

Bachelorarbeit

Automatisches Labeln von Objekten in einer Augmented Reality Umgebung

Janelle Pfeifer

Delpstraße 28

53359 Rheinbach

janelle.pfeifer@smail.inf.h-brs.de

Hochschule Bonn-Rhein-Sieg

Institute of Visual Computing

Fachbereich Informatik

Studiengang: Informatik (B.SC.)

Erstprüfer: Prof. Dr. Ernst Kruijff

Zweitprüfer: Prof. Dr. André Hinkenjann

Rheinbach, 1.10.2020

Inhaltsverzeichnis

Bachelorarbeit	1
1 Zusammenfassung	3
2 Einleitung	3
3 Grundlagen	3
3.1 Webrequests	5
3.2 Lokale und globale Koordinatensysteme in 3D Szenen	5
3.3 Kamera in 3D-Computergrafik	6
4 Design	6
5 Implementierung	7
5.1 Ein Foto aufnehmen	7
5.2 Azure Object Detection	7
5.3 Von dem Foto zum 3D Raum	7
5.4 Raycast	9
5.5 LabelCreator	10
5.6 Azure Custom Vision	10
5.7 Entwicklung der Foto-Repräsentation	12
6 Zusammenfassung	14

1 Zusammenfassung

2 Einleitung

Augmented Reality (AR) ist eine Vermischung der realen Welt mit digitalen Elementen. Es wird durch Anzeigegeräte, wie Handys, Tablets oder Augmented Reality Brillen präsentiert und bietet ein intuitives Benutzerinterface um Informationen über Objekten der realen Welt anzuzeigen. Dafür müssen Informationen über die Umgebung erfasst werden. Es ist wichtig zu bestimmen welche Objekte sich in der Umgebung befinden, die durch AR erweitert werden soll.

Es gibt mehrere Möglichkeiten reale Objekte zu erkennen. Zum einen können Markierungen in der realen Welt verwendet werden. Dabei handelt es sich um statische Bilder, beispielsweise ein Foto, oder ein QR Code, die von einer Kamera eingescannt werden. Der Marker ist einzigartig für jedes Object, das erkannt werden soll, damit sie voneinander unterschieden werden können. Der Nachteil bei diesem Vorgehen ist der Arbeitsaufwand, der damit verbunden ist, jeden Gegenstand einzeln zu Markieren und der AR Applikation die Marker bekannt zu machen.

Wenn man Markierungen in der Realen Welt umgehen möchte, kann man den Nutzer der Applikation bitten, per Geste oder per eye sight auf Objekte der Realen Welt zu weisen, die erkannt werden sollen. Dabei muss von dem Nutzer auch angegeben werden, um welches Object es sich genau handelt, damit die Applikation unterschiedliche Objekte auseinander halten kann und die korrekten Informationen mit ihnen assoziiert. Auch hier ist ein hoher Arbeitsaufwand damit verbunden alle Objekte für die Application auszuweisen.

Beide der Verfahren sind nicht auf große AR Umgebungen skalierbar, da sie sehr Arbeitsintensiv sind. Nur eine voll automatische Objekterkennung ist skalierbar.

Um diese Automatisierung zu erreichen kann Image based Object Detection aus dem Bereich der Computer Vision verwendet werden. Dabei werden Objekte in Bildern erkannt, indem nach Charakteristiken gesucht werden, die unterschiedliche Arten an Objekten auszeichnen.[1]

In dieser Thesis wird das Erkennen und Labeln von Objekten in einer AR Umgebung, mithilfe von Image based Objekt Detection, automatisiert. Dabei wird as AR Gerät "Magic Leap" verwendet.

3 Grundlagen

Object Detection

Objekt Detection ist eine Aufgabe in Computer Vision bei der Objekte in einem Bild identifiziert werden. Es wird bestimmt um welche Klasse an Objekt es sich handelt, beispielsweise ob es eine Katze oder ein Hund ist, und wo sich das Objekt befindet. Die Ausgabedaten dieser Untersuchung ist eine Liste an Objektarten und eine Liste an Bounding Boxen, die die Positionen angeben.

Object Detection kann mit Artificial Neural Networks durchgeführt werden.

Artificial Neural Networks

Artificial Neural Networks sind Machine Learning Architekturen. Sie können beispielsweise Musik, Text oder Bilder nach Mustern durchsuchen. Sie sind für keine genaue Aufgabe programmiert, sondern lernen indem sie mit Beispieldaten trainiert werden. Für jedes Beispiel gibt es ein Label, das angibt ob es das gesuchte Muster enthält oder nicht. Die Struktur des Networks verfügt über Gewichte, die Einfluss auf den Output haben. Mit jedem Trainingsbeispiel passt das Network die Gewichte an, sodass der Output dem Label des Beispiels entspricht.[1, 2]

Convolutional Neural Networks

Convolutional Neural Networks sind auf das Verarbeiten von Bildern spezialisiert. Sie nutzen aus, das Bilder viele Redundanzen und Informationsarme Bereiche haben, indem sie mit jedem Verarbeitungsschritt Informationen weglassen. So können Rechenzeit und Trainingsdaten verringert werden.[1, 2, 3]

Azure maschinelles Sehen

Microsoft Azure bietet einen Computer Vision Service an, der für Object Detection trainiert ist. Der Anwender sendet ein Bild an Microsoft, dort wird es verarbeitet und ein Ergebnis zurückgegeben.[4, 5, 6, 7]

Der Service wird über einen web request angefragt. Und sendet eine Json Datei zurück.

Über ein web endpoint per web request anfragen. Bekommt eine Json zurück. darin steht welche Obejkte darin sind, daten über das foto und objekte mit bounding box angegeben. Beispiel von Json angeben.

Azure Custom Vision

Azure bietet zusätzlich einen Computer Vision Service an, den der Nutzer Trainieren kann um bestimmte Objekte Klassifizieren zu können.[8]

was ist die Prediction? Erkläre die Iterations Erkläre was die Genauigkeit der Prediction aussagt.

Spatial Mapping

Durch Spatial Mapping wird eine 3D Abbildung einer realen Umgebung erschaffen. So können Hologramme mit der echten Welt interagieren, diese Verdecken, oder von ihr verdeckt werden.[9]

Magic Leap AR Brille

Die Magic Leap 1 Lightwear ist eine Augmented Reality Brille. Sie besitzt neun Sensoren und mehrere Kameras. Dazu gehören:

- ein Infrarot Tiefen-Sensor,
- ein Eye Tracker,

- eine Foto und Video Kamera, die im Format 16:9 mit einer Auflösung von 1920 x 1080 Pixeln aufnehmen,
- Umgebungskameras die in unterschiedliche Richtungen ausgerichtet sind. [10, 11]

Output geschieht über ein Display mit einem 50 Grad Field of View und einem Seitenverhältnis von 4:3. Input wird durch ein 6 Degree of Freedom Controller gegeben.[10, 11]

Die Magic Leap Brille läuft auf dem Betriebssystem Lumin OS. Dieses wurde für Spatial Computing entwickelt und stellt Applikationen entsprechende Funktionalitäten zur Verfügung. Beispielsweise führt das Betriebssystem Spatial Mapping durch. Mit den Sensoren und Kameras der Brille werden Daten aufgenommen und in einen Zeitlichen Zusammenhang mit der Bewegung der Brille gesetzt, um eine Rekonstruktion des Raumes zu erhalten.[12, 13, 14, 15]

Lumin OS bietet es Applikationen an,

- Raycast auf die Umgebung durchzuführen und
- ein Mesh der Rekonstruktion zu erhalten.

Lumin OS unterstützt das Verarbeiten von dem Input des Controllers und verwaltet die Privilegien der Applikationen. Dazu gehört beispielsweise der Zugriff auf die Fotokamera und das Netzwerk.[13, 16]

3.1 Webrequests

3D-Computergrafik

3.2 Lokale und globale Koordinatensysteme in 3D Szenen

In einer 3D Szenen werden die Positionen von Objekten als Matrizen in drei Dimensionalen Koordinatensystemen verwaltet. Es gibt ein globales Koordinatensystem (auch Weltkoordinatensystem oder World Space), in dem alle Objekte relativ zu einem gewählten Ursprung liegen.

Jedes Objekt hat zusätzlich ein eigenes, lokales Koordinatensystem (Objektkoordinatensystem). Dessen Ursprung liegt in dem jeweiligen Objekt. Die Position und Rotation des Objektes in dem globalen Koordinatensystem bestimmt die Relation zwischen dem globalen und dem lokalen Koordinatensystem.

Das lokale Koordinatensystem einer Kamera wird auch Camera Space genannt. Die Relation zwischen dem Camera Space und dem globalen Koordinatensystem wird in Unity durch die cameraToWorld Matrix beschrieben. Mithilfe dieser Matrix kann eine Koordinate aus dem Camera Space in die entsprechende Koordinate des globalen Koordinatensystems transformiert werden.[17]

Dazu wird die Koordinate als Vektor angegeben und mit der cameraToWorld Matrix multipliziert. Das Resultat ist ein Vektor, der eine Koordinate im globalen Koordinatensystem angibt.[17, 18]

3.3 Kamera in 3D-Computergrafik

View Frustum ist der Teilvolumen einer 3D Szene, die auf den zweidimensionalen Bildschirm abgebildet wird. Alle Objekte die von der Kamera gesehen werden, befinden sich in dem View frustum.

Clipping Plane bezeichnet eine Ebene, die den view frustum quer zur Blickrichtung begrenzt. Es gibt eine vordere und eine hintere Clippingebene. Die vordere Clippingebene liegt nahe an der Kamera dran. Alle Objekte die zwischen der Kamera und der Ebene liegen werden nicht angezeigt.

Die hintere Clippingebene limitiert wie weit Objekte entfernt sein können, bevor sie nicht mehr zu sehen sind.

4 Design

Das Ziel ist das Erkennen und Labeln von Objekten in einer AR Umgebung, durch Image Based Objekt Detection.

Wenn der Nutzer das Signal gibt, beginnt die Detection. Als Erstes wird ein Foto mit der Kamera der AR Brille aufzunehmen. Dieses Foto wird dann an Azure Object Detection und Azure Custom Vision geschickt. Als Ergebnis gibt es jeweils eine Json Datei. Darin wird angegeben, welche Arten von Objekten auf dem Foto gefunden wurden (Hund, Person, Computermouse) und wo sie sich jeweils befinden. Die Position eines Objektes ist dabei mit einer Bounding Box bestimmt. Um das Markieren von Objekten in der AR Umgebung zu vereinfachen, wird die Mitte der Box als Position des Objektes genommen.

Jedes Objekt soll an der tatsächlichen Stelle der Realen Welt markiert werden. Dafür wird in der 3D Szene der AR Umgebung eine Virtuelle Repräsentation des Fotos erschaffen. Die Fotorepräsentation muss die richtige Skalierung, Position und Rotation haben, um das räumliche Verhältnis zwischen der realen Foto-Kamera und der Umgebung nachzubilden.

Da die Foto-Kamera und das Display nahe beieinander liegen und den gleichen Blickwinkel haben, kann die Position des Displays als Repräsentation des Fotos genutzt werden. In der 3D Szene ist das Display mit der Hauptkamera gleichgesetzt. Die Clipping Plane der Kamera hat somit die gleiche Rotation und eine zumindest ähnliche Position und Skalierung wie das Foto.

Daher werden die Foto-Positionen auf Koordinaten der Clipping Plane abgebildet. Dabei werden verbleibende Positions- und Skalierungs-Unterschiede ausgeglichen. Für jedes Objekte wird so eine Koordinate auf der Clipping Plane bestimmt.

Als Nächstes wird ein Raycast von der Kamera aus durch die Clipping Plane Koordinate geschickt. Der Raycast schneidet sich mit einem Mesh, das die reale Welt abbildet. Die getroffene Position wird mit einem Schriftzug markiert. Dort befindet sich das Objekt, das auf dem Foto gefunden wurde.

Alle Objekte, die Azure Object Detection und Azure Custom Vision gefunden haben, werden so für den Nutzer in der AR Umgebung markiert.

5 Implementierung

Das Projekt wurde in Unity umgesetzt, und wurde für eine Magic Leap AR Brille entwickelt. Es wurde ein Unity Projekt Template von Magic Leap verwendet.

Zusätzlich werden einige vorgefertigte Klassen von Magic Leap verwendet. Dazu gehören MLInput, MLCamera, MLRaycast, MLPrivilegeRequestBehavior und MLSpatialMapper.

MLSpatialMapper setzt spatial mapping um und erzeugt ein Mesh, das die Umgebung abbildet.

5.1 Ein Foto aufnehmen

Das Script TakePicture implementiert das aufnehmen von einem Foto. Die Methode StartCapture wird aufgerufen, um den Process zu starten. MLCamera wird genutzt um die Kamera der Magic Leap Brille anzusteuern.

Wenn das Foto aufgenommen wurde, wird OnCaptureRawImageComplete aufgerufen. Die Daten des Bildes werden an die Klassen Azure ObjectDetection und AzureCustomPrediction weitergegeben. Dort wird die Analyse der Bilder gestartet.

5.2 Azure Object Detection

In der Methode AnalyseImage von AzureObjectDetection wird ein Web Request zusammengestellt. Azure stellt dafür einen Web endpoint zur Verfügung. In dem Request muss ein Authorization Key für den Service und die Bild-Daten enthalten sein.

Der Webrequest wird geschickt und es wird auf die Antwort gewartet. Wenn es keine Probleme mit dem Netzwerk gab, wird eine JsonDatei als String in der Antwort auf den Webrequest geschickt. Siehe Abbildung ??

Es wird für jedes gefundene Object auf dem Foto eine Bezeichnung und eine Bounding Box angegeben, in dem sich das Object auf dem Bild befindet.

Die Json Datei wird in HandleJsonResponse verarbeitet. Für den erwarteten Aufbau der Datei wurden drei Klassen geschrieben. Der Json String wird mit JsonUtility in ein DetectionResponse Object umgewandelt. Dabei werden alle gefundenen Foto-Objekte in einer Liste von DetectedObject abgelegt.

Die gefundenen Object sollen im 3D Raum mit einer Markierung ausgezeichnet werden. Dafür wird für jedes DetectedObject die Methode Cast von der Klasse PixelToWorld aufgerufen. Der Methode wird der Mittelpunkt der BoundingBox als u,v Foto-Koordinate für das DetectedObject übergeben.

5.3 Von dem Foto zum 3D Raum

In der Methode Cast wird die u,v Foto-Koordinate des Fotos und die Kameraausrichtung verwendet um das gefundene Foto-Objekt in der 3D Abbildung der Realen Welt zu lokalisieren.

Das Foto kann mit dem Display und somit mit der Clipping Plane der Hauptkamera approximiert werden. Die u,v Foto-Koordinate wird zunächst in eine x,y,z Koordinate in dem Camera Space umgewandelt. Der z Anteil gibt die Achse an, die von der Kamera nach vorne durch den View Frustum verläuft. Bei $z = -0.4$ wird der Punkt auf der Achse gewählt, auf dem die Clipping Plane liegt.

Die x und y Dimensionen beschreiben die Achsen, die horizontal und vertikal zur Clipping Plane verlaufen. Mit dem festgelegten $z = -0.4$, kann jeder Punkt auf der Clipping Plane durch x und y angegeben werden. Dazu gehören auch Punkte die außerhalb des View Frustum liegen.

Es wurden Werte für x und y ausprobiert, mit denen die Ränder des Fotos auf der Clipping Plane angegeben werden können. Dabei wurde darauf geachtet, das das Foto ein Seitenverhältnis von 16:9 hat, während das Display, und somit die Kamera, ein Seitenverhältnis von 4:3 hat. Darüber hinaus ist der Bildausschnitt des Displays kleiner. Daher liegen die Ränder des Fotos außerhalb des View Frustum.

Sind die x und y Werte für die Ränder bekannt, ergibt sich für die Achsen jeweils ein Intervall, die kombiniert alle Foto-Koordinaten auf die Clipping Plane abbilden können. Die Intervall lauten: $[-0.295, 0.2281]$ für x und $[0.1546, -0.1507]$ für y. Durch die Intervalle wird die Position und Skalierung des Fotos in Relation zu dem Display - und der Hauptkamera- berücksichtigt. Siehe Kapitel 5.7 für die Entwicklung des Cast Methode und die Ermittlung der Intervall Werte.

Es werden zwei Lineare Funktionen aufgestellt:

- Die Funktion X Bildet das Intervall für u $[0, 1920]$ auf das Intervall für x $[-0.295, 0.2281]$ ab.
- Die Funktion Y Bildet das Intervall für v $[0, 1080]$ auf das Intervall für y $[0.1546, -0.1507]$ ab.

Siehe Abbildung 1.

```

40 //Picture u and v ranges
41 private float umin = 0; //left
42 private float umax = 1920; //right
43 private float vmin = 0; //up
44 private float vmax = 1080; //down
45 //Offset Vektor x and y ranges
46 private float xmin = -0.295F; //left
47 private float xmax = 0.2281F; //right
48 private float ymin = 0.1546F; //up
49 private float ymax = -0.1507F; //down
50 private float X(float u)
51 {
52     float slope = ((xmax - xmin) / (umax - umin));
53     float b = xmin - slope * umin;
54     return slope * u + b;
55 }
56 private float Y(float v)
57 {
58     float slope = ((ymax - ymin) / (vmax - vmin));
59     float b = ymin - slope * vmin;
60     return slope * v + b;
61 }

```

Abb. 1. Funktionen X und Y

Mit den Funktionen wird für u,v eine Koordinate im Camera Space berechnet. Diese Koordinate wird dann mithilfe der cameraToWorldMatrix in die Koordinate p des globale Koordinatensystem umgewandelt. Damit wird die Position und Rotation der Kamera - und somit des Fotos - in der 3D Szene berücksichtigt. Siehe Abbildung 2.

```

28 public void Cast(float u, float v, SavedCameraState cpos,GameObject clippingPlaneMarker,
29     string objectName, bool showClippingPlane, int material)
30 {
31     //scale v,v to x,y range
32     Vector3 offset = new Vector3(X(u), Y(v), -0.4F);
33     Vector3 p = cpos.cameratoWorldMatrix.MultiplyPoint(offset);
34     Raycast.instance.StartCast(Raycast.instance.CreateRaycastParams(cpos.ctransform, p), objectName, material);
35     if (showClippingPlane)
36     {
37         GameObject sphere2 = Instantiate(clippingPlaneMarker, p, Quaternion.identity);// show point on clipping plane
38     }
39     ResultAsText.instance.Add(u + " " + v + " " + X(u) + " " + X(v) + " " + objectName + " object marked");
40 }

```

Abb. 2. Cast Methode

5.4 Raycast

Als Nächstes wird ein Raycast durch den Ursprung der Kamera und die Koordinate p gesendet. MLRaycast wird genutzt, um einen Schnittpunkt mit der Rekonstitution der Welt von Lumin OS zu bestimmen. Die Stelle die der Raycast trifft beschreibt die Position des DetectedObject im 3D Raum.

Für den MLRaycast werden zwei Parameter benötigt:

- Ein QueryParams Objekt, das Ursprung und Richtung für den Raycast beinhaltet.
 - Ursprung: Kamera
 - Richtung: Richtungsvektor von der Kamera zu der Koordinate p
- Eine Methode die aufgerufen wird, wenn der Raycast fertig ist.
 - Callback Methode: HandleOnRecieveRaycast

Siehe Abbildung 3.

```

19 public MLRaycast.QueryParams CreateRaycastParams(Transform ctransform, Vector3 target)
20 {
21     MLRaycast.QueryParams _raycastParams = new MLRaycast.QueryParams
22     {
23         // Update the parameters with our Camera's transform
24         Position = ctransform.position,
25         Direction = target - ctransform.position,
26         UpVector = ctransform.up,
27         // Provide a size of our raycasting array (1x1)
28         Width = 1,
29         Height = 1
30     };
31     return _raycastParams;
32 }

```

Abb. 3. Cast Methode

Wenn der Raycast fertig ist, wird die Methode `HandleOnReceiveRaycast` aufgerufen. Der Parameter `point` beinhaltet dabei die getroffene Koordinate. Diese an die Methode `CreateMarker` von der Klasse `LabelCreator` weitergegeben.

5.5 LabelCreator

`CreateMarker` erhält den Punkt `point` der getroffen wurde und die Bezeichnung für das `DetectedObject`. An der Koordinate von `point` wird ein Prefab `GameObject` instanziiert, das als Markierung für das `DetectedObject` in der 3D Umgebung dient.

Das Prefab besteht aus einer Kugel und einem Schriftzug, der den Namen des `DetectedObject` anzeigen soll. Dem neu instanziierten `GameObject` wird die Bezeichnung des `DetectedObject` als Schriftzug zugewiesen. Siehe Abbildung 4.

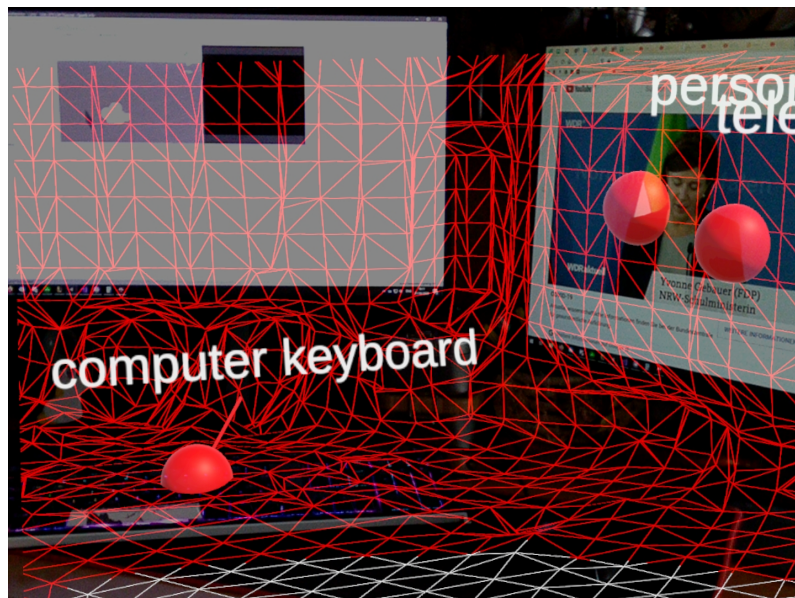


Abb. 4. Markierungen in der Welt

5.6 Azure Custom Vision

Neben der Bildanalyse mit Azure Object Detection wird auch Azure Custom Vision verwendet. Die AI wurde über die Webseite trainiert.

Die Anfrage an den Service passiert in der Klasse `AzureCustomPrediction`. Ähnlich wie bei `AzureObjectDetection` wird ein `Webrequest` erstellt mit einem `authorization key` für den Service und einem Foto als `payload`.

In der Antwort wird eine `Json` Datei zurückgeschickt, die die gefundenen `Objecte` angibt. Da die `Json` Datei ein etwas anderes Format hat, wurde eine eigene `HandleJson` Methode dafür geschrieben.

Für jedes erkannte Objekt wird die Methode `Cast` von `PixelToWorld` aufgerufen, um das Objekt in der realen Welt zu lokalisieren und zu markieren.

Das Trainieren Es wurde Probiert das Cusotm Vision Modell auf drei unterschiedliche Objekte zu trainieren. Dabei wurden vier Iterationen erstellt.

Iteration 1:

Zunächst wurde probiert Tuben von Acrylfarbe zu erkennen. Die Genauigkeit war davon war nicht so gut. Es wurden in den Fotos Acrylfarben erkannt, wo es keine gab. Sieh Abbildung 5.

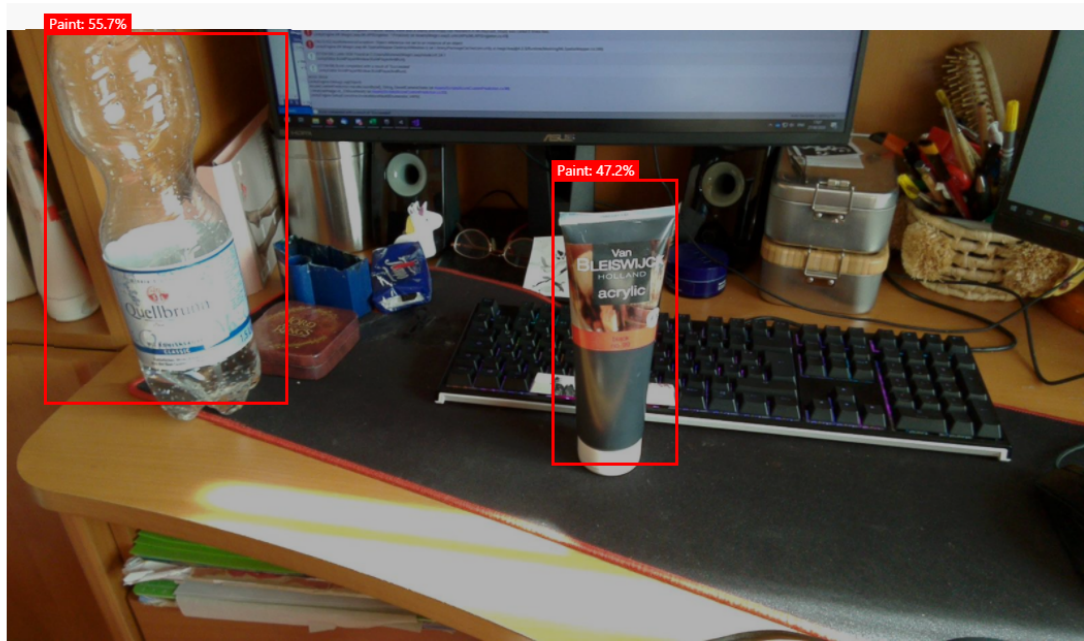


Abb. 5. Beispielfoto Iteration 1. Die Wasserflasche wurde als Farbe markiert mit 55.7 Prozent statistischer Konfidenz. Die Tatsächliche Farbtube wurde mit einer Konfidenz von 47.5 Prozent erkannt.

Iteration 2:

In der zweiten Iteration wurde probiert das Model darauf zu trainieren, eine Blaue Dose von Nivea Hautcreme to erkennen. Die Form und Farbe von der Dose ist sehr simpel, daher wurde davon ausgegangen, das sie leichter zu erkennen ist. Auch dort dar die Genauigkeit zunächt bei nur 80 Prozent. Es wurden in vielen Fotos fälschlicherweise Nivea Dosen erkannt.

Iteration 3:

In der Iteration drei wurde versucht Iteration 2 zu verbessern. Es wurden ausgewählte Training Fotos entfernt, die die Dose von einem Seitlichen Winkel zeigten, mit der mit der Erwartung das die Detektion der Dose aus dem Blickwinkel von Oben damit konsistenter wird. Zusätzlich wurden mehr Fotos von der Dose auf verschieden farbigen und gemusterten Untergrund hinzugefügt.

Die Genauigkeit der Prediction sank auf 75 Prozent.

Iteration 4:

In Iteration 4 wurden zwei Fotos von der Nivea Dose entfernt, das die Genauigkeit auf 100 Prozent steigen ließ. In der Umsetzung mit der Magic Leap Anwendung werden trotzdem häufig Objekte fälschlicherweise als Nivea Dose markiert.

Neben der Dose wurde diese Iteration darauf trainiert eine bestimmte Holz Haarbürste zu erkennen. Aufgrund von dem komplexeren, und markanten Aussehen der Bürste ist davon ausgegangen, das die Bürste besser von anderen Objekte zu unterscheiden ist. Die Bürste wurde nur mit den Borsten nach oben Photographiert.

Die Genauigkeit für die Bürste lag bei 100 Prozent. In der Umsetzung mit der Magic Leap Anwendung wird die Bürste häufig nicht erkannt, obwohl sie im Bild ist und mit den Borsten nach oben liegt. Es werden jedoch keine Objekte als Haarbürste erkannt die keine sind.

5.7 Entwicklung der Foto-Repräsentation

Um die u,v Foto-Koordinate eines gefundenen Objektes auf der Clipping Plane der Kamera zu lokalisieren wurden ein paar Herangehensweisen ausprobiert.

Das Ziel ist das setzten einer Markierung in dem 3D Raum, basierend auf der Foto-Koordinate. Das Foto beinhaltet keine Information über die Entfernung zu dem Objekt. Dafür muss ein Raycast durchgeführt werden.

Mit einer Repräsentation des Fotos in dem 3D Raum ist es möglich ein Raycast durchzuführen. Dazu muss das Foto nicht tatsächlich in dem 3D Raum vorhanden sein. Es muss jedoch mit dem Input von einer Foto-Koordinate ein Output von einer Koordinate in dem 3D Raum geben, mit dem der Raycast durchgeführt werden kann.

Die Position des Fotos hängt mit der Kamera zusammen, daher kann das Foto durch den Camera Space simuliert werden. Als erstes wurde probiert ein Sphären-Objekt an eine gezielte Koordinate in dem Camera Space zu bewegen.

Wenn die Kamera an dem Ursprung des globalen Koordinatensystem liegt und eine neutrale Rotation hat, stimmt der Camera Space mit dem globalen Koordinatensystem überein. Die Sphäre wurde in der Szene per Hand bewegt um markante Koordinaten des Camera Space abzulesen. Dabei wurden folgende Camera Space Koordinaten gefunden:

- Near Clipping Plane bei $z = -0.37$
- linker Rand bei $x = -0.153$
- rechter Rand bei $x = 0.153$
- oberer Rand bei $y = 0.1147$
- unterer Rand bei $y = -0.1147$

Die x und y Koordinaten hängen von der u,v Koordinate des Fotos ab. Es wurden linear Funktionen aufgestellt um u,v auf x,y Abzubilden. Diese Abbildung dient als Repräsentation des Fotos im 3D Raum. Zusätzlich gibt die Position und Skalierung des Fotos im Verhältnis zu der Kamera an.

Dann wurde getestet wie genau DetectedObjects in der AR Umgebung lokalisiert werden. Es wurden Testweise Fotos aufgenommen, analysiert und die DetectedOb-

jects markiert. Die entstandenen Markierungen lagen in Sichtfeld, jedoch nicht an den erwarteten Stellen.

Um dem Problem auf den Grund zu gehen, wurde ein UI Objekt erstellt, das ein Aufgenommenes Foto bei Runtime anzeigt. Das Foto wurde dann mit dem Display verglichen. Dabei viel auf, das sie ein unterschiedliches Seitenverhältnis haben und das Display einen kleinen Bildausschnitt zeigt.

Mit den Unterschieden zwischen dem Foto und dem Display zwei Möglichkeiten. Entweder wird das Foto auf das Display zugeschnitten, oder das gesamte Foto wird genommen. Dann würden auch Objekte erkannt, die außerhalb des Sichtfeldes liegen. es wurde die Entscheidung getroffen das Foto zuzuschneiden. Damit gibt es ein besseres Feedback für den Nutzer, das ein Objekt gefunden wurde.

Die Zugschneidung wurde realisiert, indem die Intervalle für u und v der Abbildungsfunktionen stärker eingegrenzt wurden. Alle Objekte die Außerhalb der Intervalle liegen werden ignoriert. Um die Intervalle zu bestimmen wurde ein Gitter zu dem Fotoanzeige UI Element hinzugefügt. Mit dem Gitter kann die u,v Position von beliebigen Stellen des Fotos abgelesen werden. Durch aufnehmen von Fotos und vergleichen mit dem Sichtfeld des Displays wurde abgelesen, bei welcher u,v Position des Fotos sich die Ecken des Displays befinden. Die Intervalle wurden dem entsprechend eingegrenzt.

Mit den durchgenommenen Veränderungen der Intervalle konnten DetectedObjects korrekt in der Umgebung lokalisiert werden. Jedoch wurden sehr häufig Objekte nicht markiert, obwohl sie im Sichtfeld des Nutzers lagen, weil deren Mittelpunkt außerhalb eines Intervalls war.

Daher wurde entscheiden doch das gesamte Foto zu verwenden und Objekte auch zu Markieren, wenn sie komplett außerhalb des Sichtfeldes lagen. Dafür wurde die Intervalle für u und v wieder auf die ursprünglichen $[0,1920]$ und $[0,1080]$ gesetzt. Die Intervalle für x und y mussten vergrößert werden.

Um die x und y Intervalle bestimme zu können wurde das Fotoanzeige UI Element Parallel zu der ClippingPlane gelget. Das Element folgt den Bewegungen der Kamera und liegt so nah an der Near Clipping Plane wie es geht und noch angezeigt wird. Das Display der Magic Leap Brille zeigt selbst solide Objekte leicht durchsichtig an. Das wurde genutzt, um Fotos aufzunehmen, mit dem UI Element anzuzeigen und mit der Realen Welt zu vergleichen. Durch ausprobieren wurde das UI Element so skaliert und verschoben, das das Foto mit der realen Welt soweit wie möglich übereinstimmt.

Dann wurden die Ränder des UI Elementes genutzt um die Intervalle für x und y zu bestimmen.

- für x: $[-0.295, 0.2281]$
- für y: $[0.1546, -0.1507]$
- Zusätzlich wurde $z = -0.4$ gesetzt. Das Ui Element musste ein wenig weiter von der Clipping Plane entfernt werden um angezeigt zu werden.

Mit diesen Intervallen für u, v, x und y konnten DetectedObjects gut lokalisiert werden und es wurden keine Objekte mehr weggelassen, von denen der Nutzer erwarten würden, dass sie markiert werden.

6 Zusammenfassung

man kann zusätzlich zum raycast auf das mesh auch noch die tiefenkamera nutzen. spatial mapping hat manchmal das mesh an der stelle noch nicht, oder ist nicht aktualisiert. besonders probleme mit wasserflaschen die nicht erkannt werden. Hinzuziehen von tiefen daten kann das lokalisieren von objekten im 3D raum noch verbessern.

Literatur

- [1] Keiron O'Shea und Ryan Nash. „An Introduction to Convolutional Neural Networks“. In: *ArXiv e-prints* (Nov. 2015).
- [2] Licheng Jiao u. a. „A Survey of Deep Learning based Object Detection“. In: *IEEE Access* PP (Sep. 2019), S. 1–1. DOI: 10.1109/ACCESS.2019.2939201.
- [3] N. Jmour, S. Zayen und A. Abdelkrim. „Convolutional neural networks for image classification“. In: *2018 International Conference on Advanced Systems and Electric Technologies (IC ASET)*. 2018, S. 397–402.
- [4] *Microsoft Azure Computer Vision*. URL: <https://azure.microsoft.com/en-us/services/cognitive-services/computer-vision/> (besucht am 12.06.2020).
- [5] *What is Computer Vision*. URL: <https://docs.microsoft.com/en-us/azure/cognitive-services/computer-vision/home> (besucht am 12.06.2020).
- [6] *Detect common objects in images*. URL: <https://docs.microsoft.com/en-us/azure/cognitive-services/computer-vision/concept-object-detection> (besucht am 12.06.2020).
- [7] *Mr und Azure 302 Maschinelles Sehen*. URL: <https://docs.microsoft.com/de-de/windows/mixed-reality/mr-azure-302> (besucht am 12.06.2020).
- [8] *Mr und Azure 302b benutzerdefinierte Vision*. URL: <https://docs.microsoft.com/de-de/windows/mixed-reality/mr-azure-302b> (besucht am 12.06.2020).
- [9] *Spatial Mapping*. URL: <https://docs.microsoft.com/de-de/windows/mixed-reality/spatial-mapping> (besucht am 17.04.2020).
- [10] *magic leap 1*. URL: <https://www.magicleap.com/en-us/magic-leap-1> (besucht am 18.09.2020).
- [11] *Glossary and Usage*. URL: <https://developer.magicleap.com/en-us/learn/guides/glossary> (besucht am 18.09.2020).
- [12] *Lumin OS Overview*. URL: <https://developer.magicleap.com/en-us/learn/guides/lumin-os-overview> (besucht am 18.09.2020).
- [13] *Magic Leap Features*. URL: <https://developer.magicleap.com/en-us/learn/guides/magic-leap-features> (besucht am 18.09.2020).
- [14] *World Rekonstruktion*. URL: <https://developer.magicleap.com/en-us/learn/guides/world-reconstruction-overview-landing> (besucht am 18.09.2020).
- [15] *1.4 Spatial Meshing - Unity*. URL: <https://developer.magicleap.com/en-us/learn/guides/meshing-in-unity> (besucht am 18.09.2020).
- [16] *App Security*. URL: <https://developer.magicleap.com/en-us/learn/guides/application-security-overview> (besucht am 18.09.2020).
- [17] *Camera.cameraToWorldMatrix*. URL: <https://docs.unity3d.com/ScriptReference/Camera-cameraToWorldMatrix.html> (besucht am 18.09.2020).
- [18] *Matrix4x4.MultiplyPoint*. URL: <https://docs.unity3d.com/ScriptReference/Matrix4x4.MultiplyPoint.html> (besucht am 18.09.2020).