

Freiform-Gesten in einem Natural User Interface mit Microsoft Kinect

Dominik Feininger
 dominik.feininger@gmail.com
 University of Applied Science Mannheim
 2013

1. Abstract

In diesem Dokument wurden Erfahrungen über die Definition von Freiform-Gesten aus dem MMI Projekt „Earth Explorer“ und der Weiterentwicklung zusammengefasst. Die Erfahrungen im MMI Projekt kommen eher aus der praktischen Richtung. Die Definition ist noch nicht funktionstüchtig implementiert und ist daher eine Zusammenfassung einer Internet-Recherche. Die Entwicklung stützt sich wie viele Beispiele in anderen Kapiteln explizit auf eine Zoomgeste¹.

2. Kennenlernen der Menschlichen Motorik

Um zu verstehen, wie Gesten erkannt werden können, muss verstanden werden, wie der Benutzer eine Bewegung ausführt. Dazu gehört auch, welche Körperteile der Benutzer simultan mitbewegt, wenn er eine bestimmte Geste ausführt. So wird offensichtlich beim Bewegen der Hände zu einer Geste die Ellenbogen, Schultern usw. auch mitbewegt.

3. Definieren von Gesten

Definieren der Gesten kann überraschend schwierig sein. Die einfachsten Gesten sind diejenigen, die zu einem einzigen Zeitpunkt geschehen, und daher nicht auf vergangene Positionen der Gelenke zurückgreifen müssen. Zum Beispiel, wenn der Benutzer seine Hand über seinen Kopf hebt. Dies kann in einem einzelnen Frame überprüft werden. Komplizierte Gesten brauchen eine Zeit, um als solche erkannt zu werden. Für eine Zoom-Geste, ist es nicht möglich aus einem einzigen Frame zu erkennen, ob eine Person zoomt oder einfach nur die Hand vor den Körper hält.

Das Programm muss in der Lage sein, relevante Informationen aus der Vergangenheit zu speichern, nur welche Informationen sind relevant? Sind die letzten 30 Frames genug Information? 30 Bilder bekommt man in nur einer Sekunde, vielleicht 60 Frames, also 2 Sekunden? Oder 5 Sekunden, also 300 Bilder? Menschen bewegen sich nicht so schnell. Vielleicht könnte man jeden fünften Frame speichern, dadurch würde die Länge bei 5 Sekunden wieder auf 60 Frames reduziert werden.

Eine bessere Idee wäre, die relevanten Informationen aus den Frames zu wählen. Für eine Zoom-Geste der Hände, die aktuelle Geschwindigkeit, wie lange sich bewegt wurde, wie weit wurden die Hände verschoben, etc. könnten nützliche Informationen sein.

Nachdem das definiert wurde, speichert man alle Informationen, die im Zusammenhang mit der Geste stehen. Doch wie definiert man diese Informationen in Zahlen? Zoomen könnte eine bestimmte Mindestgeschwindigkeit, oder eine Richtung (links / rechts nach außen), oder eine Dauer erfordern.

¹ Geste mit beiden Armen, die einem vertikalen Pinch ähnelt.

Hier sind jedoch nicht die definierten 5 Sekunden von Interesse. Diese Dauer beschreibt das absolute Minimum² das benötigt wird um davon auszugehen, dass der Benutzer zoomt. Wie bereits erwähnt, kann das nicht anhand eines einzelnen Frames bestimmt werden. Es sollte nicht allein die Implementierte Dauer bestimmen, ob eine zoom Geste als solche erkannt wird. Da die Zeitspanne in der die Geste durchgeführt wird für jeden Benutzer unterschiedlich ist, da jeder Benutzer sich unterschiedlich schnell bewegt.

Die Chancen stehen gut, dass an einem gewissen Punkt innerhalb dieser 5 Sekunden die Geste mit der Definition übereinstimmt und ein zoomen erkannt wird. Es gibt aber ein weiteres Problem, wie definiert man den Start und Endpunkt einer Geste? Die Geste kann während des Ausführens langsamer und schneller werden oder gar pausieren. Fängt die Geste nach der Pause wieder von Anfang an?

Der Punkt, den ich versuche hier klar zu machen ist, es gibt keine einfache Möglichkeit Gesten zu erkennen. Als Entwickler muss man sich in die Geste hinein denken und eine Reihe von gemeinsamen Positionen der Skelett-Daten über eine Zeit hinweg als eine Geste zu definieren.

4. Benutzererwartung

Der Benutzer sollten nicht gezwungen werden eine Bewegung für einen bestimmten Zeitraum zu wiederholen. Es ist überraschend anstrengend und nicht einfach. Der Benutzer erwartet eher wie beim Computer eine Art „Point-and-Click“, sobald geklickt wird, wird eine Reaktion des Systems erwartet. Es ist sehr unnatürlich einen Mausklick (z.B. 5sekunden) halten zu müssen bevor eine Anwendung startet.

Wiederholen einer Geste über einen Zeitraum ist okay, wenn es sich um eine stetige Aktion handelt und stetig Rückmeldung vom System gegeben wird. Somit versteht der Benutzer, dass für diese Aktion eine kontinuierliche Ausführung erforderlich ist.

5. Fakten

Dem Benutzer ist es nur begrenzt möglich seine Hände still zu halten. Beim Versuch die Hände still zu halten erkennt die Kinect Bewegungsänderungen in der dritten Nachkommastelle. Man kann von einer gewollten Bewegung sprechen, falls sich Werte an der zweiten Nachkommastelle über längere Zeit ändern^[2]. Das Koordinatensystem, ist folgendermaßen angeordnet:

X-Achse = rechts/ links
Y-Achse = hoch/ runter
Z-Achse = vor/ zurück (Tiefe)

Wobei die X-Achse sich relativ zum Körper des Benutzers verhält, die Z-Achse sich jedoch nicht verschiebt, da der Ursprung der Z-Achse in der Kinect liegt.

² Minimum Zeitdauer zum Erkennen einer Geste

6. Berechnungen

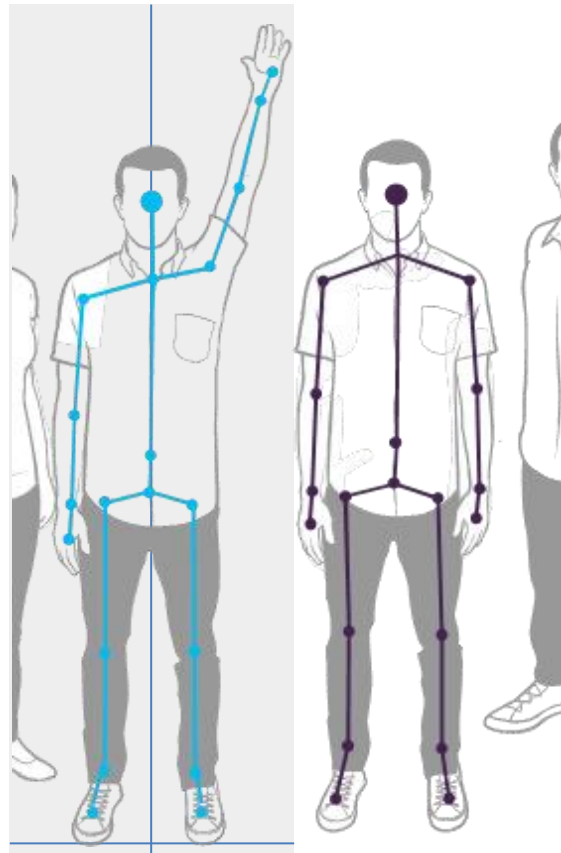


Abbildung 1 - Kinect User Skeleton³

Ausgegangen vom Körpermittelpunkt als Y-Achse befinden sich die linke Hand im negativen Bereich und die rechte Hand im positiven Bereich. Somit muss bei der Berechnung von Bewegungen jeweils Werte addiert bzw. subtrahiert werden. Der Bewegungsraum einer Hand wird sich im normalen Fall und bei optimaler Entfernung des Benutzers zur Kinect zwischen 0 und $\pm 0,8$ befinden.

7. Übersetzen von Gesten

Eine wichtige Größe bei der Übersetzung von Gesten auf ein beliebig großes Display oder Leinwand bringt einige Komplikationen mit sich. Es besteht ein Übersetzungsverhältnis von Geste zum erkennenden Gerät und zum Display. Die Bewegung die der User ausführt steht in einem bestimmten direkt proportionalen Verhältnis zu der Aktion auf dem Display. Um dieses Verhältnis genau zu bestimmen und in verschiedenen Umgebungen gleich zu halten muss eine Berechnungsmethode gefunden werden die die beteiligten Parameter verbindet und daraus eine Proportionale Umgebung schafft. Parameter dieser Gleichung sind z.B. Abstand von User zum Erkennungsgerät, Größe des Displays und Auflösung des Displays.

Theoretisch kann der Eingaberaum, in dem sich der User bewegt sehr groß sein, da keine Physische Nähe gefordert ist. Hier limitiert jedoch der Sensor zum Aufnehmen den Bewegungen den Maximalen Bereich. Weiterhin ist zu beachten, dass die Eingabe ohne eine Berührungen einer Festen

³ <http://flurfunk.sdx-ag.de/2013/03/bewegungsdrang-teil-5-die-kinect-geht.html>

Oberfläche deutlich langsamer und behäbiger abläuft als bei Interaktionen mit einer festen Oberfläche.

8. Erfahrungen

Zu allen Überlegungen aus Kapitel 3 kommen in der tatsächlichen Umsetzung weitere Konstanten die für bestimmte Gesten sinnvoll sein können. Zusätzlich zu dem Minimum an Distanz die die Hand zwischen zwei gemessenen Frames zurücklegen muss kann ein Wert, der ein Minimum für die absolute gesamte Bewegungsdistanz festlegt sinnvoll sein. Damit können kleinere Zuckungen des Benutzers und smoothie Endpunkte der Gestenrückmeldung erzeugt werden.

Grundsätzlich erzeugt die Steuerung über die Kinect eine Zeitverzögerung vom Ausführen der Geste, bis zur Anzeige auf dem Bildschirm von ca. einer viertel Sekunde⁴ pro Frame. Das kann, falls eine zu geringe Frame Rate benutzt wird zu Störungen der direkten Systemrückmeldung führen da sich die Summe aus Lack und gemessenem Frame Abstand sich auf bis zu einer halben Sekunden belaufen können.

In Abbildung 3 wird veranschaulicht, wie sich der Aktionsradius des Nutzers auf den Erkennungsbereich der Kinect auswirkt. Um bei ausgestreckten Armen einen Abstand vom Benutzer zur Kinect in Z-Richtung angeben zu können muss bekannt sein, in welche Richtung der Nutzer die Arme streckt. Es sollte, rein für die Geste des Zoomens auf ein 2D Koordinatensystem zurückgegriffen werden. Dieses Koordinatensystem, das sich in Z- und X-Richtung erstreckt kann beliebig unterteilt werden.

⁴ Gemessen durch Benutzung