0$ is 0 degree. However, it is not suitable for gesture recognition because the movement of hands for many gestures will generate similar tracks but in different rotation, and those gestures with different rotation cannot be distinguished if the data is rotated. Therefore, the step of rotating the points is abandoned in this system, and the rotation of each gesture will be left as-is. Consequently, different rotations will be distinguished, but the gestures need to be performed more restrictedly in order to be recognized (i.e. the starting point and end point should be close to where they desire to be).

The original algorithm then scales the gestures to the referenced square so that they all have the same bounding box, and translates the centroid points of the gestures to the point $(0, 0)$. This step is useful for normalizing the gestures data because the same gestures cannot be guaranteed to be performed at the same positions in the Kinect space, and the same gesture will be performed by people with different heights.

Two additional steps are needed to adapt \$1 to gesture recognition. To solve the problem of different gestures have the similar trajectories but different start positions and end positions, a step of checking constraints of the gestures is added before matching templates and calculating the scores. The constraints which are added to gestures are relative directions between hands and head, hands and elbows and etc. the constraints will be checked according to the gesture data, and only the templates which have the same constraints with the gestures will be considered necessary to be processed to the next stage. Moreover, the efficiency of the algorithm will be improved by adding constraints to gestures because only those gestures which don't pass the constraints checking will be filtered.

As mentioned by Wobbrock et al [13], the original system struggles to get an effective score for 1-D trajectory such as a line. Therefore, for dealing with gestures which will form lines, an extra step is added. The distribution of the positions will be checked to decide if it is closed to a line or not, and the step of template matching will be skipped if a gesture is decided to be a line. The results will be based on the angles between $y = 0$ and the lines formed by the points. To decide the distribution of points, correlation coefficient which is commonly used to measure the degree of linear dependence between two variables is calculated using the gesture data. The following formula is used to calculate the coefficient:

$$r_{xy} = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \sqrt{n \sum y_i^2 - (\sum y_i)^2}}$$

As the formula shows, if the distribution of the point is on the line $y = 0$ or $x = 0$, the denominator will be 0, and that will make the formula invalid. Hence, before the coefficient is calculated, all the points will be rotated until the angle between the line formed by the first point and the last point and the line $y = 0$ is 45 degree. This step makes sure the coefficient will be a float value ranged from 0 to 1. The more the coefficient value is closed to 1, the more the distribution is close to a line [17], and a minimum valued (0.98 is found suitable) is set as the threshold for deciding lines.

Gestures which are not lines will go through the Golden Section Search algorithm [18] to compare the list of hands' positions with the points in each template to get scores. To calculate the score, first, the average distance between corresponding points of candidate C and a stored template T will be get using the formula:

$$d_i = \frac{\sum_{k=1}^N \sqrt{(C[k]_x - T_i[k]_x)^2 + (C[k]_y - T_i[k]_y)^2}}{N}$$

Then the distance d_i and the length of a side of the reference square (size) will be used in the equation below to get a final score which is range from 0 to 1 using the equation:

$$Score = 1 - \frac{d_i^*}{1/2 \sqrt{size^2 + size^2}}$$

To sum up, the workflow of the \$1 recognizer in gesture recognition is listed below:

- Recording the gestures that need to be recognized as templates and saving them as in XML files which contain the positions of both hands in time manner, basic information of the gestures (e.g. name, time, length, etc.) and the constraints of the gestures.
- Loading the templates from files, packaging them and storing in a list, so that they are prepared to be compared with the real-time data of hands' positions.
- Capturing and storing the latest hands' positions into a list with fixed size. Packaging the data in the same way as the templates, and passing the processed data to the algorithm.
- Comparing the constraints of the gesture data and the templates, filtering the templates which don't meet the same requirement of the constraints.
- Calculating the correlation coefficient to decide if it is a path is a line.
- If the data is a line, calculating the angle of the line with the line $y = 0$ to decide which gesture it matches. Otherwise, calculating the score using Golden Section Search and looking for the highest score among all the templates.
- In the end, either the system finds a template which satisfies the minimum requirement for recognition, or the gesture cannot be fine a template which matches enough. After completion of the algorithm, the list that stores the real-time hands' positions will be cleared to prepare for a new gesture.

Since the \$1 is a 2-D algorithm, only x and y of the hands' positions are used for recognition, which means the paths projected on the x-y plane will be used to matching. However,

human gestures are performed in the 3D space, and not all of them can be recognized only by checking the trajectories on the x-y plane. Because how \$1 does to match the templates involves the usage of angles to find the best alignment, it needs a proper understanding and effort to adapt the algorithm to support 3D positions. It might be achieved by using matrices to rotate in 3D space and to use quaternions to replace the angles. The thought of using the 2D algorithm to match 3D human gestures alternatively was to store the position data in three copies, one for the x-y plane, one for the x-z plane and the one for the z-y plane. These three copies are three paths of the same gestures that are projected to those three planes. They are all applied to the algorithm, and the final results will be affected by the scores calculated from all three planes according to the weights. This idea is not successful after being experimented, and there are a couple of reasons why it is not a suitable solution. First, gestures will form meaningless trajectories in one or more plane which are useless for recognition. Second, a proper solution to assign the weights has yet been found.

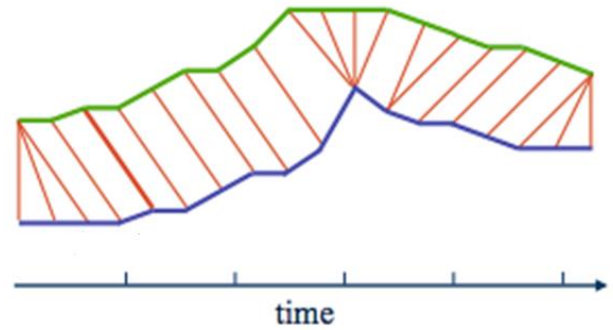


Figure 8. Comparing two sequences using DTW

The progress of implementing DTW is much simpler than implementing \$1 because of the gained knowledge of gesture recognition and the completion of the recognition system framework. For implementing the DTW, what need to do are to understand the algorithm as well as implement it and to implement a procedure of normalizing the data for DTW matching. The process of normalizing the data for \$1 is found also suitable for DTW, but DTW doesn't need the gesture and the templates have the sample amount of positions, so resample is not necessary for the gesture data, and resample is only applied to the templates when they are loaded from the files to make sure there will be a fair comparison.

The classic DTW algorithm is implemented in this project, and its working flow is the same as what is described in the previous section. The processed gesture data and each of the templates are used to calculate the distance matrices first, then the distance matrices will be used to get the warping path matrices. Because one dimensional path can be handled by DTW as long as the data is normalized correctly, the algorithm will be applied after the comparison of constraints, so the step calculating the correlation is not necessary when using DTW. As mentioned above, the way of normalizing data for \$1 cannot also be used in DTW, but this is not a perfect solution, and there is no a better way yet been found in this project. Although it works fine for most cases, it has problems in some tricky situations which will be discussed in the result section. Different approaches of normalizing the data have been experimented, but none of them could solve the problem effectively.

After getting the optimal distances for all the templates, they will be normalized using the equation which is mentioned in the last section, and all the normalized factors will be applied to the equation. The smallest value among all the templates' results will be used as the final result, and it will be converted to the range of 0 to 1 using a similar equation that is used in §1.

Because of the nature of the DTW algorithm, it could be easily adapted to support 3D gesture recognition. To calculate the warping path in 3D space is almost the same as it is in 2D space except for the equations that are used to calculate distances. By adding functions to support the calculation in 3D space, a 3D version of DTW is also achieved, and it can solve the problem of recognizing recognition using 3D trajectories and be used to compare the efficiency with 2D gesture recognition.

3.2 Rehabilitation game

The idea of the rehabilitation game is inspired by a game named Clicker Heroes [16]. The rehabilitation game will be presented in the first person view, and it only has a single scene which has the player, the monsters (bosses) and the terrain. The basic concept of the game is that players cast spell by performing gestures to beat monsters and use the diamonds dropped by the monster to upgrade spells or unlock new spells to defeat stronger enemies. The values in the game can be designed accordingly to the rehabilitation treatment for different patients. For example, only on spell can be used in one level, and the monsters' HP and the damage of the spell could lead to the result that a gesture needs to be performed a certain amount of times to win the level. When the player advances to a more difficult level, a more difficult gesture will be unlock to cast more powerful spells. Through this mechanism, the game could lead players to finish a treatment procedure with gradually added difficulties. Considering that to do the same excises repeatedly will be necessary in rehabilitation, the diamonds can also be used to upgrade low tier spells to make them more powerful than higher tier spells in later game, so that players will still choose to perform the easy gestures at some points of the game.



Figure 9. The game UI

In order to make the game more challenging, monsters are set to attack players once in a certain period of time. The health and attack of the monsters will be increased as the game goes on. The player will gain experience after killing monsters, and the health of the player will be increased and restored by leveling up. The health can also be restored by spending diamonds to buy portions. If the player is killed by the monster, the penalty will be minimum (e.g. spending some diamonds to revive) in order to avoid upsetting the players.

After defeating a certain amount of monsters in one level, the player and the camera will move around the map to reach the

position for the next level and to fight different monsters. Diamonds that dropped after the death of monsters can be collected by the player using the hand cursor on the screen, or they can be left on the ground and be collected automatically as the game moves to the next level. The game is designed to be endless, but an ultimate boss could be set in the game to be the goal that patients want to achieve, and which also means the end of the treatment process.

The movement of hand cursor is controlled by the position of user's right hand. The hand positions collected by the Kinect will be converted into related positions in screen space, and the Lerp function is used to move it smoothly. Diamonds that are dropped by monsters will be collected when the screen positions of hand cursor and diamonds overlap. The diamond counter which consists of Unity GUITexture and GUIText will reflect the change of the amount of diamonds when they are collected or consumed.

The spells in this game are made using the Particle System in Unity3D. Each spell prefab has a Particle System component to take care of the rendering part of the game object, a rigid body to simulate the physical attributes and a collider for collision detection. The spells are connected to the gestures in the configuration file which stores all the spells, their information, and related gesture names.

```
<?xml version="1.0" encoding="utf-8" standalone="yes"?>
<Spells SpellName="8" Date="Monday, June 08, 2015" TimeOfDay="11:24:27 AM">
  <Spell Name="SomeSpell11" Level="43" Damage="10000" Attribute="0" Gesture="hdel1" Number="0" Lock="0" />
  <Spell Name="SomeSpell12" Level="13" Damage="81920" Attribute="1" Gesture="hD Pool Entry Dive" Number="1" Lock="0" />
  <Spell Name="SomeSpell13" Level="9" Damage="10240" Attribute="2" Gesture="hdel10" Number="2" Lock="0" />
  <Spell Name="SomeSpell14" Level="4" Damage="800" Attribute="3" Gesture="hdel3" Number="3" Lock="0" />
  <Spell Name="SomeSpell15" Level="2" Damage="2" Attribute="4" Gesture="High Dive Power-Up" Number="0" Lock="0" />
  <Spell Name="SomeSpell16" Level="5" Damage="16" Attribute="5" Gesture="hdm7" Number="0" Lock="0" />
</Spells>
```

Figure 10. The XML file for spells

Monsters in this game have only three behaviors: idling, attacking, and being hit, and the animation is handled by the legacy Animation and Animation Clip components for simplicity. The behaviors and attributes of the monster are written in the Monster class, and the GameLogic class is in charge of spawning different monsters. The health bar of the monster is achieved by placing one texture (green) on top of another (red) which overlap each other completely when the health is full, and shows the color of green. When the monster is damaged by a spell, the Lerp function is called to gradually reduce the value of health as well as to move the green texture along the negative direction of x axis accordingly. Because the Mask component was added to the red texture, the part of green texture which is not overlapped with the red texture cannot be seen, and the red texture will gradually appear as the green texture moves away until the health bar turns to be total red when the monsters die.

The Unity Terrain component is used for the creation of the ground on which the monsters live. Since the game is design to be endless, the terrain should also be endless so that the player will always have the next destination. In this prototype, only one terrain is created manually, and it will be looped around with its copy to achieve the endlessness of terrains. The terrain is divided into four section (just like four quadrants if looks the terrain as a plane), and each section will be a battle field. The player will move to the next section after the boss in each section is defeated, and after finish all four section in one terrain, the player will move to the next terrain. To make the game more attracting, more terrains could be created or the terrains could be generated by some algorithm on run-time to make them not repeat to each other, which will give players more motivation to explore the game.

UI is also an important part of the game. There are two planes to draw the movement of hands on 2D space as well as two planes for displaying the trajectory of selected templates by using the Unity Texture2D class and SetPixel function. They make the job of testing the game much easier, and the players could know more about how they should perform the gestures. However, the disadvantage of using plane to display the hand track is that they will distract the players from the game, and force them to focus on generating desired 2D tracks to get better results. Therefore, the planes will be disabled during the gameplay.

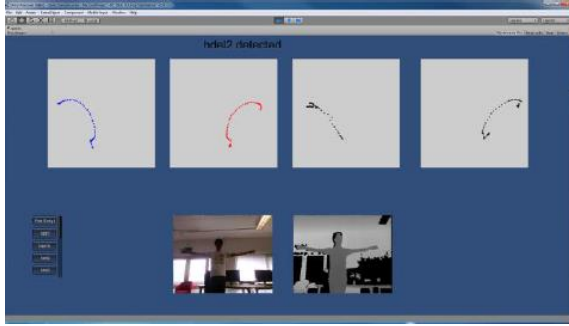


Figure 11. Planes draw gesture data

A menu for consuming the diamond and learning gestures is built based on the Unity Scroll View component. The menu can be opened using two hands swipe from right to left and closed using one hand swipe from left to right. All the available spells will be displayed in the list as UIButtons, and the distribution and positions of the spells are handled by components including Vertical Layout Group, Content Size Filter, and Layout Element. Gestures swipe down and swipe up are used to move the focus of the button to select spells, and swipe right to left using right hand will simulate the behavior of clicking the button to upgrade the spell. On the left side of the menu, an animation which indicates the related gesture to the spell is displayed for users to check how to cast spells with related gestures.

Lastly, the game data will be saved to record the progress that users made, so that they will be able to continue the process according to the save.

4. RESULT AND EVALUATION

The \$1 and DTW both have been tested using the same gesture set that is in the videos provided by the medicine school. Since the gesture data will be normalized, the distance between the user and the Kinect sensor is flexible as long as there is enough for Kinect to track the upper-limbs of the user entirely, which means the upper-limbs should be contained in the camera space no matter how they move. The test gestures were performed 20 times each, and some results show in Table 1 to demonstrate the strength and the weakness (gestures see Figure 12).

Algorithms Gestures	\$1	DTW(2D)	DTW(3D)
Gesture1	100%	100%	100%
Gesture2	0%	10%	20%
Gesture3	100%	100%	100%
Gesture4	100%	10%	10%
Gesture5	100%	80%	100%
Gesture 6	45%	55%	100%

Table 1. Test results

The recognition results are accurate overall for most of the selected gestures (some of them are not listed above), but the system performs poorly in recognizing some gestures.

Gesture 1, 2 and 3 are all recognized by the system with the perfect accuracy because the data is suitable for the algorithms to be applied.

The system fails to recognize Gesture 2. The reason why \$1 failed is that the skeleton cannot be tracked correctly when the hands overlap. For DTW, except for that reason, the gesture will be confused with Gesture 4 and 5 due to the normalization was just an approximation, and that leads to that results for those gesture are quite close, and for some unknown reason Gesture5 usually gets the highest score. This is the same reason why Gesture4 cannot be recognized correctly by DTW.

Gesture 3 and 4 generate similar tracks, but Gesture 3 can be recognized correctly while Gesture 4 cannot, this is because the constraints of Gesture 3 is unique in this gesture set and the constraints of Gesture 4 (the hands start from positions below the head and also end at positions below the head) are the same with Gesture 2 and 5.

Gesture 6 is recognized perfectly by the 3D version of DTW, while it is not that successful for \$1 and 2D DTW. The reason is that it sometimes generates circles on x-y planes while some times it is just lines, and it will still be confused with gesture 5 even it generates circles.



Figure 12. Some of the test gestures

5. CONCLUSION

In this section, the project will be summarized, and how this project could be expanded will also be discussed.

5.1 Project summary

The recognition system performs well in gesture recognition without training, and that meets the expectation of creating a simple system with acceptable efficiency.

Limitations are found for this project after running experiment. For the \$1 algorithm, it is an easy employed algorithm, and performs identically with the popular DTW algorithm. However, since it is originally developed for recognizing 2D shapes and hand writing, the performance of recognizing 3D gestures is limited. When distinguishing two gestures which will form similar tracks, the accuracy is unsatisfied. It fails to recognize some gestures which do not generate meaningful tracks, and the exploration of finding a proper solution was failed. Besides, it cannot deal with 1D tracks correctly without making changes to the Golden Section Search algorithm. The solution that is used in this project works solved this problem, but it restricts the gestures that users perform.

When DTW algorithm is used to recognize gestures in 2D space, it gives more flexibility for users to perform the gestures since it can handle 1D trajectory. The problem of DTW is an accurate way of normalizing the data for recognition has not yet been found, and the algorithms fails to recognize gestures with similar trajectories but different directions. Even after expanding DTW to support to recognize 3D trajectories, the problem still exists. On the other hand, using DTW to recognize gestures using 3D tracks instead of 2D projections improved the usability and accuracy of the system.

The rehabilitation game prototype has a completed game cycle and shows the basic concepts of the game design. There are still some features have not been added to the game due to time issue, and a lot of new elements can be added. For the accuracy of recognizing gestures, the start and the end of each gesture that users perform are set to holding for a short period of time, and that limits the use experience since users can only perform gestures according to the signals.

5.2 Further work

Although the recognition system works fine in the test environment, there is still a lot of optimization can be done. The currently system only tracks the movement of hands which are the major joints for upper limbs gestures. However, to track more joints like shoulders and elbows could help the system to recognize gestures more precisely and to distinguish gestures with similar hands movement more accurately.

Machine learning algorithms for classification could be implemented with proper training to improve the accuracy. \$1 is adapted to recognized body gestures without modifying the original 2D algorithm. It is possible to use the 3D trajectories instead of using the 2D projections with a proper understanding and modification for this algorithm,

Moreover, the template matching method for gesture recognition is not very suitable (at least if it is used as the only approach) for applications like games from my personal experience. Building a separated rule-based system and combining it with the template matching approach together can make the recognition system more practical.

Lastly, the game is only a prototype which shows the basic concept of the design and how the game will be played, and a lot of works can still be done. For example, the game system could be organized more robust, and the economy, difficulties and battles could be designed in more detail to fit for different rehabilitation treatment. A database can be established to store the information about the performance of different users, and a system can be built to analyze the patients' performances and to track the rehabilitation process.

6. REFERENCES

- [1] C. Su, C. Chiang & J.Huang, "Kinect-enabled home-based rehabilitation system using Dynamic Time Warping and fuzzy logic," *Applied Soft Computing Journal.*, vol. 22, pp. 652-666, Sep. 2014.
- [2] C. Kam Lai, P. Konrad & P.Ishwar, "A gesture-driven computer interface using Kinect," *2012 IEEE Southwest Symp. Image Analysis Interpretation*, pp. 185 – 188, Apr. 2012.
- [3] P. Turaga, R. Chellappa, Subrahmanian. V.S & O.Udrea, "Machine Recognition of Human Activities: A Survey," *IEEE Trans. Circuits Syst. for Video Technology*, vol. 18, no. 11, pp. 1473 – 1488, Nov. 2008.
- [4] G. Xu, Y. Cao &, "Action Recognition and Activity Understanding: A Review," *J. Image Graphics*, vol. 14, no. 2, pp. 189 – 195, Feb. 2009.
- [5] J. Aggarwal & Q. Cai, "Human motion analysis: A review," *Comput. Vision Image Understanding*, vol. 73, no. 3, pp. 428 – 440, Mar. 1999.
- [6] R. Ibanez, A. Soria, A. Teyseyre & Campo. M "Easy gesture recognition for Kinect," *Advances Eng. Software*, vol. 76, pp. 171 – 180, Oct. 2014.
- [7] Z. Ren, J. Yuan, J. Meng & Z. Zhang, "Robust Part-Based Hand Gesture Recognition Using Kinect Sensor," *IEEE Trans. Multimedia*, vol. 15, no. 5, pp. 1110 – 1120, Aug. 2013.
- [8] Y. Li, "Hand gesture recognition using Kinect," *2012 IEEE Int. Conf. Comput. Science Automation Eng.*, pp. 196 – 199, June. 2012.
- [9] O. Patsadu, C. Nukoolkit & B.Watanapa, "Human gesture recognition using Kinect camera," *2012 Ninth Int. Conf. Comput. Science Software Eng.*, pp. 28 – 32, May. 2012.
- [10] L. Miranda, T. Vieira, D.Martinez, T. Lewiner, A.W. Vieira & M. F. M. Campos "Real-Time Gesture Recognition from Depth Data through Key Poses Learning and Decision Forests," *2012 25th Sibgrapi Conf. Graphics, Patterns Images*, pp. 268 – 275, Aug. 2012.
- [11] E. Suma, D. Krum, B.Lange, S. Koenig, A. Rizzo & M. Bolas, "FAAST: The Flexible Action and Articulated Skeleton Toolkit," *2011 IEEE Virtual Reality Conf.*, pp. 247 – 248, Mar. 2012.
- [12] F. Kistler, B. Endrass, I.Damian, C. Dang & E. Andre, "Natural interaction with culturally adaptive virtual characters," *J. Multimodal User Interfaces*, vol. 6, no. 1-2, pp. 39 – 47, Jul. 2012.
- [13] J.O. Wobbrock, A.D. Wilson & Y.Li, "Gestures without libraries, toolkits or training: a \$1 recognizer for user interface prototypes," *User Interface Software Technology: Proc. 20th Annual Acm Symp.*, pp. 159 – 168, 2007.
- [14] B. Lange, Change. C, Suma. E, Newman. B, Rizzo. A, Bolas. M, "Development and evaluation of low cost game-based balance rehabilitation tool using the Microsoft Kinect sensor," *Annual Int. Conf. IEEE Eng. Medicine Biology Society. Conference*, 2011, vol. 2011, pp. 1831 – 1834, 2011.

- [15] L. Gallo, A.P. Placitelli & M.Ciampi, "Controller-free exploration of medical image data: Experiencing the Kinect," 2011 24th Int. Symp. Comput.-based Medical Syst., vol. 2011, pp. 1 – 6, June. 2011.
- [16] Playsaurus, 2014. Clicker Heroes. Available at: <<https://www.clickerheroes.com/>> [Accessed: June 24, 2015].
- [17] Mathsisfun, 2014. Correlation. Available at: <<https://www.mathsisfun.com/data/correlation.html>> [Accessed: June 24, 2015].
- [18] Gerald, C.F. & Wheatley, P.O., 2005. The Golden Section Search Method. Available at: <<http://homepages.math.uic.edu/~jan/mcs471/Lec9/gss.pdf>> [Accessed: June 24, 2015].
- [19] Microsoft, 2013. Microsoft Kinect - Microsoft SDK - Unity 3D. Available at: <http://wiki.etc.cmu.edu/unity3d/index.php/Microsoft_Kinect_-_Microsoft_SDK> [Accessed: June 24, 2015].
- [20] J.W. Burke, M.D.J. McNeil, D.K. Charles, P.J. Morrow, J.H. Crosbie & S.M. McDonough "Optimising engagement for stroke rehabilitation using serious games," *Visual Comput.*, vol. 25, no. 12, pp. 1085 – 1099, Dec. 2009.
- [21] C.A. Ratanamahatana & Keogh, "Everything you know about dynamic time warping is wrong," Third Workshop on Mining Temporal and Sequential Data, Aug. 2004.
- [22] S. Celebi, A. Aydin, T. Arici & T. Temiz, "Gesture Recognition using Skeleton Data with Weighted Dynamic Time Warping," Third Workshop On Mining Temporal and Sequential Data, vol. VISAPP, no. 1, pp. 620 – 625, Feb. 2013.
- [23] H. Hsu, "The Potential of Kinect in Education," *Int. J. Inform. Education Technology*, vol. 1, no. 5, p. 365, Dec. 2011.
- [24] G. Gallo, A. Placitelli & M.Ciampi, "Controller-free exploration of medical image data: Experiencing the Kinect," 2011 24th Int. Symp. Comput.-based Medical Syst., pp. 1 – 6, June. 2011.
- [25] P. Henry, M. Krainin, H. Herbst, Ren. Xf & Fox. D "RGB-D mapping: Using Kinect-style depth cameras for dense 3D modeling of indoor environments," *Int. J. Robotics Research*, vol. 31, no. 5, pp. 647 – 663, Apr. 2012
- [26] Unity Technologies, 2015. *Unity - Scripting API*. Available at: <<http://docs.unity3d.com/ScriptReference/>> [Accessed: August 27, 2015].