

Universität Augsburg
Fakultät für Angewandte Informatik
Institut für Software & Systems Engineering

Masterarbeit
Im Studiengang Informatik

Titel

Patrick Grüner

patrick.gmm91@gmail.com

Matrikelnummer: xxxxxxxx

Schneiderhansenweg 10

87700 Memmingen

4. Mai 2018



Betreuer: N.N.
Erstgutachter: Dr. Alwin Hoffmann

Inhaltsverzeichnis

1. Einleitung	2
1.1. Motivation	3
1.2. Ziel der Arbeit	4
1.3. Umfeld der Arbeit	5
1.3.1. CMORE Automotive GmbH	5
1.3.2. C.LABEL	5
2. Grundlagen	7
2.1. Annotation	7
2.1.1. Maschinelles Lernen	7
2.2. Virtual Reality	14
2.3. Aufnahme von Punktwolken mit LiDAR	17
3. Systembeschreibung	21
3.1. Komponente 1: Die VR-Brille	21
3.1.1. Warum VR die passende Technologie für C.LABEL-VR ist	22
3.1.2. Warum die Oculus Rift verwendet wurde	25
3.2. Komponente 2: Der Rechner	27
3.3. Komponente 3: Die Entwicklungsplattform	29
4. C.LABEL-VR	35
4.1. Import und Export von Daten	36
4.1.1. Architektur	37
4.1.2. HDF5-Beispiel	39
4.2. Generierung einer Punktwolke	39
4.2.1. Interne Datenstruktur	39
4.2.2. Generierung	39
4.2.3. Optimierung	39
4.3. Navigation	39
4.3.1. VR-Krankheit	39
4.3.2. Freier Flug	39

4.3.3. Teleport	40
4.4. Annotieren der Punktwolke	40
4.4.1. Einfache Pointer Annotation	40
4.4.2. Cluster Annotation	40
4.5. User Interface	46
4.5.1. Ingame Menu	46
5. Schluss	48
5.1. Zusammenfassung	48
5.2. Herausforderungen	48
5.3. Ausblick	48
A. Appendix Title	49

Abstract

Abstract goes here

Declaration

I declare that..

Acknowledgements

I want to thank...

Abbildungsverzeichnis

1.1.	Vogelperspektive auf eine LiDAR-Punktwolke, wobei die roten Punkte alle Boden- und die blauen alle Nicht-Bodenpunkte darstellen	3
1.2.	Labeling Methoden in C.LABEL	6
2.1.	Diagramm um den besten Machine Learning Algorithmus für eine Predictive Analytics-Methoden zu finden [1]	8
2.2.	Abbildung des Übertragungsprinzips von Reizen zwischen Neuronen [?]	10
2.3.	Einfaches Perzeptron nach dem Prinzip von Frank Rosenblatt	10
2.4.	Stark vereinfachte Darstellung eines KNNs zur Erkennung von Autos und Personen	12
2.5.	Vergleich zwischen den Anfängen der VR-Technologie und dem heutigen Stand. Die oberen drei Bilder zeigen das VIVED-System, das von S.S. Fisher vorgestellt wurde [2]. Die Unterer zeigen Bilder der Oculus Rift.	16
2.6.	Ablauf einer LiDAR-Messung von einem Sensor, der die Messimpulse im Uhrzeigersinn abgibt. Die Bilder stammen von [?] und 2.6b wurde um den Winkel φ erweitert.	18
2.7.	Seitliche Sicht auf eine LiDAR-Messung mit mehreren, vertikal angeordneten Lasern. Die Abbildung stammt aus der Arbeit [3] und wurde durch die Einzeichnung des Winkels θ erweitert.	19
3.1.	Die zwei potentiellen AR-Brillen	22
3.2.	Die zwei potentiellen VR-Brillen	26
3.3.	Unity Editor, der C.LABEL-VR im Play Mode ausführt.	31
4.1.	in C.LABEL-VR vom Import bis zum Export	36
4.2.	Das Import-Export-Prinzip in C.LABEL-VR	38
4.3.	Objektidentifizierung einer 3D-Punktwolke aus [4]. Die Bodenpunkte der Wolke sind blau dargestellt und die Objekte sind mit Boxen gekennzeichnet.	41

4.4. In [4] wird gezeigt, dass die Segmentierung einer Punktwolke zu einem besseren Ergebnis der Bodenpunktanalyse führt. Die obere Abbildung ist dabei ohne Segmentierung und die untere mit Segmentierung. Bodenpunkte sind blau und Objektpunkte in grün dargestellt. Die roten Rechtecke kennzeichnen den Bereich in dem die Analyse schlechter bzw. besser ist. . .	42
4.5. Vergleich der Segmentierungsmethoden zur Bodenpunktanalyse aus [4] und C.LABEL-VR. Bodenpunkte sind rot, alle anderen schwarz. Die blauen Linien stellen die Trennung der Segmente dar. Die grünen Boxen markieren den Bereich, in dem man den Unterschied der beiden Methoden bei der Bodenpunkterkennung sieht.	43
4.6. Bei zu wenigen Bodenpunkten in einem Segment passt sich die Ebene den falschen Punkten an.	43

Einleitung

Der Einsatz künstlicher Intelligenzen ist aktuell in allen Bereichen der Informatik auf dem Vormarsch. Dies gilt vor allem für die Automobilindustrie, da das Thema des autonomen Fahrens nicht ohne intelligente Algorithmen realisierbar ist. Eine wichtige Rolle bei der Entwicklung solcher Algorithmen spielt dabei das maschinelle Lernen. Dabei wird versucht, ausgehend von vielen lehrreichen Beispieldaten, die Lösung einer Aufgabe zu lernen und auf andere unbekannte Daten zu verallgemeinern. So kann ein System auch auf vorher ungewohnte Daten reagieren, was mit einer statischen Programmierung nur schwer oder gar nicht möglich ist. Der Erfolg dieses Prinzips hängt dabei stark von der Qualität der Beispieldaten ab von denen die Algorithmen lernen. Deshalb wird in die Erstellung dieser Daten sehr viel Arbeit gesteckt.

Im Bereich der Fahrerassistenzsysteme und des autonomen Fahrens ist es wichtig, dass das Auto seine Umgebung so gut wie möglich wahrnehmen kann. Algorithmen die für die Erkennung und Identifizierung des Umfeldes verantwortlich sind werden mit Daten von verschiedenen Sensoren gefüttert. Hierbei handelt es sich meist um Kamera-, Radar- oder LiDAR-Sensoren. Die Eingangsdaten dieser Sensoren müssen nun so aufbereitet werden, damit ein lernendes Computersystem etwas damit anfangen kann. Dazu werden alle, für das Anwendungsfeld wichtigen Teile mit entsprechenden Klassifikationen versehen. Diese Klassifikationen werden auch als *Labels* oder *Annotationen* bezeichnet. Zum Beispiel werden auf einem Kamerabild alle Personen und Fahrzeuge als eben diese markiert. Das Fahrzeug kann auf diese Weise lernen wie Personen und Fahrzeuge aussehen, um sie dann später im Straßenverkehr selbstständig wiederzuerkennen. Radar- und LiDAR-Sensoren liefern stattdessen Punktwolken als Eingangsdaten (vgl. Abbildung 1.1), das Prinzip bleibt allerdings das gleiche. Auch hier müssen alle wichtigen Teile, in diesem Fall Punkte, mit entsprechenden Klassen versehen werden. Dieser Vorgang nennt sich *Labeln* bzw. *Annotieren*.

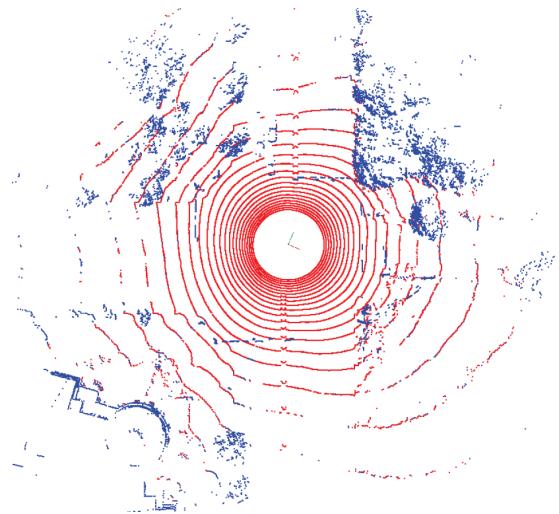


Abbildung 1.1.: Vogelperspektive auf eine LiDAR-Punktwolke, wobei die roten Punkte alle Boden- und die blauen alle Nicht-Bodenpunkte darstellen

1.1. Motivation

Der Fokus dieser Arbeit liegt auf der Annotation von Punktwolken. In einem bislang aufwändigen Verfahren (siehe Kapitel 1.3.2) müssen diese Punktwolken so aufbereitet und mit einer Bedeutung versehen werden, dass die Algorithmen an diesen Beispielen etwas lernen können. Sie müssen nach dem Lernen in der Lage sein Unterscheidungen und Erkennungen selbst durchzuführen und dabei möglichst wenige Fehler machen. Bei solchen Verfahren kommen Softwareanwendungen wie *C.LABEL* zum Einsatz. Dies ist eine Anwendung zum Annotieren verschiedener Sensordaten für den Automobilbereich und wurde von der Firma *CMORE Automotive GmbH* entwickelt. CMORE ist das betreuende Unternehmen dieser Masterarbeit und wird in Kapitel 1.3.1 vorgestellt. Auf das Programm *C.LABEL* wird in Kapitel 1.3.2 etwas genauer eingegangen. *C.LABEL* ist eine Anwendung, die zunächst die Sensordaten für den menschlichen Benutzer visualisiert. Dieser muss sie dann manuell mit unterschiedlichen Bedeutungen versehen. Wegen der großen Menge an Daten, die für das Trainieren der Algorithmen benötigt wird, muss dieses manuelle Annotieren möglichst effizient durchführbar sein.

Ein zweidimensionales Kamerabild kann sehr einfach auf einem Computermonitor dargestellt und bearbeitet werden. Eine besondere Herausforderung stellt jedoch die Verarbeitung von 3D-Daten, wie einer Punktwolke, dar. Hier stößt man mit den Möglichkeiten eines zweidimensionalen Computermonitors schnell an die Grenzen einer effizienten Darstellung und Bearbeitung. So ist es beispielsweise für die Annotierung solcher Punktwolken und ihrer Einzelmesswerte schwierig, die optimale Perspektive auf die Daten zu finden, die eine effiziente Erfassung durch den Menschen und eine entsprechende Bearbeitung zulässt. Dies erfordert von den Benutzern sehr viel Übung und Erfahrung, durch entsprechende

Drehungen der Punktwolkendarstellung sich in dieser zurechtzufinden und Messpunkte realen Objekten zuzuordnen.

Eine Alternative dazu soll die Applikation bieten, die innerhalb dieser Masterarbeit entwickelt wurde. Diese soll das Visualisieren und Annotieren von Punktwolken einfacher und effizienter gestaltet. Mittels *Virtual Reality* (VR) soll die Grenze zwischen 2-D und 3-D aufgehoben werden, sodass man sich als interaktiver Teilnehmer im dreidimensionalen Raum durch die Sensordaten navigieren, mit ihnen interagieren und sie annotieren kann. Die Applikation trägt den Namen *C.LABEL-VR*.

1.2. Ziel der Arbeit

Ziel dieser Arbeit ist es zunächst ein passendes Medium für eine solche Applikation zu finden, das nicht nur die Anforderungen der Applikation selbst erfüllt, sondern auch an einem handelsüblichen Büro-Arbeitsplatz verwendet werden kann. Anschließend soll eine Applikation erstellt werden die folgende Grundfunktionalität erfüllt:

1. Mit der Anwendung soll es möglich sein ein, bei CMORE gängiges, Datenformat einzulesen und alle nötigen Informationen daraus zu extrahieren
2. Aus den extrahierten Informationen soll eine Punktwolke generiert und innerhalb des gewählten Mediums visualisiert werden.
3. Der Benutzer muss die Möglichkeit haben durch die Punktwolke zu navigieren.
4. Die Applikation muss die Möglichkeit bieten jeden einzelnen Punkt mit einer Klassifikation versehen zu können. Dabei sollen Alleinstellungsmerkmale des gewählten Mediums benutzt werden, um eine vorzeigbare Verbesserung gegenüber der herkömmlichen Computer-Monitor-Annotation zu generieren.
5. Alle getätigten Annotationen müssen in die eingelesenen Dateien zurückgeschrieben bzw. in neue Dateien exportiert werden.

Anschließend sollen die Basisfunktionen erweitert und zusätzliche Funktionen hinzugefügt werden. Abhängig von der verbleibenden Zeit soll die Anwendung gemäß folgender Priorität erweitert werden:

1. Es sollen neue Möglichkeiten zur Annotation von Punkten implementiert werden.

2. Die Möglichkeit, eigene Klassifikationen für Punkte zu erstellen, soll hinzugefügt werden.
3. Neue Datenformate importieren.
4. Neue Navigationsmöglichkeiten implementieren.

Am wichtigsten wäre hierbei also die Applikation um effiziente Möglichkeiten zum Labeln von Punkten zu erweitern. Außerdem sollte es für den Benutzer möglich sein individuelle Klassifikationen anzulegen, um unterschiedliche Anwendungsgebiete abzudecken. Optional sind neue Navigationsmöglichkeiten und einlesbare Datenformate.

1.3. Umfeld der Arbeit

C.LABEL-VR und diese Masterarbeit wurden in Zusammenarbeit mit dem Unternehmen CMORE Automotive GmbH entwickelt. Als Grundlage für die Annotation von Punktwolken soll das CMORE-Programm C.LABEL dienen. Beides soll im Folgenden kurz vorgestellt werden.

1.3.1. CMORE Automotive GmbH

Im Mai 2011 wurde die CMORE Automotive GmbH von Richard Woller und Gregor Matenaer die CMORE Automotive GmbH in Lindau gegründet. Der Tätigkeitsbereich von CMORE liegt in der Unterstützung Entwicklung und Validierung zukünftiger Mobilitätskonzepte. Dabei bietet CMORE Lösungen und Produkte für alle Entwicklungsphasen. Dazu gehören beispielsweise komplexe Steuergerätesoftware für die Serienproduktion, Deep Learning Verfahren und Softwaretools zur hochautomatisierten Datenanalyse oder auch komplett Fahrzeugmesstechniksysteme zur Referenzdatengenerierung. Seit 2017 gibt es innerhalb des Unternehmens eine eigene Abteilung, die verantwortlich für die Distribution, Analyse und Anreicherung von Daten im Automobilbereich ist. Sie trägt den Namen C.IDS (*Integrated Data Solutions*). Ein Beispiel für solche Automobildaten sind Punktwolken von LiDAR-Sensoren, mit welchen sich auch diese Masterarbeit beschäftigt.

1.3.2. C.LABEL

C.LABEL ist ein Programm zu Annotierung von Daten im zwei- und dreidimensionalen Bereich. Im zweidimensionalen Bereich handelt es sich bei den Daten um Kamerabilder. In diesen Bildern werden wichtige Objekte mit sogenannten *Bounding Boxen* markiert, oder wichtige Bereiche werden zusammenhängend mit einer Klassifikation versehen (*Semantisches Labeln*). Im dreidimensionalen Fall werden Punktwolken eines Radar- oder LiDAR-Sensors annotiert. Hierbei müssen alle wichtigen Punkte mit einer Klassifikation

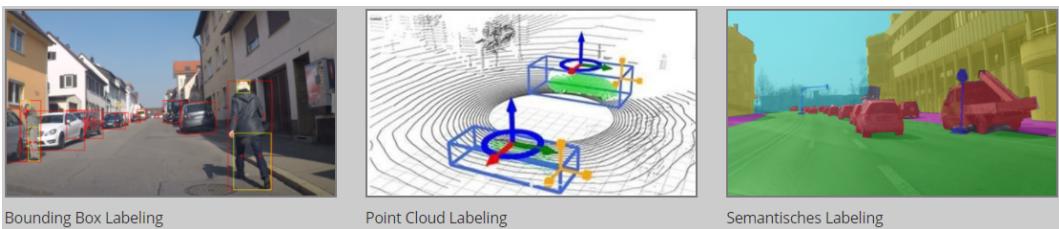


Abbildung 1.2.: Labeling Methoden in C.LABEL

versehen werden. Die Labeling Methoden sind in Abbildung 1.2 dargestellt. Das Ziel dieser Software ist es den Zeitaufwand, den die Annotation dieser Daten benötigt, zu minimieren. Dazu sind intelligente und effiziente Funktionen notwendig. Eine davon ist zum Beispiel die automatische Vorhersage von Annotationen durch stets analysierende Algorithmen. Dieses Prinzip wird *Active Learning Loop* genannt. Bei der Entwicklung von C.LABEL-VR liegt der Fokus auf der Annotierung von 3D-Punktwolken und dabei wurden Funktionen und Ansätze von C.LABEL als Basis für die Entwicklung genutzt.

Grundlagen

2.1. Annotation

TODO -Anntation im allgemeinen -Für welche zwecke wird es verwendet -Hinleitung zum maschinellen Lernen, für das Annotierte Daten notwendig sind

2.1.1. Maschinelles Lernen

Der Ausdruck Maschinelles Lernen (engl.: *Machine Learning*) ist heutzutage ein Überbegriff für Methoden, die versuchen durch Optimierung eine Funktion zu erlernen. Bei solch einer Funktion kann es sich sowohl um eine simple binäre Entscheidung auf eine Fragestellung handeln als auch um komplexere Dinge wie das Übersetzen von Texten in eine andere Sprache, das Erkennen von Personen auf Bildern, oder eben das Klassifizieren von Objekten in einer Punktwolke. Letzteres ist speziell für diese Arbeit relevant. Die Komplexität der Funktion bestimmt auch die zu verwendende Methode zur Lösung des Problems. Im Zuge des *Azure Machine Learning Studios* veröffentlichte Microsoft eine Graphik (vgl. Abbildung 2.2), mit der man den richtigen Lernalgorithmus für gewisse Vorhersage-Methoden bestimmen kann. Für die Objektklassifizierung im Automobilbereich ist es wichtig die Klasse eines Objekts mit hoher Genauigkeit aus einer großen Anzahl verschiedener Klassen zu bestimmen. Basierend auf der Graphik 2.2 ist die beste Methode dafür die Verwendung eines künstlichen neuronalen Netzes(KNN). Die Benutzung solcher Netze ist üblich, wenn es um Klassifikation von vielen verschiedenen Dingen geht und soll im Folgenden näher erläutert werden.

TODO Grafik von Bild- und Wolkenerkennung

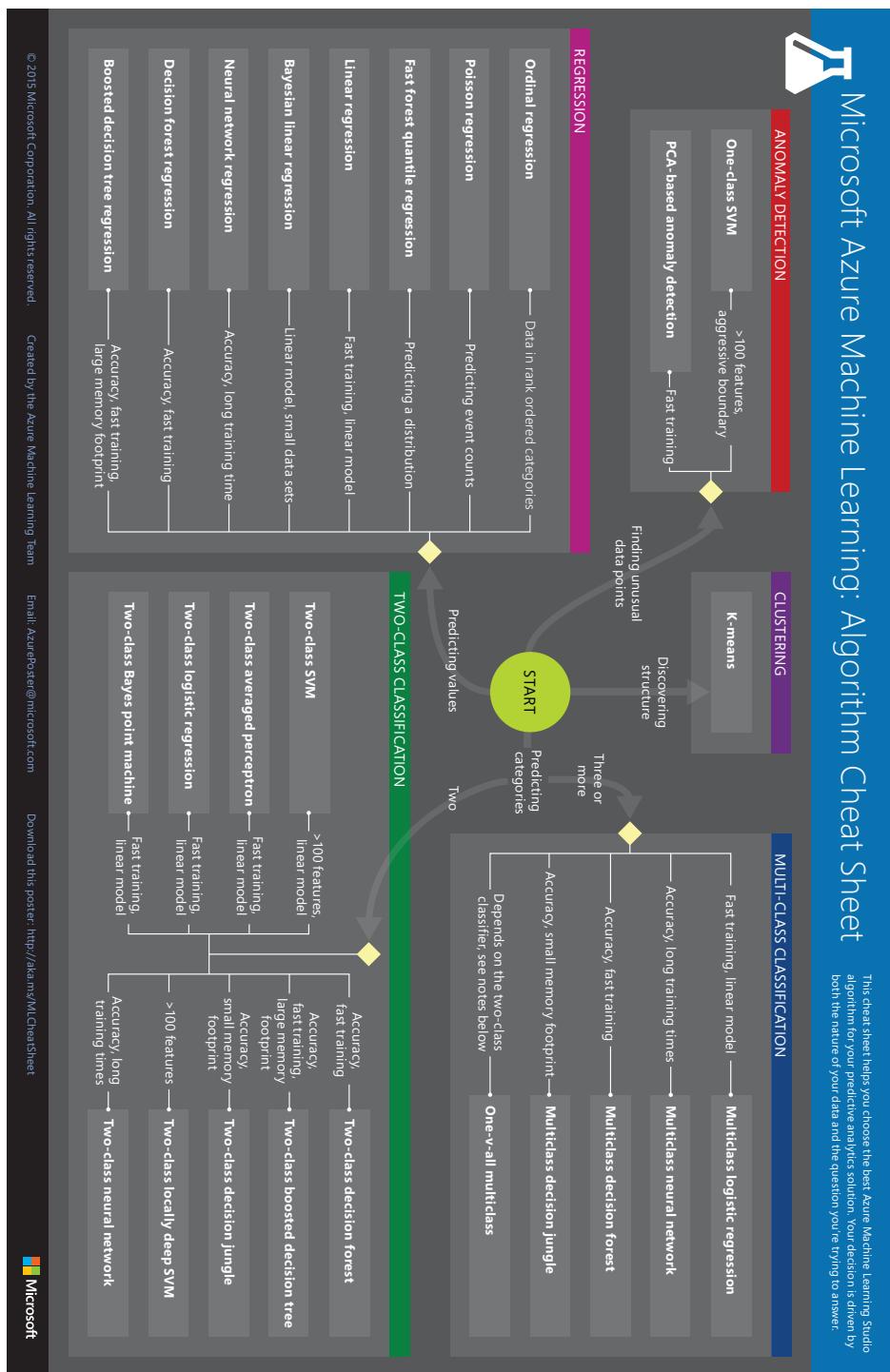


Abbildung 2.1.: Diagramm um den besten Machine Learning Algorithmus für eine Predictive Analytics-Methoden zu finden [1]

Künstliche Neuronale Netze

David Kriesel hat in [5] den Zusammenhang zwischen physischen und künstlichen neuronalen Netzen dargestellt. Seine Arbeit soll im Folgenden als Quelle dienen. Das Prinzip der künstlichen neuronalen Netze ist, wie vieles in der Informatik, einer Struktur aus der realen Welt nachempfunden. Das komplexeste neuronale Netz ist wohl das menschliche Gehirn mit etwa 86 Milliarden Nervenzellen, welche auch Neuronen genannt werden. Daher kommt auch der Begriff des neuronalen Netzes. Diese Nervenzellen sind durch Synapsen mit bis zu mehreren Tausend anderen Zellen verbunden. Diese Verbindungen bilden zusammen das Netz. Die Kommunikation zwischen den Neuronen erfolgt in der Regel über elektrische Impulse und chemische Botenstoffe. Beides wird über Synapsen zu den nächsten, verbundenen Neuronen weitergeleitet und bewirkt dort einen Reiz. Durch Übertragung von chemischen Botenstoffen kann dieser Reiz zusätzlich gewichtet werden, das heißt es kann ein stärkerer oder schwächerer Reiz ausgesendet werden. Alle ankommenden Reize werden dann am Empfangsneuron akkumuliert und ergeben den Eingangsreiz des Neurons.

Dieses, vereinfacht dargestellte Prinzip, wird auch bei den KNNs verwendet. Übersteigt dieser Reiz eine bestimmte Schwelle kommt es zum entscheidenden Ereignis: Das Neuron *feuert*. Dadurch gibt das feuernende Neurone einen Impuls an alle anderen Neuronen weiter mit dem es verbunden ist. Bei dafür vorgesehenen Neuronen oder verbänden davon, kann dieser Impuls auch ein Ereignis auslösen. Diese Ereignisse können physischer Natur sein, also beispielsweise Bewegungen, oder auch psychischer Natur, zum Beispiel das Erkennen von Dingen. Werden diese Nervenzellen also mit einem ausreichenden Impuls angesprochen, wird das entsprechende Ereignis ausgelöst.

Die entscheidende Komponente dieses Prinzips ist die adaptive Gewichtung der weitergegebenen Reize. Adaptiv deswegen, da die Gewichtung der Reize sich während eines Menschenlebens verändern können. So erkennen Kinder oft Gefahrensituation nicht, da ihr Gehirn nicht ausreichend geschult ist um diese zu erkennen. In diesem Kontext bedeutet das, dass an die *Gefahren-Neuronen* kein ausreichend starker Reiz gesendet wird. Durch sammeln von Erfahrungen optimiert sich dieser Wert, was in späteren Jahren zu einer anderen Reaktion des Gehirns führen kann als früher. Das Adaptieren dieser Gewichte auf einen passenden Wert bezeichnen wir als Lernen. Um Computersystemen einen derartigen Lerneffekt zu ermöglichen wurde dieses Biologische Prinzip nun mathematisch imitiert.

Künstliche Neuronale Netze basieren auf dem Prinzip des Perzeptrons, das von Frank Rosenblatt 1958 vorgestellt wurde [6]. Dieses Perzeptron repräsentiert die Funktion eines einzelnen Neurons auf eine mathematische Weise (vgl. 2.3). Die Eingangsdaten eines Perzeptrons bestehen aus einer n -großen Menge an Werten $x_1 \dots x_n$, welche eine gleich-große Menge an jeweiligen Gewichten $W_1 \dots W_n$ haben. Wenn das zu analysierende Ziel

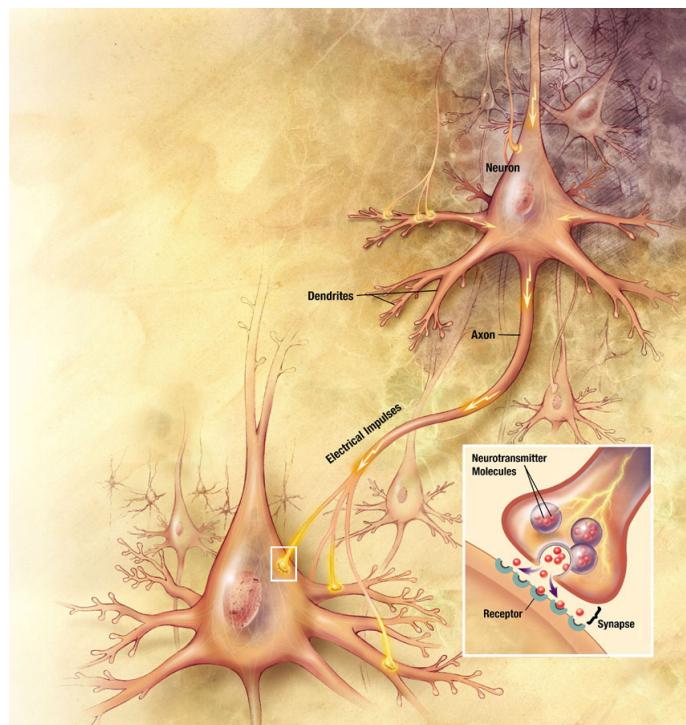


Abbildung 2.2.: Abbildung des Übertragungsprinzips von Reizen zwischen Neuronen [?]

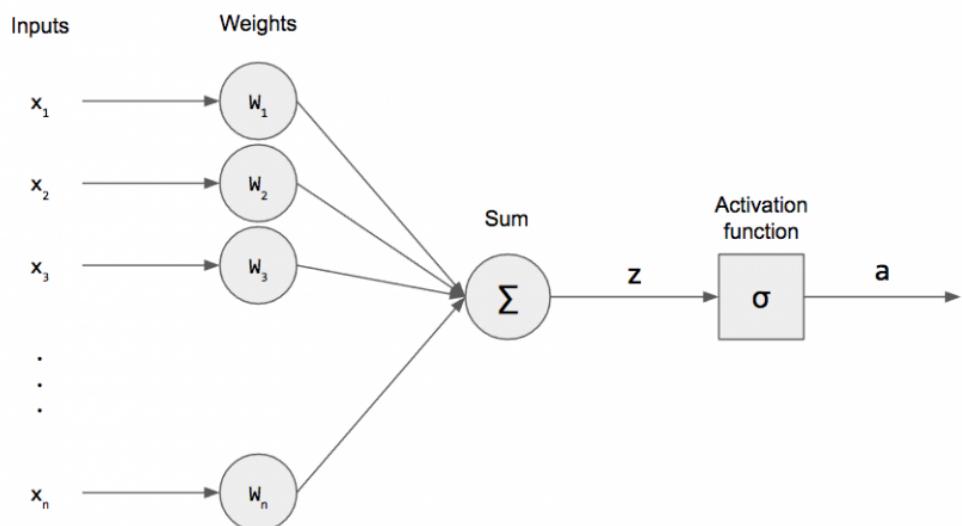


Abbildung 2.3.: Einfaches Perzeptron nach dem Prinzip von Frank Rosenblatt

beispielsweise ein Bild ist, wäre n die Menge aller Pixel des Bildes und bei Punktwolken eben die Anzahl der Punkte. Die Eingangswerte werden von einer mathematischen Funktion empfangen und zu einem skalaren Wert z zusammengefasst. Hierbei wird meist die Gewichtete Summe benutzt die folgendermaßen berechnet wird:

$$z = \sum_{i=0}^n W_i \cdot x_i \quad (2.1)$$

Im biologischen Vorbild repräsentiert diese Summe den eingehenden Reiz in eine Nervenzelle, die x -Werte sind dabei die Impulse der vorherigen Neuronen und die W -Werte repräsentieren die Gewichtung, also die Stärke der Eingangsreize. Der biologische Schwellwert wird im Perzeptron durch eine Mathematische Funktion $\sigma(z)$ dargestellt. Wie in [?] beschrieben gibt es dafür zwei beliebte Funktionen. Die *Fermifunktion* (2.2), auch Logische Funktion genannt, welche einen Ausgabe-Wertebereich von (0,1) hat und den *Tangens Hyperbolicus* ($\tanh(z)$) mit einem Wertebereich von (-1,1). Beide sind differenzierbar, was wichtig für das Lernen mit *Backpropagation* ist. Was das bedeutet wird in Abschnitt 2.1.1 näher erläutert. Wenn der Wert, den diese Funktion liefert, nicht noch durch eine weitere Ausgangsfunktion verändert wird, ist er der Ausgangswert des Perzeptrons.

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad (2.2)$$

In KNNs können mehrere dieser Perzeptrons eine sogenannte Schicht bilden. Je mehr Schichten ein neuronales Netz hat desto *tiefer* ist es, daher auch der Begriff *Deep Learning*, also tiefes Lernen. In der Regel bestehen KNNs aus mindestens 3 Schichten, nämlich einer Eingangsschicht, einer oder mehrerer Zwischenschichten und einer Ausgangsschicht (veranschaulicht in Abbildung 2.4). Die einzelnen Schichten werden im Folgenden kurz erläutert.

- **Eingangsschicht**

Die Eingabeschicht ist der Startpunkt des Informationsflusses in einem künstlichen neuronalen Netzes. Eingangssignale werden von den Neuronen dieser Schicht aufgenommen und gewichtet an alle Neuronen der ersten Zwischenschicht weitergegeben. Die Anzahl der Neuronen dieser Schicht hängt, wie schon erwähnt, von den Eingangsdaten ab. Werden Bilder mit einer Auflösung von 50x50 Pixeln in das Netz gespeist, hat dieses Netz in der Regel 50x50, also 2500 Neuronen in der Eingangsschicht.

- **Zwischenschicht**

Zwischen der Eingabe- und der Ausgabeschicht befindet sich in jedem künstlichen neuronalen Netz mindestens eine Zwischenschicht, die auch verborgene Schicht (engl.: *hidden layer*) genannt wird. Theoretisch ist die Anzahl der möglichen verborgenen Schichten in einem künstlichen neuronalen Netzwerk unbegrenzt. In der Praxis

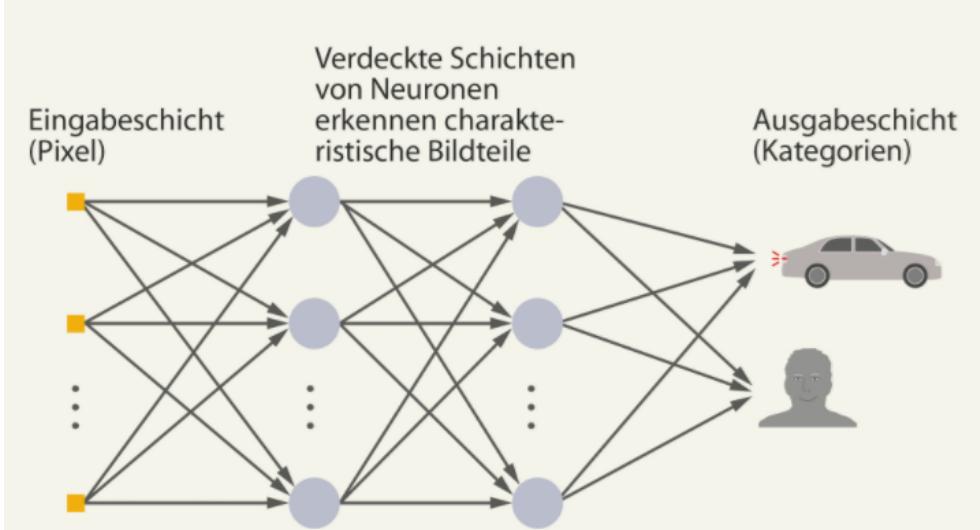


Abbildung 2.4.: Stark vereinfachte Darstellung eines KNNs zur Erkennung von Autos und Personen

bewirkt jede hinzukommende Schicht jedoch auch einen Anstieg der benötigten Rechenleistung, die für den Betrieb des Netzes notwendig ist.

- **Ausgangsschicht**

Die Ausgabeschicht liegt hinter den Zwischenschichten und bildet die letzte Schicht in einem künstlichen neuronalen Netz. In der Ausgabeschicht angeordnete Neuronen sind jeweils mit allen Neuronen der letzten Zwischenschicht verbunden. Die Ausgabeschicht stellt den Endpunkt des Informationsflusses in einem KNN dar und enthält das Ergebnis der Informationsverarbeitung durch das Netzwerk. Beim Beispiel der Objektklassifikation besteht die Ausgabeschicht nun aus n vielen Endknoten, wobei n die Menge an möglichen Klassifikationen ist. Die Eingangssumme dieser Knoten, den sie aus den vorherigen Neuronen bekommen, entsprechen der Wahrscheinlichkeit, dass es sich bei den Eingangsdaten um die jeweilige Klasse handelt. Abbildung 2.4 soll dieses Prinzip veranschaulichen.

Trainieren von künstlichen neuronalen Netzen

Damit ein Mensch eine neue Tätigkeit beherrscht muss er diese erst erlernen. In der Regel muss diese Tätigkeit dafür so oft es geht geübt werden. Bei künstlichen neuronalen Netzen ist das genauso. Das Üben bezeichnet man in diesem Kontext als *Trainieren* von Netzen. Im allgemeinen bedeutet der Vorgang des Lernens, dass ein System sich in irgendeiner Form verändert, um sich z.B. an Veränderungen in seiner Umwelt anzupassen. Die typische Veränderung von künstlichen neuronalen Netzen ist das adaptieren der Verbindungsge wichtete $W_1 \dots W_n$. Damit die Verbindungsgewichte auf einen passenden Wert verändert werden muss das Netz entsprechend oft trainiert werden. Dies geschieht nach Regeln, die

man in Algorithmen definiert – ein Lernverfahren ist also immer ein Algorithmus, den man einfach mithilfe einer Programmiersprache implementieren kann. Diese Algorithmen lassen sich grob in 3 verschiedenen Arten einteilen:

- Unüberwachtes Lernen (*Unsupervised Learning*)
- Bestärkendes Lernen (*Reinforcement Learning*)
- Überwachtes Lernen (*Supervised Learning*)

Für Objektklassifizierung wird in der Regel das überwachte Lernen verwendet, darum wird im Folgenden nur auf das dieses Verfahren eingegangen.

Der Lernprozess bei dieser Methode ist, wie der Name schon sagt, überwacht. Das kommt daher, dass der Mensch eine Menge an selbst ausgewählten Trainingsbeispielen erstellt. Das wird dann mit dieser überwachten Menge trainiert. Die Trainingsbeispiele bestehen aus einer Menge P an Eingabemustern mit jeweiliger korrekter Lösung. Das Netz kann seine eigene Ausgabe mit der richtigen Lösung des Beispiels vergleichen und somit eine Fehlerbeschreibung zurückgeben. Für das Beispiel der Punktwolkenklassifizierung bestünde dabei das Trainingsset aus vielen verschiedenen Punktwolken, die alle ein Objekt darstellen, wobei die Art des Objektes, also die Lösung der Klassifizierung, bekannt ist.

Die Werte der Gewichte $W_1 \dots W_n$ werden bei einem untrainierten Netz zu Beginn zufällig gewählt. Dadurch gibt ein solches Netz in der Regel völlig falsche Lösungen für die ersten Trainingsbeispiele aus. Das ist aber nicht weiter schlimm, da nun der Lerneffekt einsetzen soll. Durch die Bekanntheit der Lösung des Beispiels kann der Fehler zwischen dem Ergebnis des Netzes und der gewünschten Lösung genau bestimmt werden. Dieser Fehler kann anschließend durch zurückrechnen auf jedes Neuron der vorherigen Schichten projiziert werden. Auf diese Weise kann der Fehler, den jedes einzelne Neuron verursacht bestimmt werden. Dieses Zurückrechnen nennt sich *Backpropagation*. Aus dem Fehler den jedes Neuron verursacht kann schließlich der Betrag errechnet werden, um den sich das Gewicht W_i des Neurons verändern muss, damit bei zukünftigen Eingaben bessere Ergebnisse erzielt werden. Für das Verständnis des Grundprinzips ist die Formel zu Berechnung der Gewichtsänderung nicht wichtig, sie wird aber zur Vervollständigung im Folgenden angegeben:

$$\Delta w_{ij} = -\eta \delta_j o_i \quad (2.3)$$

- Δw_{ij} ist die Änderung des Gewichts w_{ij} der Verbindung von Neuron i zu Neuron j
- η ist ein fester Wert, der eine Lernrate ausdrücken soll, mit der die Stärke der Gewichtsveränderung festgelegt wird

- δ_j ist das Fehlersignal des Neurons j das sich aus der partiellen Ableitung $\frac{\partial E}{\partial net_j}$ ergibt
 - E ist dabei der quadratische Fehler der bei der Berechnung des Ergebnisses entstand, er wird durch die Formel $E = \sum_{i=1}^n (t_i - o_i)^2$ berechnet
 - * t_i steht darin für den gewünschten Zielwert
 - net_j ist die Netzeingabe $\sum_{i=1}^n x_i w_{ij}$
- o_i steht sowohl bei der Berechnung von Δw_{ij} , als auch von E für die errechnete Ausgabe des Neurons i

Zusammenfassend kann man nun sagen, dass der Lernprozess eines künstlichen Neuronalen Netzes nach folgendem Schema abläuft:

1. **Eingabe** eines Elements p aus der Trainingsmenge P in die Eingabeschicht des Netzes
2. **Vorwärtspropagierung** der Eingabe durch das Netz gesamte Netz, sodass man bei der Ausgabeschicht ein Ergebnis erhält
3. **Vergleich** des Ergebnisses mit der richtigen Lösung und Berechnung des Fehlers
4. **Zurückpropagieren** des Fehlers, um somit die Gewichte der Neuronen anzupassen

Wie unschwer zu erkennen ist basiert der Backpropagation-Algorithmus auf dem Vorhandensein eines Trainingsdatensatzes. Damit das Netz nach dem Training für neue, unbekannte Eingaben die richtige Lösung errechnet ist es elementar wichtig, dass dieser Datensatz eine hohe Qualität und Quantität hat. Qualität bedeutet in diesem Fall, dass die Trainingsdaten repräsentativ für den jeweiligen Anwendungsfall sein müssen. Ein Personenerkennungsalgorithmus der im Straßenverkehr Personen identifizieren soll liefert keine guten Ergebnisse, wenn er ausschließlich mit Bildern von Personen in Innenräumen gefüttert wird. Mit Quantität ist Anzahl der verschiedenen Elemente im Trainingssatz gemeint. KNNs verbessern ihre Ausgabe in der Regel mit größerer Menge an unterschiedlichen Trainingsdaten.

Zur Erstellung dieses Datensatzes ist es nun nötig so viele repräsentative Daten wie möglich zu sammeln und diese mit der richtigen Lösung zu annotieren. Eine Methode für eine solche Annotierung wird in dieser Arbeit vorgestellt. Dieses Kapitel sollte nun anschaulich haben warum solche Methoden so essentiell sind, um künstliche Intelligenzen zu entwickeln.

2.2. Virtual Reality

Das Konzept digitale Inhalte räumlich und interaktiv darzustellen ist nicht neu. Die Ursprünge von VR lassen bis 1965 zurückverfolgen als Ivan Sutherland „The Ultimate Display“ [7]

vorstellte. Diese Arbeit stellte die ersten Konzepte für Immersion in einer virtuellen Welt und die dafür nötigen Ein- und Ausgabegeräte dar. Er wollte damit einen Anreiz schaffen um die Grenzen der damaligen Mensch-Computer-Schnittstelle zu überwinden und neue Applikationen zu entwickeln, welche die eigene Präsenz in der virtuellen Welt miteinbeziehen [8].

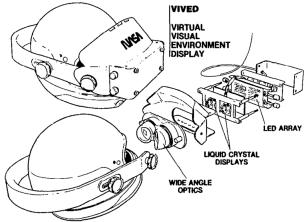
Anfang der achtziger Jahre gab es dann schon funktionierende VR-Systeme, beispielsweise das VIVED System (*Virtual Visual Environment Display*), welches von der NASA entwickelt wurde [2]. Dieses System beinhaltet eine, am Kopf befestigte Brille, mit einem integrierten 2,7" großen LCD-Display (Abbildung 2.5a). Durch Weitwinkelobjektive wird das Bild dieses Displays auf jedes Auge einzeln projiziert. Damit war es möglich Objekte, die normalerweise an einem offenen Bildschirm angezeigt werden, in einer geschlossenen und 3D-ähnlichen Sicht darzustellen (Abbildung 2.5b).

Die Projizierung auf jedes Auge vermittelt dabei einen räumlichen Eindruck, der durch das Prinzip der Stereoskopie verursacht wird [?]. Dieses Prinzip beruht darauf, dass durch zwei Augen die Umgebung aus zwei verschiedenen Blickwinkeln wahrgenommen wird. Dies kann leicht selbst festgestellt werden, indem man einen Finger zwischen die Augen hält und abwechselnd jedes Auge einzeln öffnet. Dabei ist zu erkennen, dass der Finger sich weiter rechts oder links im Blickfeld befindet, je nachdem welches Auge geöffnet ist. Durch die Zusammensetzung der zwei Blickwinkel durch das Gehirn entsteht ein für uns ein Tiefeneffekt. Dieser Effekt wird auch von VR-Brillen genutzt um 2D-Inhalte räumlich wirken zu lassen.

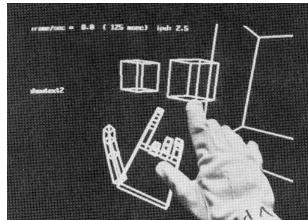
Als Eingabegerät wurde sowohl eine Spracherkennung, als auch ein Handschuh benutzt, der Bewegungssensoren an jedem Finger hatte. Somit war es dem Benutzer möglich auf eine intuitive Weise mit dem System zu Interagieren (Abbildung 2.5c), da das System die Bewegungen des Benutzers erkannte.

Am Prinzip der VR-Technologie hat sich bis heute nicht viel verändert. Die Präsentation der virtuellen Welt geschieht immer noch durch eine Brille mit integrierten Linsen und einem Display (Abbildung 2.5d). Diese beinhaltet, ebenso wie die Eingabegeräte, Bewegungssensoren, die die Position und Orientierung dieser Komponenten misst.

Die funktionelle Leistung der Komponenten ist im Laufe der Zeit jedoch deutlich besser geworden. Statt eckigen Polygonen können nun hochauflösende 3D-Modelle durch ein Display mit 2160 x 1200 Pixeln dargestellt werden (Abbildung 2.5e). Die Brille und die Eingabegeräte werden nun auch, zusätzlich zu den internen Bewegungssensoren, von externen Infrarot-Sensoren erfasst, um eine präzise Portierung der physischen Bewegungen in virtuelle zu übersetzen. Der größte Unterschied ist wohl, dass die VR-Technologie heutzutage hauptsächlich für Unterhaltungszwecke genutzt wird. Was früher mal ein Forschungsprojekt einer Raumfahrtorganisation war ist heute ein Entertainment-System, das



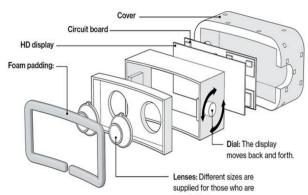
(a) VIVED Displayaufbau



(b) Sicht eines VIVED Benutzers



(c) Sicht auf einen VIVED Benutzer



(d) Rift Displayaufbau



(e) Sicht auf einen Rift Benutzer



(f) Sicht eines Rift Benutzer

Abbildung 2.5.: Vergleich zwischen den Anfängen der VR-Technologie und dem heutigen Stand. Die oberen drei Bilder zeigen das VIVED-System, das von S.S. Fisher vorgestellt wurde [2]. Die Unteren zeigen Bilder der Oculus Rift.

in viele Privathaushalten zu finden ist.

Ein Grund dafür war die Entwicklung der zwei, derzeit geläufigsten VR-Brillen *Oculus Rift* und *HTC Vive*. Ein paar geschichtliche Eckdaten davon wurden in [9] gesammelt. 2010 wurde der erste Prototyp der Oculus Rift von Palmer Luckey entworfen. Dieser hatte eine horizontale Sichtweite von 90°, was zur damaligen Zeit im Verbrauchermarkt einzigartig war. 2013 hatte Valve einen Durchbruch bei der Entwicklung von *Low-Persistence*-Displays, den sie öffentlich geteilt haben. Dies war sehr hilfreich für alle Entwickler, denn Oculus beispielsweise benutzte seitdem bei all ihren nachfolgenden Brillen. Low-Persistence (zu deutsch: *geringe Ausdauer*) bedeutet in dem Falle die angezeigten Bilder nur wenige Millisekunden anzuzeigen um somit Objekte, die sich Bewegen schärfer darstellen zu können. 2014 wurde Oculus Rift für zwei Milliarden Dollar von Facebook gekauft und beschäftigen mittlerweile circa 400 Mitarbeiter allein für die Virtual Reality Entwicklung. Daran kann man das Potential sehen, das große Konzerne in der VR-Technologie sehen, denn kurz danach kündigte nämlich auch Sony eine VR-Brille für die Playstation 4 an.

Die, aus der Großproduktion der führenden VR-Konzerne resultierende, Massentauglichkeit und die Erschwinglichkeit der Technologie machen diese aber auch für die Industrie interessant. Vor allem im Zuge der *Industrie 4.0* ist VR für Digitalisierung des industriellen Umfelds von Nutzen. Es werden Leistungsfähige und qualitativ immer bessere Lösungen zur virtuellen Datenerfassung und -visualisierung entwickelt. Mittels VR lassen sich Arbeitssituationen simulieren und gefahrlos testen. Industriearbeiter können in Umgebungen

agieren, die in der Form noch gar nicht existieren. So können zum Beispiel Arbeitsumfelder und -prozesse effizient geplant werden. Schadens- und Störfälle lassen sich dank VR schnell und in Gänze erfassen und aus der Ferne warten. Es können auch Dinge räumlich begehbar gemacht werden, die es in der physischen Welt nicht gibt, wie beispielsweise Sensordaten. Im Rahmen dieser Arbeit werden Daten eines Lidarsensors als Punktwolke dargestellt, durch die man sich dann Bewegen kann.

2.3. Aufnahme von Punktwolken mit LiDAR

Punktwolken, wie sie in dieser Arbeit behandelt werden, bestehen aus vielen einzelnen Abstandsmessungen, wobei jeder Punkt einer dieser Messungen entspricht. LiDAR (*Light Detection and Ranging*) ist eine Methode zur Abstandsmessung mit Hilfe von Licht, welche in der Automobilbranche häufig verwendet wird. In den letzten Jahrzehnten wurde der LiDAR-Sensor aber auch in vielen anderen Bereichen zur herkömmlichen Methode, wenn es um die Messung von Abstandsinformationen geht. Im Bereich Archeologie benutzten Johnson und Oiumet diese Technologie um die Landschaft von New England in den USA zu analysieren [10]. Die Analyse von Landschaften ist durch LiDAR aber nicht nur auf dem Planeten Erde möglich. In [11] wurde LiDAR benutzt um Oberflächenmessungen auf dem Mars durchzuführen. In der Robotik wird es häufig verwendet um Hindernisse für mobile Roboter zu identifizieren, wie in [12].

Für die Messung der Distanz wird in der Regel ein Laser benutzt. Der LiDAR-Sensor sendet mit Hilfe eines solchen Lasers einen Lichtimpuls aus und misst die Zeit t bis das Signal wiederkehrt. Anhand dieser Zeit ist es möglich die Distanz d zu bestimmen, die vom Licht zurückgelegt wurde. Dazu nimmt man die Formel 2.4, in der c für die Geschwindigkeit des Lichts steht.

$$d = \frac{c \cdot \Delta t}{2} \quad (2.4)$$

Mit dem Impuls eines einzelnen Lasers bekommt man eine einzelne Abstandsmessung für den Zeitpunkt des Impulses. Führt man aber viele dieser Messungen innerhalb einer 360° Drehung durch bekommt man Abstandsinformationen des gesamten Umfeldes. Eine Möglichkeit eine 360° Umfeldmessung zu ermöglichen ist es den Sensor um seine eigene Achse zu drehen, wie es in Abbildung 2.6 gezeigt wird. In der oberen Reihe dieser Abbildung ist der LiDAR-Sensor zu sehen, dessen Lichtimpuls durch einen Spiegel (graues Objekt) gezielt auf die Umgebung gerichtet wird. Dieser Spiegel dreht sich um die eigene Achse um die Lichtimpulse nach und nach auf das Umfeld zu verteilen. In der mittleren Reihe ist die Vogelperspektive auf den Sensor und seine Umgebung abgebildet. Der Sensor ist, wie in der oberen Abbildung, in blau dargestellt und die Objekte des Umfeldes in grün. Die unteren Abbildungen zeigen, ebenfalls in der Vogelperspektive, die Abstandsinformationen, die

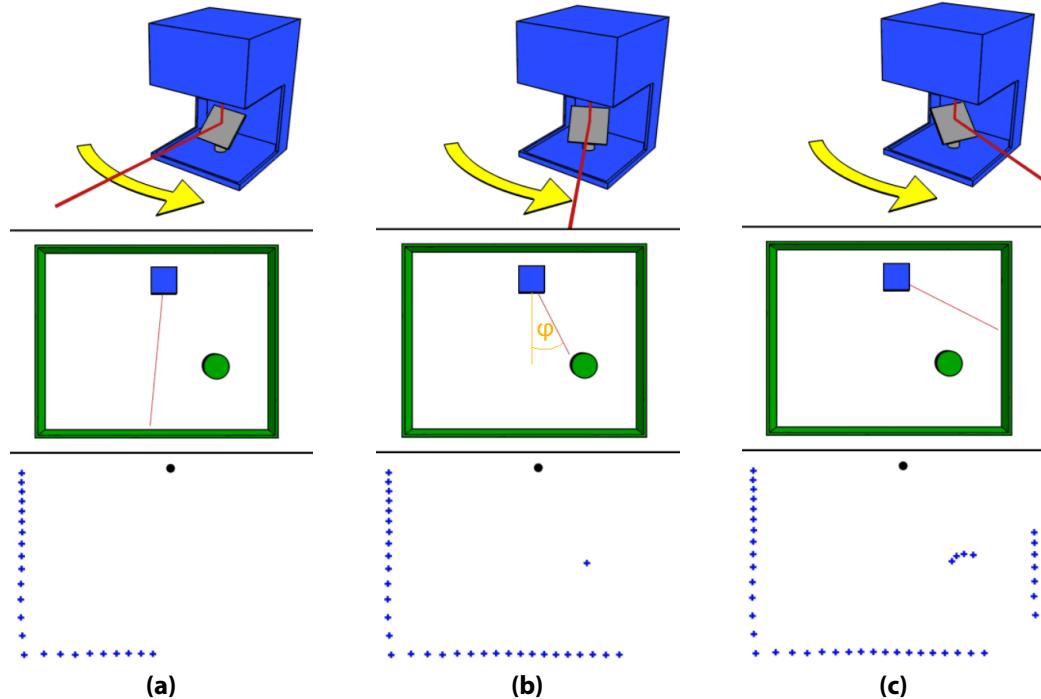


Abbildung 2.6.: Ablauf einer LiDAR-Messung von einem Sensor, der die Messimpulse im Uhrzeigersinn abgibt. Die Bilder stammen von [?] und 2.6b wurde um den Winkel φ erweitert.

der Sensor liefert. Diese werden in Form von blauen Punkten dargestellt und bilden eine zweidimensionale Punktwolke, in der man den Rahmen des Umfelds und das kreisförmige Objekt erkennen kann.

Diese Methode liefert also eine gute horizontale Rundumsicht, welche aber innerhalb der vertikalen Sichtweite keinerlei Informationen bietet. Dies genügt in der Regel für simple Abstandserkennungen und Detektion von Hindernissen auf Höhe des Sensors. Möchte man jedoch Aussagen über das detektierte treffen oder braucht Informationen über andere vertikalen Ebenen, beispielsweise den Boden, muss man diese Methode noch erweitern. Dafür werden mehrere Lichtimpulse mit verschiedenen vertikalen Winkeln ausgesendet. Diese treffen die Umgebung dann mit unterschiedlicher Höhe womit die daraus resultierende Punktwolke eine dritte Dimension erhält. In Abbildung 2.7 ist die seitliche Sicht auf einen LiDAR-Sensor mit mehreren vertikalen Lasern zu sehen.

Wie schon erwähnt werden die erhaltenen Messdaten in eine Struktur konvertiert, die sich Punktwolke nennt. Dafür werden die Abstandsinformationen innerhalb eines bestimmten Koordinatensystems dargestellt. In der Regel wird dabei das *Sphärische Koordinatensystem* verwendet. Darin wird jeder Punkt durch einen Abstand r von einem Referenzpunkt und zwei dazugehörigen Winkel repräsentiert. Einer dieser Winkel ist dabei horizontal (Winkel φ) zu interpretieren und der andere vertikal (Winkel θ).

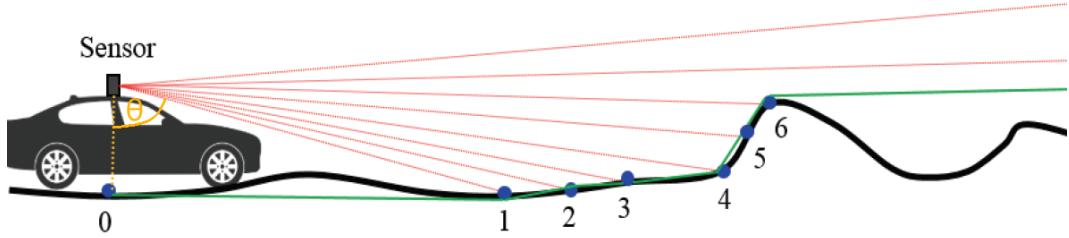


Abbildung 2.7.: Seitliche Sicht auf eine LiDAR-Messung mit mehreren, vertikal angeordneten Lasern. Die Abbildung stammt aus der Arbeit [3] und wurde durch die Einzeichnung des Winkels θ erweitert.

Nimmt man den Punkt „1“ aus der Abbildung 2.7 als Beispiel, entspräche θ dem vertikalen Winkel. Dieser gibt den Winkel zwischen der gelben, gestrichelten Linie und der roten Linie an, die den Sensor und „1“ verbindet. Diese rote Linie repräsentiert den Abstand des Punktes vom Referenzpunkt, welcher die Position des Sensors entspricht. Der horizontale Winkel ist der Rotationswinkel des Sensors um seine eigene Achse (vgl. φ in Abbildung 2.6b). Dieser wird vom Sensor bei jedem abgegebenen Impuls gemessen. Zusammenfassend kann also durch den vertikalen Winkel θ des Lichtimpulses, den horizontalen Winkel φ des Sensors und dem gemessenen Abstand r zum Sensor, der Punkt „1“ genau im Raum definiert werden.

Viele Anwendungen, die mit Punktwolken zu tun haben arbeiten jedoch mit dem kartesischen Koordinatensystem, so auch C.LABEL-VR. In diesem Koordinatensystem wird jeder dreidimensionale Punkt durch drei Koordinaten angegeben, nämlich der x -, y - und z -Koordinate. Ausgehend von den drei Werten des sphärischen Koordinatensystems können die drei kartesischen Koordinaten durch die Formeln in 2.5 berechnet werden. Auch hier steht r für den Abstand zwischen dem Referenzpunkt und dem Punkt den man definieren will. θ steht für den vertikalen Winkel und φ für den horizontalen Winkel, also den Rotationswinkel des Sensors um die eigene Achse.

$$x = r \sin \theta \cos \varphi \quad (2.5a)$$

$$y = r \sin \theta \sin \varphi \quad (2.5b)$$

$$z = r \cos \theta \quad (2.5c)$$

Durch die eben erklärte Definition der Messdaten in einem Koordinatensystem können diese Daten nun mit geeigneten Programmen (C.LABEL) als Punktwolke dargestellt werden. Darin ist dann jede einzelne Messung die der Sensor getätigt hat sichtbar. Somit

kann der Mensch die Daten visuell analysieren und beispielsweise Zusammenhänge wie Objekte entdecken. Damit auch Computersysteme solche Objekte in diesen Daten erkennen können müssen sie erst lernen wie diese aussehen. Dazu muss eine Vielzahl an Punktwolken mit Hilfe eines LiDAR-Sensors aufgenommen werden und anschließend eine, vom Anwendungsfall abhängige, Anzahl an Punkten in dieser Wolke mit entsprechenden Klassifikationen versehen werden (Prinzip der Annotation aus Kapitel 4.4). Mit den annotierten Daten können nun Algorithmen, wie beispielsweise künstliche neuronale Netze (vgl. Kapitel 2.1.1), lernen welche Zusammenhänge es in solchen Daten gibt und diese dann in nicht annotierten Daten wiedererkennen. Das Ziel von C.LABEL-VR ist es solche Daten in der virtuellen Realität zu visualisieren und Funktionen zur Verfügung zu stellen diese in der virtuelle Umgebung zu annotieren.

Systembeschreibung

Das Hauptthema mit dem sich diese Arbeit beschäftigt ist das Virtual Reality Programm C.LABEL-VR. Damit dieses Programm entwickelt und betrieben werden kann bedarf es aber einiger erwähnenswerter Komponenten, welche das System um die Applikation herum repräsentieren. Dazu gehört vor allem eine Virtual Reality Brille. Warum hierfür die Oculus Rift genommen wurde und warum Virtual Reality für solch eine Applikation geeigneter ist als vergleichbare Technologien wird in Kapitel 3.1 erläutert. Der Computer der für den Betrieb der VR-Brille benutzt wurde und was bei der Auswahl dessen Komponenten wichtig war ist in Kapitel 3.2 beschrieben. Mit den eben genannten Komponenten lässt sich eine VR-Applikation zwar ausführen, jedoch nicht entwickeln. Dafür ist zusätzlich eine Entwicklungsumgebung notwendig, genauer gesagt eine Spiel-Engine (*Game Engine*). Welche Engine in dieser Arbeit verwendet wurde und welche grundlegenden Konzepte dieser wichtig für die Entwicklung waren, wird in Kapitel 3.3 beschrieben.

3.1. Komponente 1: Die VR-Brille

Wie schon erwähnt, war die praktische Hauptaufgabe dieser Arbeit die Realisierung einer Applikation zur Annotation von Punktwolken in einer dreidimensionalen Umgebung. Hierfür wurde die Virtual Reality Technologie verwendet, welche im Kapitel 2.2 näher erläutert wird. Diese ist jedoch nicht die einzige Technologie die für diesen Anwendungsfall in Frage kam. Auch die *Augmented Reality* (AR) bietet einige interessante Konzepte um solch eine Applikation zu entwickeln. Warum hierbei aber VR zu bevorzugen ist, wird im folgenden Abschnitt 3.1.1 erläutert. Im anschließenden Abschnitt 3.1.2 wird gezeigt welche VR-Brillen für die Entwicklung in Frage kamen und warum die Oculus Rift aus diesen ausgewählt wurde.



(a) Microsoft Hololens AR-Brille

(b) Meta 2 AR-Brille

Abbildung 3.1.: Die zwei potentiellen AR-Brillen

3.1.1. Warum VR die passende Technologie für C.LABEL-VR ist

Um das Labeling von Daten, insbesondere Punktwolken, im dreidimensionalem Raum zu verwirklichen, gibt es zwei wesentliche Technologien die dafür relevant erscheinen. Neben VR, das schon vorgestellt wurde (Kapitel 2.2), gibt es noch die Augmented Reality Technologie. Unter AR versteht man die direkte oder indirekte Sicht auf die reale, physische Welt, welche durch digitale Inhalte erweitert wird. Dies wird vor allem durch Smartphones oder AR-Brillen realisiert. Bei einem Smartphone hat man beispielsweise einen indirekten Blick auf die reale Umgebung durch das Display, welches ein Live-Bild oder ein aufgenommenes der Kamera anzeigt. Diesen Bildern können nun digitale Inhalte hinzugefügt werden. Bei einer Brille (vgl. Abbildung 3.1) sieht man durch die Gläser direkt auf die physische Welt. Hierbei fungieren die Gläser als Display, welche in der Lage sind holographische Objekte anzuzeigen. Dies führt zu einer, für den Benutzer, sehr immersiven Vermischung der realen und der digitalen Welt. Üblicherweise ist es bei diesen Brillen sogar möglich durch Gesten, die mittels Händen durchgeführt werden, mit den Hologrammen zu interagieren.

Eignung von Augmented Reality für 3D-Annotation

Der, im Rahmen dieser Arbeit, entwickelte Prototyp zur 3D-Datenannotation ist nicht nur als Nutzungssoftware geplant, sondern dient auch als Anschauungsmaterial, beispielsweise für Messen. Zieht man diese Tatsache in Betracht, eignet sich Augmented Reality sehr gut für diesen Zweck. Es wirkt durch die Vermischung von realen und digitalen Inhalten nämlich sehr futuristisch. Zudem ist das direkte einblenden holographischer Inhalte, durch eine AR-Brille, zum Zeitpunkt der Erstellung dieser Arbeit, noch wenig verbreitet (Microsoft Hololens wurde nur wenige Tausend mal verkauft [13]), was bei unwissenden Benutzern zu einem beeindruckenden Effekt führt. Des Weiteren bietet die Augmented Reality Technologie viele Möglichkeiten, die für zukünftige Labeling-Methoden relevant sein könnten. Den heutigen AR-Brillen ist es durch Tiefenkameras möglich Objekte und deren Distanzen zur Brille wahrzunehmen. Folglich könnte ein zukünftiger Ansatz sein, eine Punktwolke der Umgebung, wie sie in Abbildung ?? zu sehen ist, mit einer AR-Brille zu erstellen. Diese kann anschließend vor Ort, in der Umgebung in der die Punktwolke

erstellt wurde, klassifiziert werden.

Ein wichtiger Punkt der ebenfalls betrachtet werden muss ist die Bedienung der Brille bzw. die Interaktion mit den digital dargestellten Objekten und Informationen. Die Steuerung von AR-Brillen erfolgt handelsüblich über Gesten, welche mit der Hand bzw. den Händen getätigten werden. Ob diese Art der Interaktion für den hier gewünschten Anwendungsfall passend wäre, lässt sich im Voraus schwer feststellen. Denkbar wäre, dass sich der Benutzer durch intuitive Gesten, die man auch bei realen Objekten nutzt, schnell an die Bedienung des Tools gewöhnen würde. Andererseits könnte die Selektion der Elemente einer Punktwolke zu ungenau sein, da die Brille die Position der Hand nicht richtig interpretiert bzw. diese nicht richtig zur Position der digitalen Inhalte interpretiert. Eine genaue Einschätzung der Bedienung für solch einen Anwendungsfall kann allerdings nur gegeben werden indem man einen Labeling-Prototypen erstellt und ihn anschließend über längere Zeit testet. Dies ist im Zeitrahmen einer Masterarbeit jedoch nicht möglich. Falls sich nämlich die Steuerung als ungeeignet herausstellt bleibt nicht genug Zeit für eine Neuentwicklung auf einer anderen Plattform.

Darüber hinaus gibt es natürlich auch Aspekte die gegen die AR-Technologie sprechen. Die Klassifizierung einer Punktwolke mit dieser Technik ist an einem normalen Büro-Arbeitsplatz nicht möglich, zumindest nicht in einem Maßstab in dem es sinnvoll wäre. Die holographischen Punkte kollidieren, wenn sie darauf programmiert sind, mit der Umgebung und so würde Darstellung des Umfeldmodells verfälscht werden. Auch wenn sie das nicht tun, dann wäre die Umgebung selbst immer noch ein Hindernis für den Benutzer, denn man möchte sich ja durch die Punktwolke bewegen.

Des Weiteren sind bei AR-Brillen die vorherrschenden Lichtbedingungen ein großer Faktor, welche die Funktionalität gewisser Anwendungsfälle beeinflussen können. Wird beispielsweise eine Lichtquelle zu sehr von einem Objekt reflektiert, kann dieses von der Tiefenkamera der Brille nicht mehr richtig erfasst werden. Somit stimmt das Umfeldmodell nicht mehr mit der Realität überein. Zudem wird nicht nur die Funktionalität von Lichteinflüssen gestört, sondern auch die Darstellung der Hologramme. Bei zu viel Licht sind diese deutlich schwerer zu erkennen, vergleichbar mit der Nutzung eines Laptops im Freien, bei dem der Bildschirminhalt wegen zu heller Umgebung ebenfalls schlecht erkennbar ist.

Für die Entwicklung selbst ist es ebenfalls hinderlich das es, zum Zeitpunkt der Erstellung dieser Arbeit, wenig Dokumentationen und Hilfestellungen, zum Beispiel in Entwicklerforen, gibt. Dies ist jedoch üblich für Technologien, die noch nicht lange auf dem Markt sind. Zu guter Letzt ist noch der finanzielle Faktor zu berücksichtigen. Durch die Aktualität der Technologie ist diese auch sehr teuer. Der Preis für eine AR-Brille kann dabei den Betrag von 5.000 Euro überschreiten [14].

Eignung von Virtual Reality für 3D-Annotation

Die Virtual Reality Technologie bietet den großen Vorteil, beliebig große Räume virtuell begehbar zu machen, ohne sich in der realen Welt selbst bewegen zu müssen. Dies ermöglicht die Bewegung durch Punktewolken, die im dreidimensionalen Raum dargestellt sind, an jedem üblichen Arbeitsplatz eines Büros. Des weiteren wirkt die Steuerung der VR-Brillen mittels den zugehörigen Controllern zum Teil sehr ausgereift. Die Positions- und Bewegungserkennung des Benutzers wird als „*äußerst präzise*“ beschrieben [15]. Dies ist sehr wichtig um die Handhabung der Anwendung für den späteren Endnutzer so angenehm wie möglich zu machen.

Für die Entwicklung von VR-Applikationen selbst ist zu sagen, dass es mittlerweile sehr umfassende Dokumentationen, Anleitungen und Beiträge zu den jeweiligen Plattformen und Software Development Kits (SDKs) gibt. Dies hängt damit zusammen, dass sich die Virtual Reality Technologie schon etwas länger auf dem Markt befindet und die Zahl der Entwickler für VR-Applikationen stetig steigt. Dies ist vermutlich nicht zuletzt der Tatsache zu verdanken, dass die Preise für VR-Brillen gesunken sind. Der Preis einer Oculus Rift beträgt zum Zeitpunkt dieser Arbeit 449 Euro. Sowohl die zahlreichen Informationen für Entwickler als auch der niedrige Preis der Hardware sprechen, neben den zuvor genannten Aspekten, für die Wahl von Virtual Reality als Plattform für die Entwicklung von 3D-Labeling.

Was mit der VR-Technologie nicht ohne Weiteres funktioniert, ist die physische Begehung einer Punktewolke. Zwar bietet die Technik die Möglichkeit durch zusätzliche Sensoren die Bewegungen des Benutzers in die virtuelle Welt zu übersetzen, jedoch funktioniert dies nur auf begrenztem Raum und ist unmöglich an einem üblichen Büro-Arbeitsplatz durchführbar. Der reale begehbarer Raum muss nämlich durch diese Sensoren abgesteckt werden und ist somit durch deren Reichweite und Genauigkeit begrenzt. Ein weiter Negativaspekt ist, dass viele Benutzer von VR-Brillen über Übelkeit klagen. In manchen Tests sind es mehr als 50 Prozent aller Teilnehmer, vor allem wenn es sich um Frauen handelt [16].

Entscheidung

Für die Wahl des passenden Mediums, auf dem das dreidimensionale Labeling entwickelt werden soll, sind mehrere Faktoren entscheidend. Diese Faktoren ergeben sich danach, wie sehr sich etwas für die genannte Anwendung eignet. Zunächst ist es wichtig eine Punktewolke deutlich und nach Möglichkeit im richtigen Maßstab darzustellen. Hier kommt die Augmented Reality-Technologie sicherlich an ihre Grenzen. Grund dafür ist die ungenaue Interaktion mit der Umgebung und die lichtempfindlichen Darstellung digitaler Inhalte. In der virtuellen Realität dagegen kann eine solche Wolke problemlos bei jeder Bedingung angezeigt werden, da die reale Umgebung nicht berücksichtigt wird. Sogar der Maßstab

kann gut eingehalten werden, da dem virtuellen Raum keine Grenzen gesetzt sind, sodass beliebig große Umfeldmodelle an jedem Arbeitsplatz dargestellt werden können.

Auch die Steuerung spricht für den Einsatz von Virtual Reality. Zwar ist das Konzept der Greif-Steuerung der Meta 2 ein guter Ansatz, da diese aber noch nicht zuverlässig und genau funktioniert, ist sie für die Selektion vieler Objekte nicht gut geeignet. Dagegen bietet die Steuerung durch Controller bei VR-Brillen deutlich mehr Möglichkeiten eine gute Bedienung für die gewünschte Anwendung zu entwickeln. Die Tasten eines solchen Controllers können vom Entwickler nämlich nach Wunsch programmiert werden wogegen die Gesten für AR-Brillen nicht verändert bzw. nicht neu erstellt werden sollten [17]. Zudem gibt es, durch die deutlich frühere Markteinführung, viel mehr Dokumentationen und Hilfestellungen im Bereich Virtual Reality, was die Entwicklung deutlich erleichtert. Aufgrund der eben genannten Vorteile von VR gegenüber AR wird für die Entwicklung, im Rahmen dieser Arbeit, **Virtual Reality** verwendet.

3.1.2. Warum die Oculus Rift verwendet wurde

Die, in dieser Arbeit, entwickelte Methode zur 3D-Datennotation ist vorwiegend für einen normalen Arbeitsplatz im Büro konzipiert, also einen Platz an dem man einen Desktop-Computer zur Verfügung hat. Deswegen macht es keinen Sinn VR-Brillen in Betracht zu ziehen die von einem mobilen Gerät betrieben werden (*Samsung Gear VR*), da diese deutlich weniger Leistung und eine schlechtere Bedienung haben als Geräte für einen Desktop-PC. Auch die Nutzung der *Playstation VR* ist nicht sinnvoll, da zum Betrieb eine Playstation 4 nötig ist, welche in den wenigsten Unternehmen vorhanden sein dürfte.

Für das Klassifizieren im dreidimensionalen Raum ist eine gute Darstellung und Bedienung notwendig. Bei VR-Brillen, die einen Desktop-PC als Recheneinheit haben wird beides geboten. Durch die hohe Leistung eines performanten Computers kann der dreidimensionale Raum hochauflösend dargestellt werden. Auch die genaue, für Spiele konzipierte, Steuerung über Controller kommt einem Labeling-Tool zugute. Zum Zeitpunkt der Erstellung dieser Arbeit, gibt es zwei relevante VR-Brillen aus dieser Sparte, die *Oculus Rift* und die *HTC Vive*.

Oculus Rift

Mit 470 Gramm ist die Oculus Rift (vgl. Abbildung 3.2b) deutlich leichter als andere Modelle, wie beispielsweise die HTC Vive mit 600 Gramm oder die Playstation VR mit 610 Gramm [18]. Dies ist für längeres Arbeiten mit ihr sehr vorteilhaft, da schwere Brillen oft schon nach kurzer Zeit störend beim Tragen sind. Des weiteren lässt sich die Brille problemlos an jedem handelsüblichen Büro-Arbeitsplatz verwenden. „*Oculus empfiehlt, die*



(a) HTC Vive AR-Brille mit Controllern

(b) Oculus Rift VR-Brille mit Controllern

Abbildung 3.2.: Die zwei potentiellen VR-Brillen

*Kameras im Abstand von zwei Metern nebeneinander zu platzieren – bei einem klassischen PC-Arbeitsplatz also links und rechts neben dem Monitor. Das funktioniert in der Praxis gut, man kann sich mit diesem Aufbau ungefähr jeweils einen Schritt in alle Richtungen bewegen“ [19]. Auch in Sachen Steuerung ist die Oculus Rift anderen Modellen überlegen. Durch die sogenannten *Oculus Touch* Controller, mit denen man virtuell nicht nur Greifen, sondern auch Daumen und Zeigefinger bewegen kann, gelingen ihr „*deutlich feinere Bewegungen*“ als beispielsweise der HTC Vive [19]. Dies bietet für die Entwicklung einer Lösung zur Markierung vieler Objekte zahlreiche Möglichkeiten.*

Was bei der Oculus Rift nicht so gut funktioniert, wie bei anderen Brillen ist das Roomscaling, also das Erfassen der Bewegungen des Benutzers durch einen größeren Raum. Dieser Aspekt erschwert somit das räumliche Begehen einer Punktwolke. Für dieses Verfahren gibt es bei der Rift zwei Methoden. Einmal mit zwei Kameras (kleine Überwachungsfläche), die sich im Abstand von drei Metern in zwei Raumecken gegenüberstehen und eine mit drei Kameras (größere Überwachungsfläche). Beide Methoden funktionieren nicht einwandfrei und auch nicht so gut wie bei der HTC Vive [19].

HTC Vive

Die HTC Vive (vgl. Abbildung 3.2a) punktet mit dem ausgeklügelten Lighthouse Tracking System, welches in Zusammenarbeit mit VALVE entwickelt wurde [20]. Das System kann die Kopfbewegungen des Benutzers durch mehr als 70 Sensoren auf bis zu ein Zehntel eines Grades messen [21]. Die Sensoren sind dabei sowohl in die Brille integriert als auch extra positioniert. Auch die Controller werden durch dieses System erfasst, was zu einer genauen Übertragung der physischen Handbewegungen in die virtuelle Welt führt. Dies kommt der Auswahl zahlreicher Objekte, wie es bei der Klassifizierung von Punkten notwendig ist, sehr entgegen. Auch eine gute Darstellung ist aufgrund zweier Displays gewährleistet, die jeweils eine Auflösung von 1080 auf 1200 Pixel haben. Bei diesen Displays kommt es weder zu Verzerrungen des Bildes („*lens distortion*“ [21]) noch zu Pixelfehlern. Ebenfalls nützlich, vor allem am Arbeitsplatz, ist die Bluetooth-Funktion der Vive. Damit ist es möglich sein

Smartphone mit der Brille zu verbinden und somit eingehende Emails oder Ähnliches zu lesen und zu beantworten, ohne die Brille absetzen zu müssen.

Der große Unterschied zwischen Vive und Rift liegt bei der Bedienung durch die unterschiedlichen Controller. Hier hat die HTC-Brille das Nachsehen, weil sie durch Anatomie der Controller weniger Möglichkeiten zur Interaktion in der virtuellen Umgebung bietet. Die sogenannten *Wands* bieten nämlich keinerlei Funktionen die es ermöglichen eine Hand nachzuahmen, um beispielsweise einzelne Finger zu bewegen und somit virtuelle Objekte greifen zu können. Im Hinblick auf eine durchdachte Steuerung für die Auswahl und Markierung vieler Punkte in einer Punktewolke gibt es somit weniger Möglichkeiten diese zu entwickeln.

Entscheidung

Bleibt noch die Frage zu klären welche Brille verwendet werden soll. Diese Frage ist schwerer zu beantworten als die vorherige. Die zwei potentiellen VR-Brillen, also Oculus Rift und HTC Vive, sind beide geeignete Kandidaten. Die Vive punktet durch gutes Roomscaling und gutes Erfassen der Brille und der Controller, wogegen die Oculus Rift guten Tragekomfort und gut durchdachte Controller-Bedienung bietet.

Wichtig für die Entwicklung einer 3D-Datennotation ist aber nicht, ob man sich physisch, also mittels Roomscaling durch die Punktewolke bewegen kann, sondern wie man die klassifizierbaren Punkte markieren bzw. auswählen kann. Da hierfür die Oculus Touch Controller mehr Möglichkeiten bieten wird für die Entwicklung eines VR-Labeling-Tools die **Oculus Rift** verwendet.

3.2. Komponente 2: Der Rechner

Für die Berechnungen, die notwendig sind um Inhalte in der Oculus Rift anzuzeigen, ist nicht die Brille selbst verantwortlich. Diese Aufgabe übernimmt ein separater Computer. Die Brille fungiert dementsprechend nur als Anzeigegerät. Auch die Sensoren sind an den PC angeschlossen um das Erfassen der Brille und der Controller zu gewährleisten. Vor allem die Grafikberechnungen sind sehr aufwendig sind, da die Darstellung von VR-Inhalten zwei mal berechnet werden (für jedes Auge extra) und nicht nur einmal, wie es an einem Bildschirm der Fall ist. Deshalb kann nicht jeder handelsübliche PC oder Laptop dazu verwendet werden eine solche VR-Brille zu betreiben. Was ein passender Computer für die Nutzung von Virtual Reality bieten muss wird im Folgenden erläutert.

Wichtig sind zunächst die Basiskomponenten wie ein schneller Prozessor und eine leistungsstarke Grafikkarte. Aber auch ein schneller Arbeitsspeicher mit ausreichender Kapazität ist wichtig damit notwendige Daten so schnell wie möglich zur Verfügung stehen.

Für die Wahl, welche Komponenten eingesetzt werden sollen, kann sich an den empfohlenen Spezifikationen des Herstellers orientiert werden (Oculus Rift Recommended System Specs [22]). Diese sind jedoch vorwiegend für den Endnutzer gedacht, bei dem es wichtig ist ein einzelnes Programm für die Brille auszuführen. Für die Entwicklung ist dagegen mehr Leistung notwendig. Soll zum Beispiel eine erstellte Applikation getestet werden, muss nicht nur die Applikation selber ausgeführt werden, sondern auch die Entwicklungsumgebung und zusätzliche Fehlerbehebungs- und Analyse-Anwendungen. An dieser Stelle darf also nicht gespart werden.

Sind die Basiskomponenten ausgewählt, sollte noch geprüft werden ob noch weiteres Zubehör wichtig ist um diese zu betreiben. Da gerade Prozessoren mit hoher Taktrate viel Wärme produzieren, muss stets für ausreichende Kühlung gesorgt werden. Bei der Rechner-Konfiguration für diese Masterarbeit wurden beispielsweise, neben einem hochwertigen Prozessorkühler, noch zwei weitere Gehäuselüfter verbaut, um die Abwärme aus dem inneren des Computers heraus zu leiten.

Zu guter Letzt ist es wichtig ein passendes Mainboard auszuwählen, auf dem alle Komponenten verbaut werden. Für die Oculus Rift werden vier USB-Anschlüsse empfohlen, welche das Mainboard haben muss [22]. Ebenfalls muss das Board genug Platz bieten um die gewählten Komponenten zusammen anbringen zu können. Leistungsstarke Grafikkarten und Prozessorlüfter können zum Teil sehr groß ausfallen. Auch für weitere Aufrüstungen sollte das Board genügend Platz bieten, falls es, zu einem späteren Zeitpunkt, notwendig wäre zusätzliche Komponenten (zweite Grafikkarte) hinzuzufügen. Unter Berücksichtigung der eben genannten Aspekte wurde folgende Rechner-Konfiguration für die Entwicklung verwendet:

Bezeichnung der Komponente	Verwendete Komponente
Prozessor:	Intel® Core™ i7-7700K
Grafikkarte:	Gainward GeForce GTX 1070 Phoenix
Mainboard:	ASUS PRIME Z270-A
Arbeitsspeicher(RAM):	HyperX DIMM 16 GB DDR4-2400 Kit
Festplatte:	Samsung 850 Pro 2,5" 512 GB
Netzteil:	Cooler Master G550M 550W
Prozessorlüfter:	Noctua NH-D9L
Gehäuselüfter:	2x Coolink SWiF2-1200 120x120x25
Gehäuse:	Cooler Master N300

Tabelle 3.1.: Rechner-Konfiguration für die Entwicklungen dieser Arbeit

3.3. Komponente 3: Die Entwicklungsplattform

Für die Entwicklung von Virtual Reality Anwendungen werden in der Regel sogenannte Spiel-Engines (*Game Engines*) verwendet. Dies sind komplexe Mehrzweckwerkzeuge für die Erstellung von Multimedia-Inhalten und Videospielen. Eine Spiel-Engine bietet Funktionen, welche die wichtigsten Bereiche der Spielentwicklung abdecken. Eine davon ist die Bildsynthese aus Rohdaten, welche meist geometrische Beschreibungen im zwei- oder dreidimensionalen Raum sind. Die Bildsynthese ist in der Fachsprache besser bekannt als *Rendering*. Der Prozess des Renderings berechnet für jedes angezeigte Bild die Objekte, welche vom Benutzer aus sichtbar sind, deren Oberflächen anhand von ihren Materialeigenschaften und die Lichtverhältnisse des Bildes, die sich unter anderem durch die indirekte Beleuchtung zwischen Körpern äußert. Weitere Funktionen, welche diese Plattformen zu Verfügung stellen, sind beispielsweise Physikalisch Berechnungen (z.B. Schwerkraft), Animationen von Objekten und künstliche Intelligenzfunktionen.

Für die Entwicklung von VR-Inhalten ändert sich an den Funktionen, welche die Engine bieten muss nicht all zu viel, da die meisten Anwendungen schon dreidimensional sind. Die meisten Elemente der Bildsynthese müssen allerdings, wie schon im vorherigen Abschnitt 3.2 erwähnt, zweimal berechnet werden, nämlich für jedes Auge einzeln. Im Folgenden soll erläutert werden warum für diese Arbeit die Unity Engine als Entwicklungsplattform ausgewählt wurde. Anschließend wird diese kurz vorgestellt.

Vergleich zwischen Unreal und Unity

Für die Erstellung von Virtual Reality Software gibt es zwei relevante Entwicklungsumgebungen, die gut für das Entwickeln von Virtual Reality geeignet sind. Der Grund dafür ist, dass sie schon seit längerem die Erstellung von VR-Anwendungen unterstützen und dies auf umfassende Weise. Die zwei Entwicklungsumgebungen sind *Unreal Engine 4* und *Unity Engine*. Beide sind für nicht kommerzielle Zwecke zunächst kostenlos benutzbar. Beide Plattformen bieten alles nötige um VR-Applikationen zu entwickeln, weshalb die Wahl zwischen ihnen eher subjektiv, aufgrund der eigenen Bedürfnisse und Erfahrungen, zu fällen ist.

In [23] wurden diese beiden Entwicklungsplattformen gegenübergestellt. Im Allgemeinen kann man sagen, dass die Unreal Engine mehr Möglichkeiten für professionelle Spieleentwicklung bietet, da viel Wert auf visuelle Qualität gelegt wird. Beispielsweise wird durch integrierte Post Processing-Techniken und gute Shader das erzeugte Bild besser dargestellt als in Unity. Sie bietet neben dem Scripting auch die Gelegenheit, durch C++ Programmierung Module der Engine zu verändern bzw. neue hinzuzufügen. Das gestalten von User Interfaces ist mit dem *UMG UI Designer* ebenfalls gut gelöst .

Die Unity Engine bietet vor allem Anfängern einen leichten Einstieg in die Welt der Grafik- und Spieleprogrammierung. Im Editor können ganz einfach sogenannte *Game Objects* erstellt werden, denen dann Funktionalitäten und Eigenschaften angehängt werden können. Dies geschieht meist in Form von Skripten, welche oft schon vorgefertigt zur Verfügung stehen. Wenn dies nicht der Fall ist können diese angepasst oder neu erstellt werden. Die große Community der Unity Engine ist hierbei eine große Hilfe.

Entscheidung

Ein Vorteil den die Unity Engine, im Falle dieser Arbeit, bietet ist Nutzung von C# als Skript-Sprache. Da C.LABEL (vgl. 1.3.2) ebenfalls in C# programmiert wurde können Komponenten, wie etwa das Einlesen von Daten, leicht in die VR-Anwendung übernommen werden. Des Weiteren habe ich selbst, im Rahmen eines Universitätsprojekts, mit dieser Plattform gearbeitet. Ein längeres einarbeiten in die Entwicklungsumgebung wäre demnach nicht nötig. Die Vorteile der Unreal Engine sind, wie schon erwähnt, die visuelle Präsentation und die dabei gebotene Performance. Diese sind vor allem für das Erstellen aufwändiger Landschaften und deren Inhalten von Vorteil.

Für die Entwicklung der angestrebten Anwendung spielen jedoch Dinge wie Objekt- und Landschaftsdarstellung keine große Rolle. Viel wichtiger ist sowohl die Möglichkeit der Kompatibilität zur C.LABEL-Anwendung, als auch die vorhandene Entwicklungserfahrung mit einer Plattform. Gerade letzteres ist auf Grund des begrenzten Zeitrahmens einer Masterarbeit von Vorteil, da ohne große Einarbeitung mehr Zeit zum Entwickeln von Funktionalität und Optimierungen bleibt. Aus diesen Gründen wird im weiteren Verlauf dieser Arbeit die **Unity Engine** als Entwicklungsumgebung verwendet.

Entwickeln mit Unity

Die Unity Engine ist eine der am häufigst benutzten Spiel-Engines auf dem Markt. Gerade im Bereich der mobilen Geräte wurden 34 Prozent der 1000 beliebtesten Spiele mit Unity erstellt [24]. Sie ist ebenfalls bei Neueinsteigern und Hobbyprogrammierern beliebt, da sie über umfangreiche Lernmaterialien und eine große Entwicklergemeinschaft verfügt. Die Unity Engine unterstützt sowohl alle gängigen großen Plattformen wie PC, Playstation und XBOX, als auch die kleinen mobilen Geräte wie Android, iOS, Nintendo 3DS und PS Vita. Darüber hinaus kann man mit ihr auch Anwendungen für, in der Applikationsentwicklung, weniger gebräuchliche Plattformen wie Samsung SMART TV, oder tvOS erstellen. Wichtig für diese Arbeit ist die umfangreiche Unterstützung von Virtual Reality Brillen, speziell der Oculus Rift. Um die Vorgänge von C.LABEL-VR, die in Kapitel 4 beschrieben werden, besser zu verstehen, werden im Folgenden die Grundlagen für die Entwicklung mit Unity etwas näher erläutert.

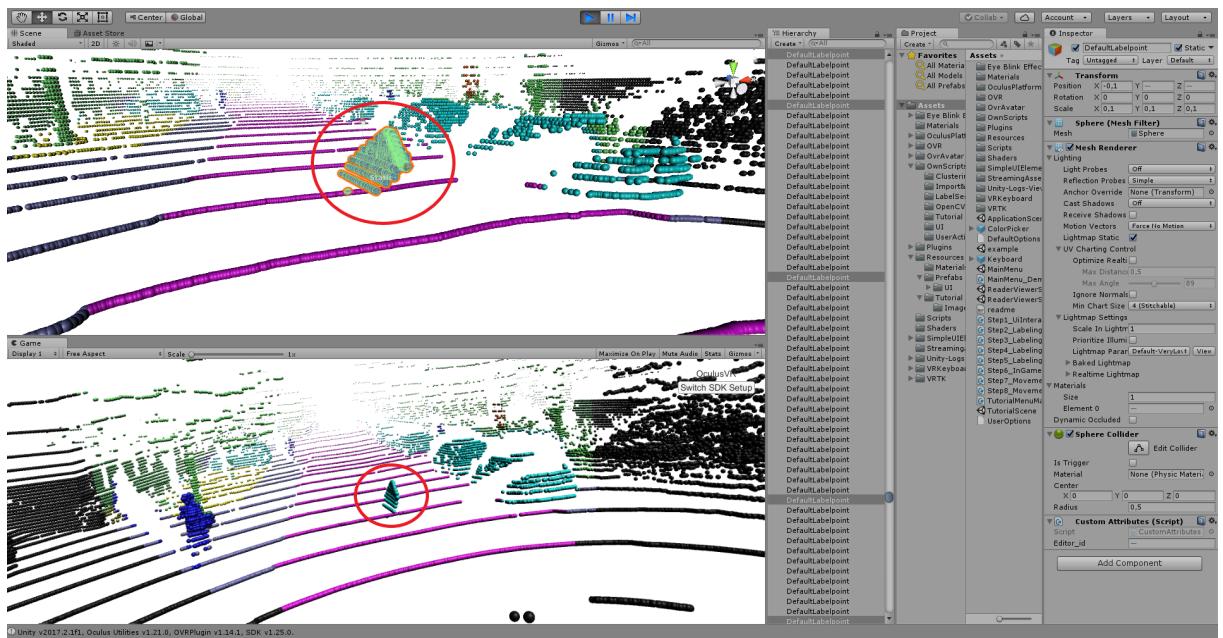


Abbildung 3.3.: Unity Editor, der C.LABEL-VR im Play Mode ausführt.

Bei der Entwicklung mit Unity werden Anwendungen mittels Szenen konzipiert. Eine Anwendung besteht aus mindestens einer Szene, in der Regel sind es aber mehrere. Jede Szene kann dabei eigene Objekte und Logiken beinhalten und es kann je nach Programmierung des Entwicklers zwischen den Szenen gewechselt werden. Ein klassischer Fall dafür wäre eine Applikation mit einer Hauptmenü- und einer Anwendungsszene. Als erstes wird die Hauptmenü-Szene gestartet und von dieser aus dann, beispielsweise durch die Betätigung eines Knopfes, die Anwendungsszene. In größeren Spielen ist es üblich jeden abgetrennten Bereich im Spiel mit einer extra Szene zu erstellen. Dies dient der Übersicht und der Wartbarkeit, die in Szenen mit zu vielen Objekten nicht mehr gegeben wäre.

Die Konzeption der einzelnen Szenen geschieht, wie in der Regel in bei jeder Spiel-Engine, über einen Editor (siehe Abbildung 3.3). Die Konzeption beinhaltet beispielsweise die Erstellung von Landschaften und deren Objekten, Anordnung der Objekte und das Testen der Szenen. Für die Erstellung von komplexen 3D-Objekten werden externe Programme wie *Blender* verwendet, deren Erzeugnisse dann in den Unity Editor importiert werden können [25]. Die Erstellung von simplen primitiven Objekten wie Kugeln oder Würfeln ist aber auch direkt im Unity Editor möglich.

Mit der Anordnung der Objekte ist zum einen die Platzierung der Objekte innerhalb der Szene gemeint, welche beim Beginn der Szene die Startposition der Objekte ist. Diese kann aber während des Ablaufs der Szene durch Programmierung geändert werden. Zum anderen ist die Anordnung der Objekte innerhalb von Hierarchien gemeint. Dabei werden *Eltern-Kind*-Beziehungen zwischen den Objekten festgelegt, was zur Folge hat, dass Änderungen bei den Eltern sich auch auf die Kinder auswirken. Erstellt man zum Beispiel ein leeres Objekt „Mensch“ und ordnet diesem Kopf, Arme, Beine und Oberkörper zu, so kann man

das Objekt „Mensch“ bewegen und alle Körperteile bewegen sich mit.

Das Testen einer oder mehrerer Szenen ist mit dem *Unity Play Mode* möglich. Startet man diesen Modus wird die aktuelle Szene, die im Editor geöffnet ist, gestartet ohne die gesamte Applikation vorher kompilieren zu müssen. Das spart dem Entwickler Zeit da zum einen die Kompilierzeit wegfällt und zum anderen die Applikation nicht immer von der Startszene aus gestartet werden muss. Der Editor in der Abbildung 3.3 befindet sich im Play Modus. In diesem Editor sind zwei große Fenster abgebildet, in denen die laufende Szene zu sehen sind. Das obere Fenster zeigt die Editor-Sicht. Innerhalb dieser können Objekte während des Ablaufs der Szene ausgewählt und geändert werden. Im roten Kreis sind markierte Objekte zu sehen, die nach links verschoben wurden. Im unteren Fenster ist die Benutzer-Sicht zu sehen, welche die Szene aus der Sicht des Benutzers zeigt. Dadurch kann überprüft werden wie sich die getätigten Änderungen aus der Editor-Sicht auf die Tatsächliche Sicht des Benutzers auswirken. Auch hier sieht man im roten Kreis die Objekte die im Editor verschoben wurden.

Die erwähnten Objekte, welche sich in den einzelnen Szenen befinden, werden in Unity *Gameobjects* genannt. Ein Gameobject ist die übergeordnete Klasse für alle Arten von Objekten in Unity, beispielsweise 2D- und 3D-Objekte, aber auch Elemente von Benutzeroberflächen. Jedes Objekt in Unity ist also ein Gameobject und es ist möglich diesen Objekten bestimmte Komponenten zuzuordnen. Mit diesen Komponenten lassen sich Eigenschaften und Verhaltensweisen der Objekte festlegen. Jedes Gameobject in Unity hat standardmäßig eine, nicht entfernbare, Komponente die sich *Transform* nennt. Diese legt die Position, Orientierung und Größe des Objekts in der Szene fest. Sichtbare Objekte müssen beispielsweise eine *Mesh Renderer*-Komponente haben, welche die Materialeigenschaften des Objektes festlegt und somit der Unity Engine signalisiert, dass das Objekt bei der Bildsynthese berücksichtigt werden muss. Es kann aber auch unsichtbare Objekte geben, die lediglich eine Logik-Komponente haben, welche beispielsweise Aktionen anderer Objekte überwacht.

Die Komponenten legen die Eigenschaften und Verhaltensweisen von Objekten mit Hilfe eines Skriptes fest, welches mit einer bestimmten Programmiersprache erstellt wird. In Unity stehen die Sprachen C# und JavaScript zur Verfügung. Das Implementieren dieser Skripte wird in der Fachsprache als *Scripting* bezeichnet. Dadurch kann der Programmierer, wie schon erwähnt, die Eigenschaften und das Verhalten von Objekten festlegen und beeinflussen. Zur Veranschaulichung dieses Vorgangs dient die Auflistung 3.1. Darin ist ein minimalistischer Grundaufbau eines Skriptes zu sehen, dass ein Objekt unter einer Bedingung größer werden lässt.

Auflistung 3.1: Grundaufbei eines Gameobjects in Unity

```
using UnityEngine;

public class TestObject : MonoBehaviour
{
    public static TestObject Instance { get; set; }
    public bool EnableGrowing { get; set; }

    private void Awake()
    {
        Instance = this;
    }
    private void Start()
    {
        Instance.EnableGrowing = true;
    }
    private void Update()
    {
        if(Instance.EnableGrowing)
            Instance.transform.localScale += new Vector3(1, 1, 1);
    }
}
```

Ein Skript besteht in der Regel aus einer Klasse, wie beispielsweise *TestObject* in Auflistung 3.1, die von der Unity-Basisklasse *MonoBehaviour* erbt. Diese Vererbung führt dazu, dass der Klasse *TestObject* diverse Funktionen zugeordnet werden, die von der Unity Engine zu gewissen Zeitpunkten aufgerufen werden. Dadurch kann das Verhalten des Objekt gezielt gesteuert werden. Diese Funktionen werden jedoch nur aufgerufen, wenn sie innerhalb der erbenden Klasse auch implementiert werden. All diese Funktionen sind in der Dokumentation der Unity Engine aufgelistet und beschrieben [26]. Die Klasse *TestObject* implementiert drei von diesen Funktionen, nämlich *Awake()*, *Start()* und *Update()*.

Awake wird als erstes aufgerufen, sobald alle Objekte in Unity geladen sind und kann somit zur Initialisierung von Referenzen verwendet werden bevor die Szene startet. Im Beispiel wird der statischen Instanz von *TestObject* die Referenz des geladenen Skripts *TestObject* zugewiesen. Dieses Vorgehen nennt sich *Singleton-Pattern* und ist eine übliche Methode zum verwalten von Referenzen in Unity [27]. Dadurch ist es für andere Skripte möglich, ohne Instanziierung der Klasse, auf diese zuzugreifen.

Die Methode *Start* wird vor dem ersten Aufruf von *Update* ausgeführt, also kurz vor dem Start der Szene. Hier können somit noch Werte zugewiesen werden. *EnableGrowing* wird im Skript *TestObject* auf den boolschen Wert „wahr“ gesetzt. Dieser kann dann von anderen Skripten durch das eben beschriebene Singleton-Pattern geändert werden

(*TestObject.Instance.EnableGrowing = false*).

Die Update-Methode ist eine Art Endlosschleife die von der Unity Engine immer wieder aufgerufen wird solange das Skript in einem aktiven Zustand ist. Sie repräsentiert die, in der Fachsprache oft genannte, *Game-Loop*. Alle Vorgänge die während der Laufzeit der Applikation ablaufen sollen müssen also in dieser Funktion implementiert werden. Im Beispiel wird die Größe des Objektes erhöht, solange die Variable EnableGrowing den boolschen Wert „wahr“ hat. Es gibt darüber hinaus noch einige weitere dieser Methoden innerhalb der MonoBehaviour-Klasse, beispielsweise die Methode *OnEnable()*, die aufgerufen wird wenn das Skript vom inaktiven in den aktiven Zustand wechselt. Mit Hilfe all dieser MonoBehavior-Funktionen lassen sich also Objekte und deren Verhalten manipulieren. Auf diese Weise werden gewünschte Vorgänge in Unity implementiert.

In diesem Abschnitt wurden nun grundlegende Vorgänge beschrieben, die wichtig sind um Anwendungen mit Unity zu erstellen. Darunter waren zum Beispiel das Konzipieren von Szenen im Editor und das Scripting, also das Programmieren von Logik für die Anwendung. Damit soll ein Einblick in den Prozess der Entwicklung von C.LABEL-VR gegeben werden. Die Beschreibung der Funktionen von C.LABEL-VR selbst und deren Konzepte sind im folgenden Kapitel 4 beschrieben.

4

C.LABEL-VR

Dieses Kapitel beschreibt die Virtual Reality Applikation *C.LABEL-VR* und ist somit der Hauptteil dieser Arbeit. Das Ziel dieser Applikation ist es die Annotierung von Punktwolken, wie sie beispielsweise in C.LABEL möglich ist, innerhalb einer dreidimensionalen Umgebung zu realisieren. Dazu müssen zunächst Datenformate eingelesen werden, die Informationen über Punktwolken enthalten. Anschließend müssen aus diesen Daten Punktwolken erzeugt werden. Der Vorgang des Imports wird in Kapitel 4.1 näher beschrieben, die Erzeugung der Punktwolken in 4.2. Um sich in der virtuellen Umgebung durch diese Wolken bewegen zu können, wurden diverse Möglichkeiten zur Navigation entwickelt. Auf die Funktion dieser Möglichkeiten und deren Auswirkungen auf das Befinden des Menschen (*Virtual Motion Sickness*) wird in Abschnitt 4.3 eingegangen.

Anschließend geht es um die Annotierung der generierten Punktwolken. Annotierung bedeutet in diesem Kontext, dass die einzelnen Punkte der Wolken mit bestimmten Klassifikationen versehen werden, welche der Art des Objektes entsprechen, dem die Punkte zugehören. Ist der Punkt beispielsweise Bestandteil eines Autos, so wird er mit der Klasse *Auto* versehen. Um diese Aufgabe erfüllen zu können wurden mehrere Arten der Annotierung entwickelt. Welche dies sind und wie sie realisiert wurden, wird in Abschnitt 4.4 gezeigt. Das letzte Kapitel 4.5 beschäftigt sich mit der Benutzerschnittstelle der Applikation. Dieses wird auch *User Interface* oder kurz *UI* genannt. Dabei wird hauptsächlich auf die Funktionen des Menüs eingegangen welches man während des Labelns aufrufen kann. Ebenfalls wird erläutert welche Herausforderungen es bei der Erstellung von UIs in der virtuellen Realität gibt und wie diese bewältigt wurden. Zur Veranschaulichung aller Funktionen von C.LABEL-VR ist in der Abbildung 4.1 der grobe Workflow in Form eines Diagramms abgebildet.

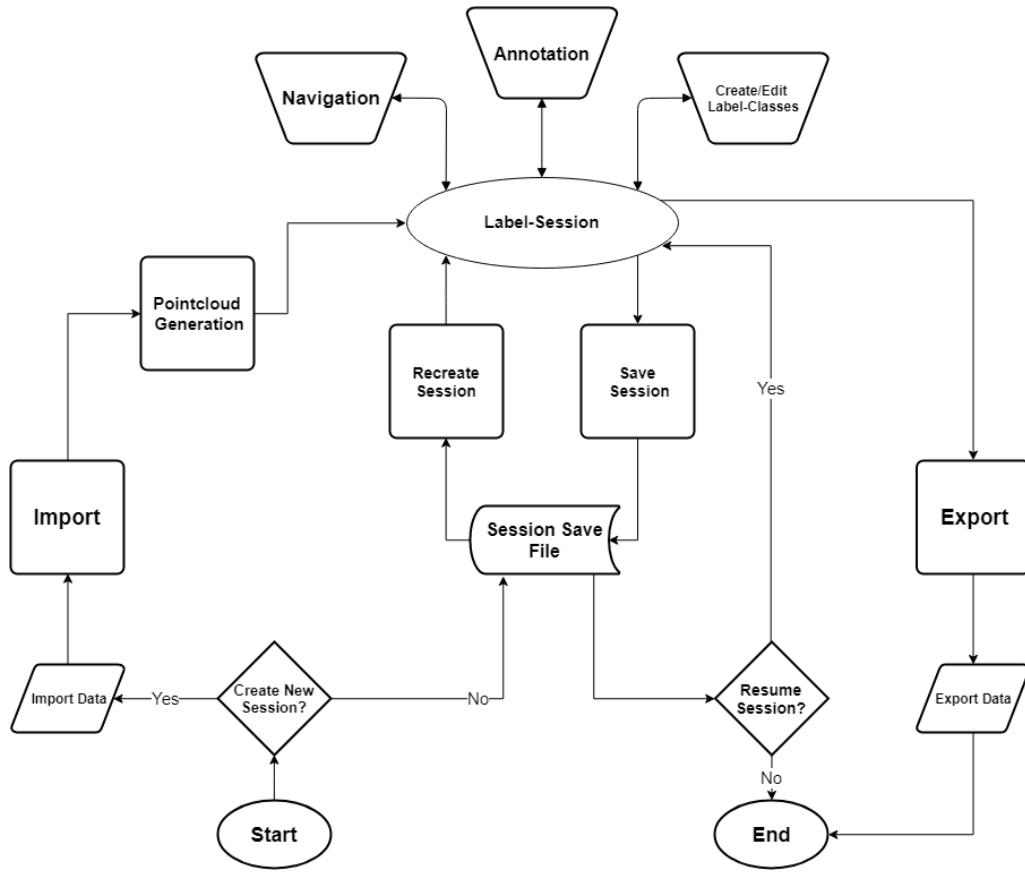


Abbildung 4.1.: in C.LABEL-VR vom Import bis zum Export

4.1. Import und Export von Daten

Umfeldmodelle, welche in Form einer Punktwolke vorliegen, werden meist von einem Radar-, oder, wie in Kapitel ?? gezeigt, von einem LiDAR-Sensor aufgenommen. Die Informationen über diese Modelle, zum Beispiel die Positionen der einzelnen Punkte, werden in einem passenden Datenformat abgespeichert. Um welches Format es sich dabei handelt kann sehr unterschiedlich sein. In der Regel hängt es davon ab zu welchem Zweck die aufgenommenen Daten gebraucht werden und wie umfangreich diese Daten sind. Am häufigsten werden sie dafür verwendet um Lernalgorithmen und neuronale Netze (siehe Abschnitt 2.1.1) zu trainieren. Automobilhersteller haben für solche Vorhaben natürlich unterschiedliche Konzepte und somit können auch die Formate der Daten unterschiedlich sein. Um die Applikation C.LABEL-VR für mehrere dieser Hersteller benutzbar zu machen, muss das Programm also mit verschiedenen Datenformaten umgehen können.

Innerhalb dieser Arbeit wurden mit zwei verschiedenen Arten von Daten gearbeitet. Das erste Format ist das PCD-Format (*Point Cloud Data*). Die Bibliothek PCL(*Point Cloud Library*) benutzt das Datenformat **PCD** um Informationen über Punktwolken zu verwalten. PCL ist eine sehr mächtige C++-Bibliothek, die in der Regel von jedem benutzt wird, der performante Algorithmen für Punktwolken entwickeln möchte. Auch C.LABEL generiert

aus eingelesenen Daten PCD-Files um Funktionen der Point Cloud Library für die Wolken zu benutzen.

Das zweite Format ist das, von der HDF-Group bereitgestellte, HDF5-Format (*Heterogeneous Data Format*). Es ist ausgelegt zum schnellen Erstellen und Auslesen von komplexen und großen Datenstrukturen. Welche Elemente diese Strukturen enthalten können wird später in Kapitel 4.1.2 erklärt. Zur schnellen Bearbeitung großer Datenmengen ist dies ein gängiges Format und wird von den meisten Kunden von CMORE benutzt um Sensordaten zu verwalten. Deshalb wird es auch in C.LABEL-VR verwendet.

4.1.1. Architektur

Wichtig ist es also sowohl den Import, als auch den Export so modular wie möglich zu gestalten, dass C.LABEL-VR ohne große Änderungen an der Hauptapplikation um neue Datenformate erweitert werden kann. Das Prinzip, das dafür entwickelt wurde, ist in Abbildung 4.2 dargestellt und wird im Folgenden näher erläutert. Die angesprochenen Erweiterungen um neue Datenformate werden als *Addons* bezeichnet. Jedoch handelt es sich dabei nicht immer nur um eine Erweiterung eines, von Grund auf, neuen Datenformats. Bei Komplexen Datenformaten wie HDF5 kann der Entwickler die interne Struktur der Datei selbst bestimmen. Jede unterschiedliche Struktur muss auf unterschiedliche Weise eingelesen werden und braucht somit unterschiedliche Addons. Um die Erweiterbarkeit der Applikation sicherzustellen wurde der Import und der Export von der Hauptapplikation abgekapselt.

Dafür wurde eine interne Datenstruktur (kurz IDS) eingeführt, welche alle nötigen Informationen enthält, die für das Labeling notwendig sind. Sie wird im späteren Kapitel 4.2.1 genauer vorgestellt. Der Vorteil dabei ist, dass die Hauptapplikation dadurch stets auf die gleiche Datenstruktur zurückgreifen kann. Sie bleibt somit unabhängig vom importierten Datenformat. Wird ein neues Import-Format eingeführt muss an der Hauptapplikation also nichts verändert werden. In der Abbildung 4.2 erkennt man, wie die Hauptapplikation (*C.LABEL-VR Session*) nur mit internen Daten arbeitet und somit von den Import- und Exportdaten getrennt ist.

Jedes Import-Addon hat also zunächst die Aufgabe, seine jeweiligen Daten einzulesen und alle Informationen, die für die IDS notwendig sind, daraus zu extrahieren. Aus diesen Informationen können anschließend die entsprechenden Punktwolken generiert werden. Die Generierung wird in Kapitel 4.2.2 erklärt. Beim Export ist es notwendig, dass die Struktur der Daten gleich bleibt, die vom Import-Addon eingelesenen wurden. Es sollen lediglich die Labeling-Informationen aus der IDS hinzugefügt werden. Aus der IDS kann das Export-Addon diese Informationen extrahieren und anschließend in die eingelesenen Daten exportieren. Die importierten Daten werden so allerdings überschrieben. Möchte

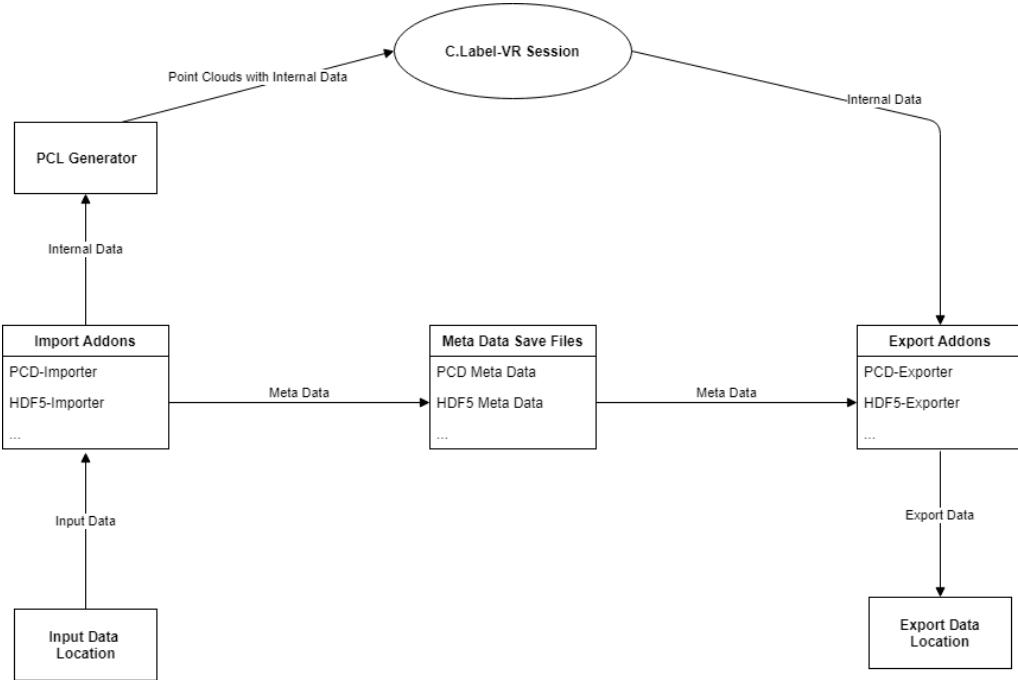


Abbildung 4.2.: Das Import-Export-Prinzip in C.LABEL-VR

man die Daten nicht überschreiben, müssen die Daten neu erstellt und mit den Labeling-Informationen versehen werden.

Die Importierten Daten, beispielsweise HDF5, enthalten allerdings oft deutlich mehr Informationen, als für die interne Struktur notwendig. Diese Zusatzinformationen werden *Metadaten* genannt. Für den Export bedeutet dies, dass man aus der internen Datenstruktur von C.LABEL-VR, die Struktur der anfangs eingelesenen Daten nicht wieder rekonstruieren kann, da die *Metadaten* fehlen würden. Um dieses Problem zu lösen müssen die zusätzlichen Daten separat gespeichert werden. Dazu wurde eine erweiterbare Datenstruktur angelegt. In dieser kann man für jede Import-Struktur eine Metadatenstruktur anlegen. Das entsprechende Export-Addon kann die Metadaten dann verwenden um, zusammen mit den internen Daten, das gewünschte Exportformat zu rekonstruieren.

Zusammenfassend ist also folgendes zu tun, um die Benutzung ein neues Datenformates oder einer neuen Datenstruktur zu gewährleisten: Zunächst muss eine Importfunktion programmiert werden, die alle Informationen aus den gewünschten Daten ausliest. Anschließend muss diese Funktion das genormte, interne Datenformat (Kapitel 4.2.1) erstellen und alle nötigen Informationen aus den eingelesenen Daten darin speichern. Daraus können die entsprechenden Punktwolken dann generiert werden können. Falls es Daten gibt, die über diejenigen in der IDS hinausgehen, muss eine Metadatenstruktur angelegt werden, in welche die zusätzlichen Daten abgelegt werden können. Zuletzt ist eine Exportfunktion zu implementieren, welche die Labeling-Informationen aus der internen Datenstruktur in die eingelesenen Daten exportieren kann. Außerdem muss sie in der Lage sein aus den

internen- und den Metadaten das Ausgangsdatenformat wiederherzustellen, um die Daten als neue Files exportieren zu können. Im folgenden Kapitel soll dieses Prinzip konkret an einem HDF5-Beispiel gezeigt werden.

4.1.2. HDF5-Beispiel

HDF5 Daten

TODO Evtl. den Import und Exportvorgang anhand eines minimalistischen HDF5-Beispiel erklären

4.2. Generierung einer Punktewolke

TODO Kurze einleitung

4.2.1. Interne Datenstruktur

TODO -Wie sieht die interne Datenstruktur aus die Konzipiert wurde um die Importierten Daten zu vereinheitlichen. -Was passiert mit den Daten, die über diese Struktur hinaus gehen (Metadaten)

4.2.2. Generierung

TODO Wie wird aus der internen Struktur eine Punktewolke in Unity ezeugt?

4.2.3. Optimierung

TODO Was ist bei der Erzeugung der Punktewolken zu beachten, damit die vielen Tausend Punkte der Wolke keinen massiven FPS-Einsturz verursachen

4.3. Navigation

TODO kurze Einleitung

4.3.1. VR-Krankheit

TODO kurze Erläuterung von Motion Sickness

4.3.2. Freier Flug

TODO Erklärung der FreeFly-Modus und des Funktionsprinzips dahinter

4.3.3. Teleport

TODO Erklärung des Teleport Modus und des Funktionsprinzips dahinter

4.4. Annotieren der Punktwolke

TODO kurze Einleitung

4.4.1. Einfache Pointer Annotation

TODO Erklärung der Pointer Annotation und des Funktionsprinzips dahinter

4.4.2. Cluster Annotation

In einer Punktwolke mit mehreren tausend Punkten wird es für den Benutzer auf Dauer mühsam jeden Punkt einzeln zu Annotieren. Bei diesem Problem soll die *Cluster Annotation* Abhilfe schaffen. Diese Art der Annotation soll selbstständig diejenigen Punkte in einen Cluster aufnehmen und sie mit der gleichen Klasse labeln, die zum gleichen Objekt gehören. Der Benutzer wählt hierfür, wie im vorherigen Kapitel 4.4.1 einen Punkt mit Hilfe des Pointers aus. Ausgehend davon werden umliegende Punkte untersucht, ob sie zum gleichen Objekt gehören wie der ausgewählte Startpunkt. Ist dies der Fall werden sie in den Cluster aufgenommen. Werden keine Punkte mehr aufgenommen bleibt der Endcluster übrig, welcher dann einem Objekt entspricht, zum Beispiel einem Auto.

Wissenschaftler beschäftigen sich intensiv mit solchen Verfahren zur Objektidentifizierung in dreidimensionalen Punktwolken. Ansätze, wie sie in [4], [28] oder [29] vorgestellt werden sind nur wenige von vielen. Die meisten haben aber ein ähnliches Vorgehen. Zuerst werden all diejenigen Punkte in der Wolke entfernt, die zum Boden gehören. So bleiben anschließend nur noch Punkte übrig, welche zu Objekten wie beispielsweise Autos, Häuser oder Straßenschilder gehören. Die übrigen Punkte können nun mit diversen Algorithmen zu Clustern zusammengefasst werden. Ein Beispiel für das Ergebnis eines solchen Verfahrens bietet die Abbildung 4.3. Im Folgenden wird vorgestellt wie die Cluster-Analyse in C.LABEL-VR funktioniert.

Bodenpunkt-Analyse

Wie schon erwähnt werden bei vielen Verfahren der Cluster-Analyse von Punktwolken zuerst die Bodenpunkte identifiziert. Dieses Verfahren wird auch in dieser Arbeit angewendet. Der Grund dafür ist, dass die Bodenpunkte nicht in das Suchverfahren nach Objektpunkten aufgenommen werden dürfen. Sie würden nicht nur die Suche nach echten Objektpunkten

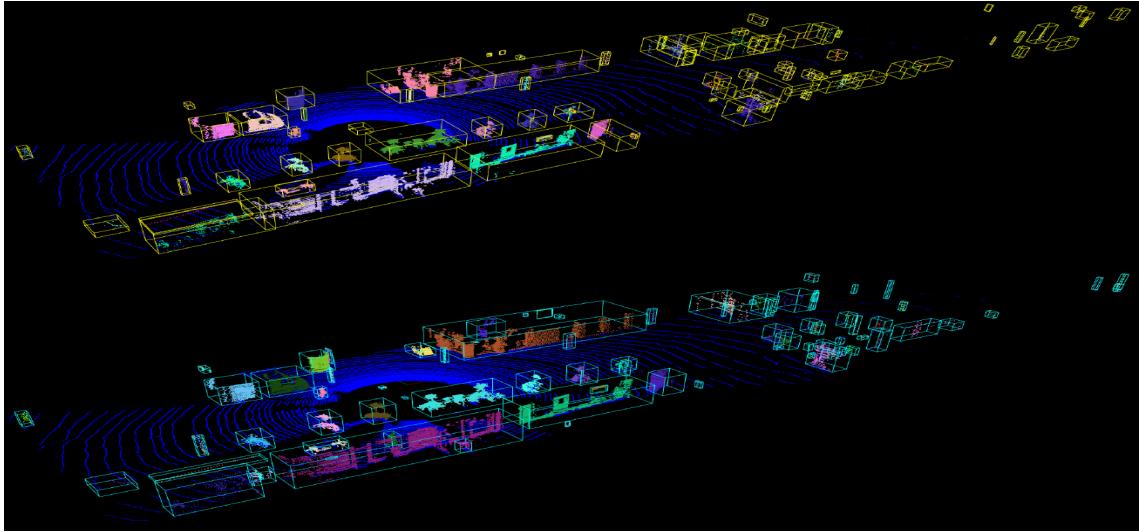


Abbildung 4.3.: Objektidentifizierung einer 3D-Punktwolke aus [4]. Die Bodenpunkte der Wolke sind blau dargestellt und die Objekte sind mit Boxen gekennzeichnet.

erschweren sondern auch das Ergebnis der Objektcluster verfälschen. Als Erklärung soll folgender, simpler Cluster-Algorithmus dienen:

1. Definiere einen Startpunkt p .
2. Ausgehend von p , suche alle Punkte innerhalb von Radius x .
3. Jeder gefundene Punkt ist ein neuer Startpunkt p .
4. Führe 2. aus bis es keine neuen Punkte p mehr gibt.

Alle Objekte die auf dem Boden stehen, also deren Punkte nahe an Bodenpunkten sind, würden mit dem obigen Algorithmus die Bodenpunkte mit in den Cluster aufnehmen. Dies kann sogar dazu führen, dass 2 Objekte, die mittels Bodenpunkten miteinander verbunden sind zu einem Cluster zusammengefügt werden. Darum müssen die Bodenpunkte von den Objektpunkten getrennt werden um eine exakte Objekterkennung zu garantieren.

Die Segmentierung des Bodens wird in C.LABEL-VR nach dem Einlesen des internen Datenformates ausgeführt, also zwischen Import Addons und PCL Generator (siehe Abbildung 4.2). Beim Erstellen der Punktwolken ist also schon bekannt welcher Punkt ein Bodenpunkt ist und welcher nicht. Dazu wird im internen Datenformat das Attribut für die Bodenpunkt kennzeichnung mit dem entsprechenden Wert versehen (vgl. Abschnitt 4.2.1). Wie diese Werte ermittelt werden ist im Folgenden erklärt.

Als Basis für die Bodenpunktanalyse in C.LABEL-VR wurde das *Ground Plane Fitting*-Verfahren (GPF) aus [4] verwendet. Ist im weiteren Verlauf dieser Arbeit die Rede von GPF, ist immer der Ansatz der eben zitierten Arbeit gemeint. Ziel dieses Verfahrens ist es ein mathematisches Modell einer Ebene zu errechnen, welches den Boden repräsentieren soll.

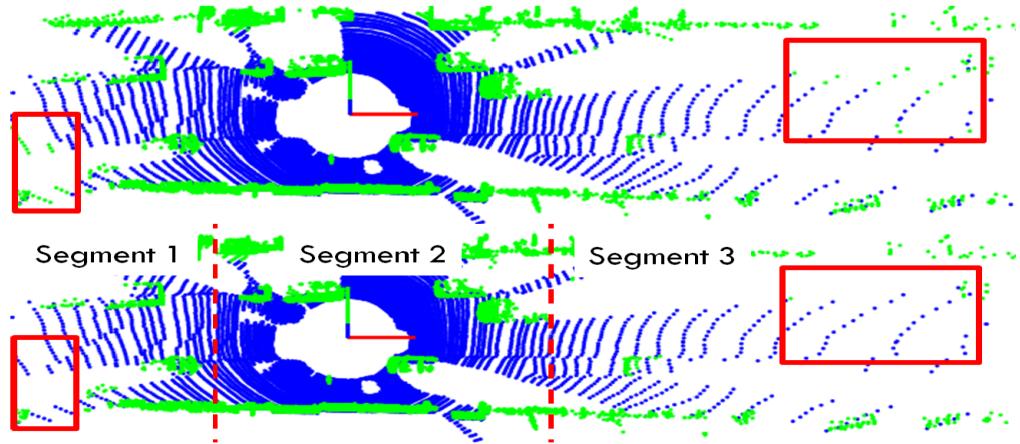


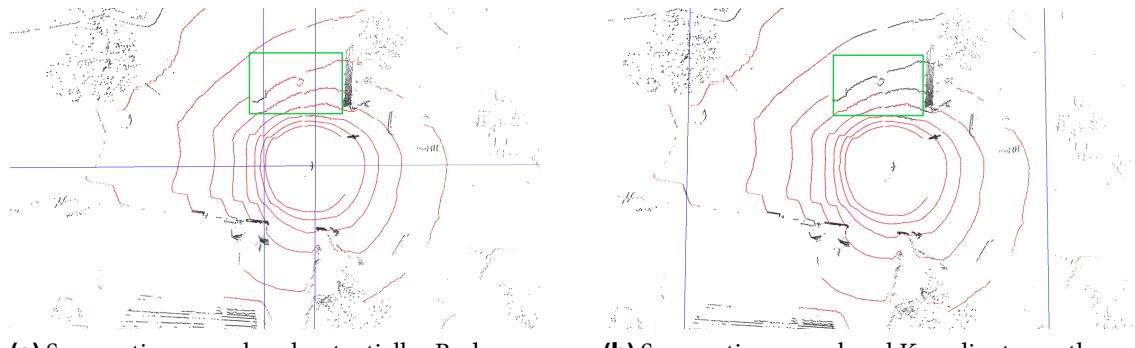
Abbildung 4.4.: In [4] wird gezeigt, dass die Segmentierung einer Punktwolke zu einem besseren Ergebnis der Bodenpunktanalyse führt. Die obere Abbildung ist dabei ohne Segmentierung und die untere mit Segmentierung. Bodenpunkte sind blau und Objektpunkte in grün dargestellt. Die roten Rechtecke kennzeichnen den Bereich in dem die Analyse schlechter bzw. besser ist.

Anschließend kann von jedem Punkt der Abstand zu dieser Ebene ermittelt werden. Ist der Abstand eines Punktes kleiner als ein definierter Wert, so handelt es sich bei ihm um einen Bodenpunkt.

Die Punktwolke, die analysiert werden soll, wird dabei zunächst in eine Anzahl an N_{segs} Segmenten aufgeteilt. Der Algorithmus zur Identifizierung der Bodenpunkte wird für jedes dieser Segmente ausgeführt. Der Grund für diese Segmentierung ist die mögliche Unebenheit der Bodenfläche. Soll eine wellige Oberfläche von einer Ebene repräsentiert werden kommt es zu ungenauen Ergebnissen. Wird aber für jedes Segment der Punktwolke eine eigene Ebene ermittelt so verbessert sich das Ergebnis der Bodenpunktanalyse, wie in Abbildung 4.4 zu sehen ist.

Beim GPF-Verfahren wurden Punktwolken in 3 große Segmente entlang der x -Achse eingeteilt, da diese Achse der Fahrtrichtung des Autos entspricht. Damit wurden für die Testdaten, die CMORE für diese Arbeit zur Verfügung stellte, keine zufriedenstellenden Ergebnisse erzielt (vgl. Abbildung 4.5b). Das Problem ist, dass sich durch eine zu simple Segmentierung entweder zu viele oder zu wenige Bodenpunkte in einem Sektor befinden können. Bei zu vielen Bodenpunkten kann es sein, dass der Boden nicht linear verläuft. Auf der linken Seite des Autos könnte sich beispielsweise eine Anhöhe befinden und auf der rechten Seite ein flaches Gelände. In diese Unebenheit kann keine richtige Ebene eingesetzt werden. Bei zu wenigen Bodenpunkten kann es sein, dass der Algorithmus die Ebene an Objekte annähert. Handelt es sich bei solchen Objekten beispielsweise um eine senkrechte Wand, können die vielen hohen Punkte die Ebene senkrecht werden lassen (siehe Abbildung 4.6).

Darum wurde in dieser Arbeit ein Verfahren entwickelt, dass die Segmente nicht anhand



(a) Segmentierung anhand potentieller Bodenpunkten

(b) Segmentierung anhand Koordinaten entlang der x -Achse

Abbildung 4.5.: Vergleich der Segmentierungsmethoden zur Bodenpunktanalyse aus [4] und C.LABEL-VR. Bodenpunkte sind rot, alle anderen schwarz. Die blauen Linien stellen die Trennung der Segmente dar. Die grünen Boxen markieren den Bereich, in dem man den Unterschied der beiden Methoden bei der Bodenpunkterkennung sieht.

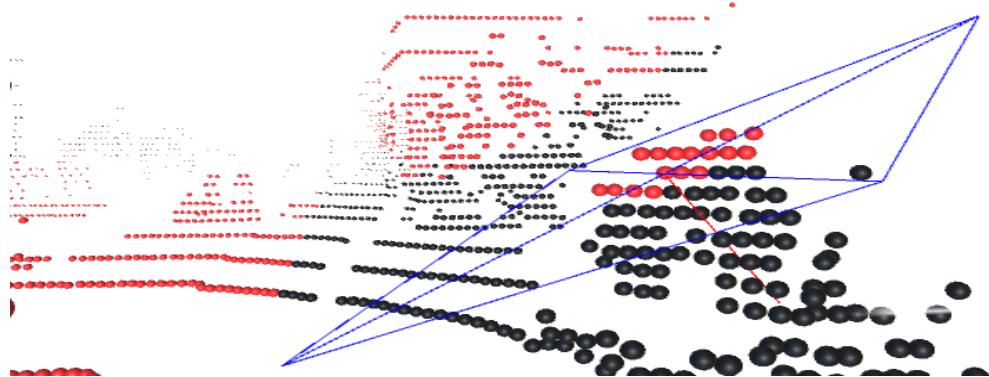


Abbildung 4.6.: Bei zu wenigen Bodenpunkten in einem Segment passt sich die Ebene den falschen Punkten an.

ihrer Größe einteilt sondern nach der Anzahl der potentiellen Bodenpunkte, welche die Segmente enthalten (siehe Algorithmus 1). Potentielle Bodenpunkte sind diejenigen Punkte, welche sich in einem bestimmten Höhenbereich befinden. Der Höhenbereich wurde auf -0.5m bis 0.5m festgelegt. Das bedeutet alle Punkte, deren Höhenwert sich innerhalb des Wertebereichs befindet, sind potentielle Bodenpunkte. Das Segmentierungsverfahren erstellt nun die Segmente so, dass sich in jedem Segment gleich viele dieser potentiellen Bodenpunkte befinden.

Dazu werden alle Punkte des gesamten Segments nach ihrem x -Wert sortiert und anschließend wird durch diese Menge an sortierten Punkten iteriert. Bei jeder Iteration wird geprüft ob der aktuell betrachtete Punkt ein potentieller Bodenpunkt ist oder nicht. Ist das der Fall wird eine Zählvariable erhöht. Übersteigt der Wert dieser Variable die Anzahl an potentiellen Bodenpunkten, die in einem Segment sein sollen, werden zukünftig betrachtete Punkte in ein neues Segment aufgenommen. Die Zählvariable wird anschließend wieder auf Null gesetzt um im neuen Segment die potentiellen Bodenpunkte erneut zu zählen.

Algorithm 1 Algorithmus zur Einteilung einer Punktwolke in N_{segs} Segmente

```

1: Inputs:
     $P$  : Menge aller Punkte
     $N_{segs}$  : Anzahl der Segmente
2: Initialize:
     $S$  : Menge aller Segmente mit ihren Punkten
     $P_{pot}$  : Menge aller potentiellen Bodenpunkte
     $C_{segs}$  : Zähler des aktuellen Segments
     $C_{points}$  : Zähler der Bodenpunkte des aktuellen Segments
     $N_{gpps}$  : Anzahl der potentiellen Bodenpunkte pro Segment
3: MainLoop:
4:  $P_{sorted} = \text{SORTASCENDINGONXVALUE}(P);$ 
5:  $P_{pot} = \text{GETPOINTSBETWEENHEIGHTVALUES}(-0.5, 0.5, P_{sorted});$ 
6:  $N_{gpps} = \left\lceil \frac{|P_{pot}|}{N_{segs}} \right\rceil;$ 
7:  $C_{segs} = C_{points} = 0;$ 
8: for  $i = 1 : |P_{sorted}|$  do
9:   if  $p_i.yValue \geq 0$  then
10:     $S[C_{segs}] \leftarrow p_i;$ 
11:   else
12:     $S[C_{segs} + 1] \leftarrow p_i;$ 
13:   end if
14:   if  $p_i \in P_{pot}$  then
15:      $C_{points} ++;$ 
16:   end if
17:   if  $C_{points} \geq N_{gpps}$  then
18:      $C_{points} = 0;$ 
19:      $C_{segs} += 2;$ 
20:   end if
21: end for

```

Zusätzlich zu der Segmentierung aus [4] wurden die Segmente noch durch die y -Achse

geteilt, um den vorher angesprochenen Problemfall der unterschiedlichen Geländebeschaffenheiten links bzw. rechts vom Auto auszugleichen. Dazu wird bei jeder Iteration der y -Wert des betrachteten Punktes geprüft. Punkte mit einem größeren Wert als Null werden in ein anderes Segment aufgenommen als solche mit kleinerem Wert. Das Ergebnis und der direkte Vergleich mit der normalen GPF-Methode ist in Abbildung 4.5a zu sehen.

Nach dem die Segmentierung der Punktfolge vorgenommen wurde, muss nun die Bodenebene für jedes Segment berechnet werden. Die Definition dieses Vorgehens ist in Algorithmus 2 gegeben. Auch dafür wurde die GPF-Methode aus [4] als Basis hergenommen. Zuerst werden gewisse Startpunkte ermittelt, welche auch *Seedpoints* bzw. *Seed-Punkte* genannt werden. Laut GPF werden diese Punkte durch den *LPR*-Wert (*lowest point representative*) ermittelt. Um diesen Wert zu berechnen wird zunächst eine Anzahl an N_{LPR} Punkten definiert. Danach werden die N_{LPR} tiefsten Punkte des Segments genommen und der Durchschnitt der Höhenwerte dieser Punkte berechnet. Alle Punkte des Sektors die eine maximale Höhendifferenz T_{init} zu diesem Wert haben sind Seed-Punkte.

Dieses Vorgehen ist allerdings anfällig für Rauschen in Form von sehr tiefen Punkten. Befinden sich im Segment also fälschlicherweise Punkte, die tief unter dem Boden sind, verfälschen diese die richtige Berechnung des Durchschnittshöhenwerts. Deshalb wurden in C.LABEL-VR für die Berechnung Durchschnittshöhenwerts LPR nicht die N_{LPR} tiefsten Punkte hergenommen, sondern alle potentiellen Bodenpunkte die sich in dem Segment befinden. Aus diesem Durchschnittshöhenwert und der maximal zulässigen Differenz T_{init} werden, wie schon beschrieben, die Startpunkte ermittelt. Die Formale Beschreibung dieses Vorgehens ist ebenfalls in Algorithmus 2 zu sehen (GetInitialSeedPoints).

Mit diesen Startpunkten kann nun die erste von N_{iter} Schätzungen für die Bodenebene berechnet werden. Eine Ebene E kann definiert werden durch einen Vektor \vec{n} der senkrecht auf der Ebene steht (Normalenvektor) und einem Punkt p der auf der Ebene liegt. Ziel ist es nun diese beiden Komponenten so zu berechnen, dass die dadurch definierte Ebene bestmöglich in Seed-Punkte passt. Für den Punkt p kann man den einfach Durchschnitt \hat{s} aus allen Elementen der Menge aller Seed-Punkte S nehmen. Dieser Punkt repräsentiert den Ursprung der Ebene und bietet sich deshalb an weil man ihn bei späteren Berechnungen wiederverwenden kann (siehe Gleichung 4.1).

Die Berechnung des Normalenvektors \vec{n} ist dagegen deutlich aufwändiger. Zunächst wird die Streuung aller Punkte $s \in S$ untersucht. Dazu wird eine Kovarianzmatrix C des Ranges R^{3x3} berechnet, welche die Streuung der Seedpoints repräsentiert. Diese Matrix kann nun auf geometrische Eigenschaften untersucht werden, beispielsweise welcher Vektor die kleinste Streuung repräsentiert und somit als \vec{n} verwendet werden kann.

$$C = \sum_{i=1:|S|} (s_i - \hat{s})(s_i - \hat{s})^T \quad (4.1)$$

Im Allgemeinen gilt für eine Matrix Folgendes: „Eine Matrix definiert durch $y = A * x$ eine lineare Abbildung $R^m \rightarrow R^n$. Der Hauptberuf einer Matrix ist, aus einem Vektor einen anderen zu machen. Im Allgemeinen ändert die Matrix dadurch Richtung, Betrag (und sogar Dimension) eines Vektors. Matrixzerlegungen spalten die Aktion der Matrix in leichter zu durchblickende Einzelschritte auf“ [30]. Aus diesen Einzelschritten können geometrische Informationen extrahiert werden. Bei der GPF-Methode wird C durch die Singulärwertzerlegung gespalten. Dieses Verfahren wird auch *Singular Value Decomposition* (SVD) genannt. Wie in [31] gezeigt, wird bei SVD eine Matrix $A \in R^{m \times n}$ in 2 orthogonale Matrizen $U \in R^{m \times m}$ und $V \in R^{n \times n}$ und eine diagonal Matrix Σ zerlegt, sodass die Gleichung 4.2 erfüllt wird. Der Normalenvektor \vec{n} der gesuchten Ebene ist gegeben durch die dritte Spalte der Matrix U [32].

$$A = U\Sigma V^T \quad (4.2)$$

Die gesuchte Ebene ist nun durch den Punkt p und den Vektor \vec{n} nun eindeutig definiert. Nun kann der Abstand jedes Punktes des Segments zu seiner orthogonalen Projektion auf der definierten Ebene berechnet werden. Dieser Abstand wird dann mit einem Schwellwert T_{plane} verglichen, welcher zu Beginn der Berechnung definiert werden muss. Ist der Abstand eines Punktes kleiner als der Schwellwert so handelt es sich bei ihm um einen Bodenpunkt. In C.LABEL-VR wird solch ein Punkt vorläufig als Bodenpunkt markiert und in die Menge an neuen Seed-Punkten aufgenommen. Mit diesen neuen Seed-Punkten wird anschließend die nächste von N_{iter} Berechnungen der Bodenebene getätigt(vgl. Algorithmus 2).

Cluster-Analyse

TODO Erklärung des Cluster-Algorithmus (Recursive Radius Search)

4.5. User Interface

TODO -kurze Einleitung -Erklärung der Architektur hinter dem UI-System

4.5.1. Ingame Menu

TODO Erklärung der Einstellbaren Parameter im Ingame Menü

Algorithm 2 Algorithmus zur Identifizierung von Bodenpunkten eines Punktwolken-Segments

- 1: **Inputs:**
 \mathbf{P} : Menge aller Punkte der Wolke
- 2: **Initialize:**
 \mathbf{P}_{seeds} : Menge aller Seed-Punkte
 N_{iter} : Anzahl der Iterationen
 T_{init} : Distanzgrenzwert der initialen Seed-Punkte
 T_{plane} : Distanzgrenzwert zur geschätzten Ebene
- 3: **MainLoop:**
 - 4: $\mathbf{P}_{seeds} = \text{GETINITIALSEEDPOINTS}(\mathbf{P}, T_{init});$
 - 5: **for** $i = 1 : N_{iter}$ **do**
 - 6: $plane = \text{FITPLANEINTOPOLY}(P_{seeds});$
 - 7: $\text{CLEAR}(\mathbf{P}_{seeds});$
 - 8: **for** $j = 1 : |\mathbf{P}|$ **do**
 - 9: $distance = \text{GETDISTANCETOPLANE}(plane, p_k);$
 - 10: **if** $distance < T_{plane}$ **then**
 - 11: $p_k.groundpoint = true;$
 - 12: $\mathbf{P}_{seeds} \leftarrow p_k$
 - 13: **else**
 - 14: $p_k.groundpoint = false;$
 - 15: **end if**
 - 16: **end for**
 - 17: **end for**
 - 18:
- 19: **GetInitialSeedPoints:**
- 20: **Initialize:**
 \mathbf{P}_{low} : Alle tiefen Punkte aus \mathbf{P}
 H_{avg} : Durchschnittshöhe aller Punkte aus \mathbf{P}_{low}
- 21: **for** $i = 1 : |\mathbf{P}|$ **do**
- 22: **if** $p_i.height < 0.5 \ \& \ p_i.height > -0.5$ **then**
- 23: $\mathbf{P}_{low} \leftarrow p_i$
- 24: **end if**
- 25: **end for**
- 26: $H_{avg} = \text{GETAVERAGEHEIGHT}(\mathbf{P}_{low});$
- 27: **for** $i = 1 : |\mathbf{P}|$ **do**
- 28: **if** $p_i.height < T_{init}$ **then**
- 29: $\mathbf{P}_{seeds} \leftarrow p : i;$
- 30: **end if**
- 31: **end for**

Schluss

5.1. Zusammenfassung

5.2. Herausforderungen

5.3. Ausblick

A

Appendix Title

Literaturverzeichnis

- [1] Microsoft, “Machine learning – cheat sheet für algorithmen für microsoft azure machine learning studio,” Internetquelle(abgerufen am 21.03.2018), Dec. 2017.
- [2] S. S. Fisher, M. McGreevy, J. Humphries, and W. Robinett, “Virtual environment display system,” in *Proceedings of the 1986 workshop on Interactive 3D graphics - SI3D '86*. ACM Press, 1987.
- [3] P. Chu, S. Cho, S. Sim, K. Kwak, and K. Cho, “A fast ground segmentation method for 3d point cloud,” *Journal of information processing systems*, vol. 13, no. 3, pp. 491–499, 2017.
- [4] D. Zermas, I. Izzat, and N. Papanikolopoulos, “Fast segmentation of 3d point clouds: A paradigm on LiDAR data for autonomous vehicle applications,” in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, may 2017.
- [5] D. Kriesel, *Ein kleiner Überblick über Neuronale Netze*, 2007. [Online]. Available: erhältlich auf <http://www.dkriesel.com>
- [6] C. V. D. Malsburg, “Frank rosenblatt: Principles of neurodynamics: Perceptrons and the theory of brain mechanisms,” in *Brain Theory*. Springer Berlin Heidelberg, 1986, pp. 245–248.
- [7] I. E. Sutherland, “The ultimate display,” in *Proceedings of the IFIP Congress*, 1965, pp. 506–508.
- [8] G. Enrico, “Virtual reality: Past, present and future,” *Studies in Health Technology and Informatics*, vol. 58, no. Virtual Environments in Clinical Psychology and Neuroscience, pp. 3–20, 1998.
- [9] C. Gammage, “Exploring oculus rift: A historical analysis of the ‘virtual reality’paradigm,” *ART 108: Introduction to Games Studies*, 2017.
- [10] K. M. Johnson and W. B. Ouimet, “Rediscovering the lost archaeological landscape of southern new england using airborne light detection and ranging (LiDAR),” *Journal of Archaeological Science*, vol. 43, pp. 9–20, mar 2014.

- [11] D. E. Smith, “The global topography of mars and implications for surface evolution,” *Science*, vol. 284, no. 5419, pp. 1495–1503, may 1999.
- [12] S. K. Reddy and P. K. Pal, “Segmentation of point cloud from a 3d LIDAR using range difference between neighbouring beams,” in *Proceedings of the 2015 Conference on Advances In Robotics - AIR '15*. ACM Press, 2015.
- [13] P. Gothard, “Microsoft hololens: Firm admits sales figures are in the ’thousands,” Internetquelle (abgerufen am 24.11.2017), Jan. 2017. [Online]. Available: <https://www.theinquirer.net/inquirer/news/3003380/microsoft-hololens-firm-admits-sales-figures-are-in-the-thousands>
- [14] Microsoft, “Microsoft hololens choose the option that best suits your purpose.” Internetquelle (abgerufen am 27.11.2017), Nov. 2017. [Online]. Available: <https://www.microsoft.com/de-de/hololens/buy>
- [15] R. Berg, “Touch-controller machen oculus rift zur besten vr-brille,” Internetquelle (abgerufen am 24.11.2017), Dec. 2016. [Online]. Available: <https://www.welt.de/wirtschaft/webwelt/article160455301/Touch-Controller-machen-Oculus-Rift-zur-besten-VR-Brille.html>
- [16] J. Munafo, M. Diedrick, and T. A. Stoffregen, “The virtual reality head-mounted display oculus rift induces motion sickness and is sexist in its effects,” *Experimental Brain Research*, vol. 235, no. 3, pp. 889–901, dec 2016.
- [17] W. M. R. D. Forum, “Creating a new gesture,” Internetquelle (abgerufen am 28.11.2017), Apr. 2016. [Online]. Available: <https://forums.hololens.com/discussion/549/creating-a-new-gesture>
- [18] J.-K. Janssen, “Oculus rift im test: Virtual reality für die massen,” Internetquelle (abgerufen am 27.11.2017), Mar. 2016. [Online]. Available: <https://www.heise.de/ct/artikel/Oculus-Rift-im-Test-Virtual-Reality-fuer-die-Massen-3151909.html>
- [19] ——, “Oculus touch im test: So echt können sich händende in vr anfühlen,” Internetquelle (abgerufen am 27.11.2017), Dec. 2017. [Online]. Available: <https://www.heise.de/ct/artikel/Oculus-Touch-im-Test-So-echt-koennen-sich-Haende-in-VR-anfuehlen-3552295.html>
- [20] R. Artus, “Lighthouse von valve erklärt,” Internetquelle (abgerufen am 28.11.2017), Nov. 2017. [Online]. Available: <https://vrjump.de/lighthouse-erklaert>
- [21] L. Painter, “Htc vive review,” Internetquelle (abgerufen am 28.11.2017), Aug. 2017. [Online]. Available: <http://www.techadvisor.co.uk/review/wearable-tech/htc-vive-review-2017-3635648/>

- [22] O. R. Support, “Was sind die empfohlenen mindestsystemvoraussetzungen, die für den betrieb der oculus rift erforderlich sind?” Internetquelle (abgerufen am 29.11.2017), Nov. 2017. [Online]. Available: <https://support.oculus.com/170128916778795/>
- [23] P. Gora and L. Leibetseder, “Unreal vs unity: Ein vergleich zwischen zwei modernen spiele-engines,” Favoritenstrasse 9-11/186, A-1040 Vienna, Austria, Oct. 2016, 1. [Online]. Available: https://www.cg.tuwien.ac.at/research/publications/2016/Przemyslaw_Gora_2016_UVU/
- [24] U. Technologies, “Die weltweit führende software der game-industrie,” Internetquelle (abgerufen am 24.04.2018), Apr. 2018. [Online]. Available: <https://unity3d.com/de/public-relations>
- [25] Unity, “Importing objects from blender,” Internetquelle (aufgerufen am 04.05.2018), May 2018. [Online]. Available: <https://docs.unity3d.com/560/Documentation/Manual/HOWTO-ImportObjectBlender.html>
- [26] ——, “Monobehaviour,” Internetquelle (abgerufen am 01.05.2018), May 2018. [Online]. Available: <https://docs.unity3d.com/ScriptReference/MonoBehaviour.html>
- [27] C. Games, “Unity3d tips 2 – the singleton pattern,” Internetquelle (abgerufen am 01.05.2018). [Online]. Available: <http://clearcutgames.net/home/?p=437>
- [28] M. Hoy, J. Dauwels, and J. Yuan, “Efficient ground object segmentation in 3d lidar based on cascaded mode seeking,” in *2017 IEEE 20th International Conference on Intelligent Transportation Systems (ITSC)*, Oct 2017, pp. 1–6.
- [29] S. Goga and S. Nedevschi, “An approach for segmenting 3d lidar data using multi-volume grid structures,” in *2017 13th IEEE International Conference on Intelligent Computer Communication and Processing (ICCP)*, Sept 2017, pp. 309–315.
- [30] E. H. Clemens Brand, “Matrixzerlegungen. Überbestimmte systeme,” Vorlesungsskript, Montanuniversität Leoben, Mar. 2014.
- [31] D. Kalman, “A singularly valuable decomposition: The svd of a matrix,” *College Math Journal*, vol. 27, pp. 2–23, 1996.
- [32] I. Söderkvist, “Using svd for some fitting problems,” Lulea University of Technology.