

1 Einführung

1.1 Intension

Zur Bewertung der Qualität des Unterrichtes, wird die Aufmerksamkeit der Schüler verwendet, da zwischen beiden ein Zusammenhang besteht. Allerdings ist der Parameter Aufmerksamkeit recht schwierig zu erfassen, wodurch verschiedene Verfahren verwendet werden. Unter anderem Fragebögen die ein Schüler selbst ausfüllen sollen oder die Auswertung durch einen Beobachter der bewertet ob ein einzelner Schüler Aufmerksam (on-Task) oder nicht (off-Task) ist.

Für die Bewertung ob on/off-Task werden Kriterien festgelegt, wie Blickrichtung, Körperhaltung und Tätigkeit, dann wird die Person beobachtet wie diese sich verhält.

Bei der Videostudie zur Wirksamkeit des Unterrichtsprozesses [App15] wurden die Kriterien Blickkontakt zum legitimen Sprecher oder Objekt, Aktive Beteiligung an der Aufgabe, keine Ausübung anderer Tätigkeiten, keine Motorische Unruhe und keine Themenferne Kommunikation festgelegt. Dann wird immer in einem 1min Intervalle der Schüler beobachtet und bewertet. Sollte drei oder mehr Punkte erfüllen, gilt der Schüler als on-Task.

Diese Art der Auswertung ist recht einfach, allerdings gibt es Interpretationsfreiheiten, gerade bei den Bewertungen der Tätigkeiten, die von jedem Beobachter anders gewertet werden. Außerdem ist es sehr zeitintensiv, so werden alleine zum anschauen des Videos für jeden Schüler bei einer Klassen (30 Personen) 30 min gebraucht und sollte jeder Schüler mehrfach beobachtet werden entsprechend mehr. Sollten recht wenige Zyklen durchgeführt werden, so wird das gesamte Verhalten eines Schülers in der Unterrichtsstunde mit nur wenigen beobachteten Minuten beschrieben, die Auswertung benötigt allerdings entsprechend weniger Zeit.

Durch eine zu geringe Auswertungsrate kann nur eine Aussage über den gesamten Unterricht gemacht werden und nicht über einzelne Übungen oder ähnliches, auch die Bewertung eines einzelnen Schülers ist nur schwer möglich. [App15]

1.2 Problemstellung

Die aktuellen Verfahren zur Analyse von Aufmerksamkeit im Unterricht werden meist äquivalent zu Abschnitt 1.1 durchgeführt. Der jeweilige Schüler wird von einer Person beobachtet und dann nach bestimmten Kriterien bewertet. In die Bewertung fließt allerdings auch die Meinung/Auffassung des Bewerters ein und ist daher nicht perfekt objektiv.

Diese Art der Auswertung ist recht ungenau und Arbeitsintensiv, da sie von einer Person ausgeführt wird. Ein wichtiger Parameter ist die Blickrichtung des einzelnen Schülers, da sie meist dorthin gerichtet ist, wo auch die Aufmerksamkeit liegt.[HR92]

Ziel ist en nun mit möglichst geringem Aufwand an Hardware eine Bestimmung der Blickrichtung einer ganzen Klasse vorzunehmen. Die Messung soll den Unterricht möglichst wenig beeinträchtigen, wodurch Eye-Tracking Brillen nicht verwendet werden, wegen den Kosten und der Ablenkung. Auch der Aufbau soll recht einfach und für Laien anwendbar sein, somit wird nur eine festmontierte Kamera vor der Klasse eingesetzt.

Für diese Anforderungen soll nun ein Verfahren entwickelt werden, mit dem es möglich ist das Film-

material von einer gesamte Klasse auf einmal auszuwerten, um von allen Personen die Blickrichtungen bzw. die Gesichtsorientierung während einer Schulstunde zu bestimmen.

1.3 Gesetzte Bedingungen der Anwendung

Damit der Unterricht, wie im Szenario der Problemstellung beschreiben Abschnitt 1.2, möglichst wenig beeinflusst wird, ergeben sich folgende Randbedingungen:

- Brillen, Kontaktlinsen und ähnliches sind erlaubt.
- Die üblichen Bewegungen im Unterricht wie Sprechen, Kopfdrehungen usw. der Schüler sind gestattet.
- Es soll gleichzeitig auf Distanzen zwischen $2.5 - 8m$ zur Kamera auf einer Breit von $6m$ funktionieren.
- Möglichst alle Blickrichtungen der Schüler sollen so exakt wie möglich erfasst werden.

Ein deutsches Klassenzimmer hat $55 - 65m^2$, da noch Abstand zur Tafel usw. beachtet werden muss ergibt sich, wenn sich die Kamera an der Tafel befindet, einen Abstand von $2.5 - 8m$ zwischen Kamera und Schüler auf einer Breiten von $6m$. Somit muss der Linsenwinkel mindestens 100° betragen, damit alle im Fokus sind.

Außerdem soll die Anwendung auf schon vorhanden Aufnehmen eines Unterrichtes arbeiten, bei denen oben genannten Bedingungen erfüllt sind.

To Do

Quelle für Klassenzimmer

1.3.1 Randbedingungen der Anwendung

Des weiteren werden folgende Annahmen gemacht:

- Die Szene ist Innerhalb eines Gebäudes, mit ausreichend gleichmäßiger Beleuchtung.
- Die Überführung zwischen Welt- und Kamerakoordinatensystem bekannt.
- Die Kamera befindet sich vor der Klasse, so das die Hauptblickrichtung der Schüler in ihrem Fokus liegt.
- Die Gesichter sind komplett sichtbar und nicht verdeckt.

Natürlich sind auch alle inneren Parameter der Kamera bekannt.

1.4 Hardware

Als Messinstrument wird nur eine einzelne Farbkamera eingesetzt. Das Videomaterial der Schulklassen wurde mit einer unbekannten Videokamera aufgezeichnet, daher sind nur die Parameter des Filmes (640×480 mit $??Hz$) bekannt.

Für die Messungen im Versuch wurde die Explorer 4K Action Camera verwendet, sie besitzt eine 170° Weitwinkel-Linse mit großem Feald of View. Mit der 2.7K Einstellung wird ein 2688×1520 Video mit 30FPS aufgezeichnet. Sie wurde fest montiert.

1.5 Software

Für die Umsetzung werden folgende Software-Elemente aus fremder Quelle eingesetzt.

1.5.1 ElSe

Ellipse Selection for Robust Pupil Detection in Real-World Environments, ein Algorithmus zur Bestimmung der Pupille im Bild. Der Ursprüngliche ElSe-Algorithmus ist für Graubilder mit Infrarotbeleuchtung ausgelegt, wurde für diese Anwendung auf Farbbilder modifiziert.

Entwickelt von der Uni Tübingen. [WF16]

1.5.2 MTCNN Face Detection

Multi-task Cascaded Convolutional Networks, ein Algorithmus zur Detektion von Gesichtern und Bestimmung von 5 Gesichts-Landmarks in Farbbildern. Dabei werden drei CNN auf einer Bildpyramide angewendet um so zuverlässig Gesichter mit verschiedenster Größe zu erkennen.

[KZ15]

1.5.3 OpenCV

Open Source Computer Vision, ist eine C/C++ Bibliothek von Algorithmen zur Bildverarbeitung in Echtzeit unter der BSD Lizenz (Berkeley Software Distribution)

[Wik15][BK08]

1.5.4 OpenFace

Ein Open-Source Echtzeitverfahren auf Basis von CLNF zur Bestimmung und Analyse von Gesichtsmerkmalen in Grau-Bildern und Videos. Dabei werden 68 signifikante Punkte im Gesicht bestimmt und auf Basis jener Position und Orientierung ermittelt.

Entwickelt von der University of Cambridge [TB16]

2 Theorie & Grundlage

2.1 Grundlagen

Die Gesichtserkennung ist Teil der Bildverarbeitung und wird ständig weiterentwickelt. Darunter fallen neben der Detektion des Gesichtes auch seine Analyse wie Orientierung oder das erkennen von Mimiken und Übereinstimmungen.

2.1.1 Künstliches neuronales Netz

Ein künstliches neuronales Netz besteht aus miteinander verknüpften künstlichen Neuronen. Jedes Neuron erhält Eingangswerte, diese werden individuelle Gewichtet, mittels einer Übertragungsfunktion zusammengefasst und durch eine Schwellenwertfunktion das Ergebnis bestimmt.

Um die Parameter der Neuronen zu bestimmen, werden sie zufällig initialisiert und dann so angepasst, dass sie zu einer gegebenen Eingabe das gewünschte Ergebnis liefert und der Fehler über dem gesamten Trainingsdatensatz minimal ist.

To Do

Quelle

2.1.2 Convolutional Neural Network (CNN)

CNN ist eine Weiterentwicklung der neuronalen Netzen und werden zur Klassifizierung verwendet, unter anderem im Bereich Bild- und Spracherkennung. Dies wird durch eine gewichtete Faltung der Eingabe erreicht und sind state of the art bei vielen Anwendungen.

Durch die Faltung werden die Information aus den umliegenden Punkten eines Bereiches zusammengefasst und komprimiert an die nächste Schicht weitergegeben, um in der untersten Schicht alle vorhandenen Informationen zusammenzuführen. Der Faltungskern kann ja nach Anwendung beliebig gestaltet sein, so ist eine Glättung durch einen Gauß-Kernel oder Kantendetektion durch einen Kirsch-Operator möglich.

Ein CNN kann in zwei Bereiche aufgeteilt werden, der Feature Extraktion in welcher durch verschiedene Kernel und Komprimierung die Eingabeinformationen zur Klassifizierung, dem zweiten Bereich, aufbereitet. Gelernt werden können die Kernel an sich und die jeweiligen Bewertungen.

To Do

Quelle

Bild

2.1.3 Constrained Local Model (CLM)

Dies ist ein Verfahren um mehrere Punkte eines Objektes zu lokalisieren. Dabei wird eine Wahrscheinlichkeitskarte für jeden einzelnen Punkt erstellt, wo dieser sich aufhalten kann auf Basis eines

Trainingsdatensatzes. Nun wird versucht für das Bild, auf welchem gerechnet werden soll, für jeden Punkt den maximalen Wert zu erreichen zwischen passendem Farbverlauf und Wahrscheinlichkeit. Dieser Art der Bestimmung von Punkten mit Positionsabhängigkeiten ist ziemlich zuverlässig und dennoch dynamisch genug um auch mit kleinen Veränderungen klar zu kommen. Dies ist Wichtig, bei der Detektion von verschiedenen leicht verformbaren Objekten wie Gesichter und daher zuverlässiger als das Active Appearance Model (AAM).

To Do

Quelle
AAM
zur Detecton der Landmarks

2.1.4 PDM & GAVAM

Mit Point Distribution Model (PDM) können verformbare Objekte recht gut modelliert werden. Dabei wird die durchschnittliche Form \bar{X} bestimmt und eine Matrix P von Eigenvektoren ermittelt, um die möglichen Deformierungen darzustellen.

$$X = \bar{X} + P \cdot b$$

Somit kann durch einen Skalierungsvektor b alle möglichen Formen X des Objektes dargestellt werden. Zur Vereinfachung reicht es, die signifikantesten Eigenvektoren in P auf zu nehmen und dennoch X ausreichend genau zu beschreiben.

Ist bekannt welche Art der Verformung durch den Eigenvektor dargestellt ist, z.B. eine bestimmte Orientierung, so kann anhand des Skalierungsvektors die Rotation des berechneten Objektes bestimmt werden, siehe Generalized Adaptive View-based Appearance Model (GAVAM). Eine Problematik bei dieser Art der Bestimmung der Rotation entsteht, wenn neben der Verschiebung der Landmarks durch die Rotation, auch eine Deformierung des Objektes stattgefunden hat und somit keine eindeutige Lösung gefunden werden kann. Dies ist eine Problematik wenn auf Gesichtern gerechnet wird, da immer eine Veränderung der Mundwinkel oder Augenlider vorhanden ist.

[Wik17][Kyb07][MWM08]

2.1.5 Non-maximum suppression (NMS)

Ein Verfahren um Kanten in einem Bild exakter zu bestimmen. Dabei wird der Farbwert des Pixels mit dem umliegenden verglichen und sollte es nicht maximal sein auf Null gesetzt.

Auf diese Weise bleibt nur noch ein Kantenpixel übrig.

To Do

Quelle

2.2 MTCNN Face Detection

Bei Multi-task Cascaded Convolutional Network handelt es sich um ein Verfahren dass bei der Detektion von Gesichtern auch deren Ausrichtung berücksichtigt wird, um so bessere Ergebnis zu erzielen.

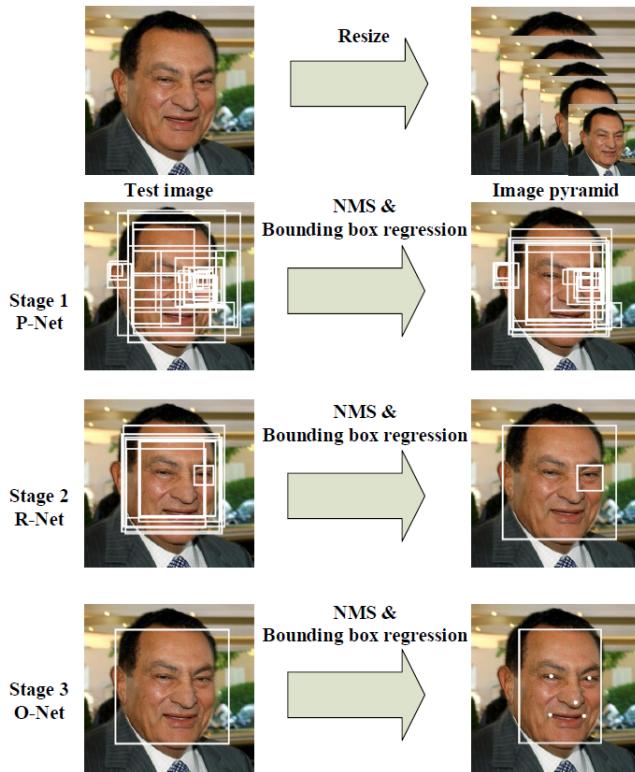


Abbildung 2.1: Darstellung des Funktionsablaufes von MTCCN[KZ15]

2.2.1 Anforderungen

Sein Einsatzgebiet ist die Vorverarbeitung eines Frames für die spätere Auswertung. Somit soll dieser Schritt von einem möglichst robusten Verfahren zur Detektion von Gesichtern durchgeführt werden. Dabei wird auf recht großen Bild gearbeitet mit verhältnismäßig kleinen und verschiedenen großen Gesichtern.

Außerdem sollte das Verfahren ausreichend schnell sein, da es sich hierbei nur um ein Vorverarbeitungsschritt handelt und zur Beschleunigung der späteren Berechnung beitragen soll.

2.2.2 Die 3 Stufen der Verarbeitung

Für die gute Detektionsqualität sorgt die dreistufige Verarbeitung auf der Bildpyramide. Bei der Bildpyramide handelt es sich um ein in verschiedenen Größen skaliertes Bild, damit der Gesuchte Inhalt in der gewünschten Auflösung abgebildet ist, ohne dass etwas über den Inhalt bekannt ist. Dies ist von Vorteil, damit das CNN auf eine feste Größe von Gesichtern optimiert werden kann, um neben dem möglichen Farbverläufen, durch die Skalierung das Lernen nicht zusätzlich zu erschweren.

Stufe 1

Beim ersten Verarbeitungsschritt werden alle Bereiche eines Bilds gesucht, in denen möglicherweise ein Gesicht zu sehen ist. Dazu wird zuerst ein einfaches CNN eingesetzt und die Ergebnisse, die sich sehr stark überlappen, zusammengefasst.

Für die Detektion wird von einem CNN, dem sogenannte Proposal Network (P-Net), eingesetzt und

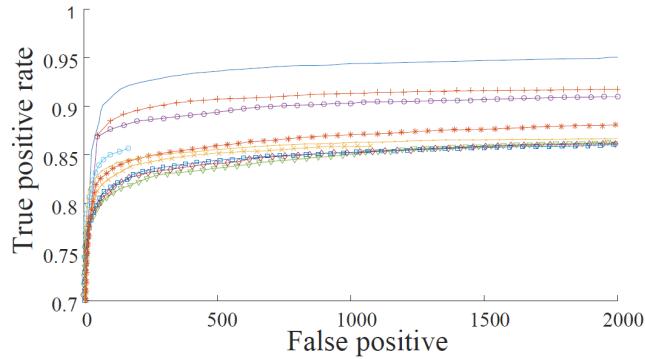


Abbildung 2.2: normale blaue Linie[KZ15]

sehr viele Bounding-Boxen ermittelt. Diese werden nun mit einem NMS ausgedünnt, um die am stärksten überlappenden Bounding-Boxen zusammen zu fassen.

Stufe 2

Anschließend werden die möglichen Bereiche mittels eines weiten CNN analysiert, damit alle Nicht-Gesichtsbereiche erkannt und entfernt werden können.

Dies wird von dem Refine Network (R-Net) übernommen und anschließend die möglichen Bounding-Boxen mittels NMS weiter reduziert.

Stufe 3

Der letzte Schritt wird von einem deutlich genaueren CNN übernommen, um ein Gesicht zu detektieren, dem sogenannten Output Network (O-Net). Womit die resultierenden exakten Boxen und 5 Landmarks ermittelt werden.

2.2.3 Qualität

MTCNN Face Detection ist bei der Zuverlässigkeit im Vergleich zu anderen bekannten Verfahren überlegen, siehe Abbildung 2.2 und zudem Echtzeit fähig. Im Test-Datensatz sind auch Gesichtern mit einer Größe von 20×20 enthalten und wurden erfolgreich erkannt.

Somit sind alle Anforderungen erfüllt um mit diesem Verfahren den vorhanden Frame für die nachfolgenden Berechnungen vorzubereiten, daher wird es auch hier eingesetzt.

2.3 Skalieren von Bildern

Da die Berechnungen meist auf recht kleinen Bildausschnitten ausgeführt wird, müssen diese für weitere Rechenschritte hochskaliert werden, damit es von OpenFace besser verarbeitet wird.

Dabei ist es wichtig, dass die Gesichtsmerkmale möglichst gut rekonstruiert werden, um die entsprechenden Landmarks zu bestimmen.

2.3.1 Nearest-Neighbor

Dieses Verfahren verwendet als neuer Farbwert, den gleichen Wert wie das nächstgelegene Pixel. Dadurch werden nur die ehemaligen Pixel größer und das Gesicht wirkt sehr Kantig, da keine neuen



Abbildung 2.3: Die ursprüngliche Abbildung von Lena betrug 100 Pixel Kantenlänge und beim Schachbrett 48 Pixel, beide wurden mittels Nearest-Neighbor auf 512 Pixel vergrößert und bei Lena die Differenz bestimmt

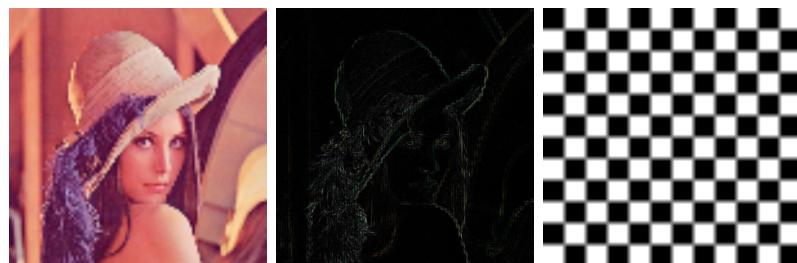


Abbildung 2.4: Die ursprüngliche Abbildung von Lena betrug 100 Pixel Kantenlänge und beim Schachbrett 48 Pixel, beide wurden mittels linearer Interpolation auf 512 Pixel vergrößert und bei Lena die Differenz bestimmt

Farbwerte bestimmt werden, siehe 2.3.

2.3.2 Linear

Um den neuen Farbwert zu ermitteln, wird zwischen den nächst gelegenen umliegenden Pixel linear Interpoliert, wodurch weitere Farbwerte entstehen. Das Ergebnis ist gleichmäßiger als Neares Neighbor, und dennoch ein recht einfaches Verfahren. Die Kanten wirken allerdings unscharf, siehe 2.4.

2.3.3 Bicubic

Um den Farbwert zu ermitteln, werden die umliegenden 4×4 Pixelwerte betrachtet um den Farbverlauf als eine Funktion 3. Grades zu bestimmen. Somit werden feinere Details besser dargestellt als beim linearen Verfahren und Kanten bleiben eher erhalten. Allerdings kann es durch den bestimmten Verlauf auch zum Überschwingen kommen, wodurch Fehlfarben entstehen können. Ein Beispiel ist in 2.5 zu sehen. [Wik16a]

2.3.4 Lanczos

Dieser Filter basiert auf einer Sinc-Funktion über einen Bereich, um so eine Bewertung der benachbarten Pixelwerte zu erhalten. Somit ergibt sich der neue Farbwert aus den bewerteten umliegenden Pixeln, wobei durch hie Fehler entstehen können, siehe 2.6. Die Funktion kann und wird für die An-



Abbildung 2.5: Die ursprüngliche Abbildung von Lena betrug 100 Pixel Kantenlänge und beim Schachbrett 48 Pixel, beide wurden mittels bikubischem Verfahren auf 512 Pixel vergrößert und bei Lena die Differenz bestimmt



Abbildung 2.6: Die ursprüngliche Abbildung von Lena betrug 100 Pixel Kantenlänge und beim Schachbrett 48 Pixel, beide wurden mittels Lanczus-Verfahren auf 512 Pixel vergrößert und bei Lena die Differenz bestimmt

wendung auf einen 8×8 Bereich begrenzt. [Wik16b]

$$L(x) = \begin{cases} \frac{\sin(\pi x)}{\pi x} \cdot \frac{\sin(\pi \frac{x}{a})}{\pi \frac{x}{a}} & \text{wenn } -a < x \geq a, a \leq 0 \\ 1 & \text{wenn } x = 0 \\ 0 & \text{sonst} \end{cases}$$

2.4 OpenFace

Die Aufgaben von OpenFace ist die Analyse der Gesichtes basierend auf Bildern. Dabei sind für die Anwendung nur Kameraparameter bekannt und keinerlei Zusätze wie eine Tiefenbild oder Infrarotbeleuchtung der Szene vorhanden. Dabei ist für die Anwendung die Kopfposition (Translation und Orientierung) und Blickrichtung von Interesse, da mit ihnen zurückrechnet werden kann wohin die Person schaut.

OpenFace kann neben den Landmarks auch die Position, Blickrichtung und Gesichtsmerkmale bestimmen, basierend auf einem einfachen Bild. Sollte ein Video als Quelle fungieren, so kann OpenFace auch lernen. Somit sind die Resultate basierend auf Videos besser als auf einfachen Bildern.

2.4.1 Verarbeitungsschritte

Der Rechenaufwand ist so ausgelegt, dass alle Berechnungen auf einer Webcam in Echtzeit ausgeführt werden können, dies ist im aktuellen Fall nicht notwendig, da es sich um eine nachträgliche Auswertung handelt. Durch den Aufbau sind nur recht kleine Farbbilder der Gesichter in einem Video vorhanden wodurch eine Auswertung erschwert wird.

Gesichts-Landmarks: Detektion und Verfolgung

Für die Bestimmung und Tracking der Landmarks wird ein Conditional Local Neural Fields (CLNF) eingesetzt. Dabei handelt es sich im Grunde um ein Constrained Local Model (CLM) nur mit verbesserter Patch Experts und Optimierungsfunktionen.

Zu Beginn werden verschiedene initiale Hypothesen aus der dlib-Bibliothek verwendet und die Passende ausgewählt. Bei den initiale Hypothesen handelt es sich um verschiedene Gesichtsorientierungen auf denen verschiedene Netze gelernt wurden. Dies ist zwar langsamer, aber auch exakter als eine einfache Hypothese. Wird ein Tracing auf Videos durchgeführt, so wird als initiale Hypothese das Ergebnis aus dem letzten Frame verwendet. Sollte das Tracing scheitern, so wird das CNN reseted um Neu zu beginnen.

Die beiden Hauptkomponenten ist das Point Distribution Model (PDM) zur Erfassung der Anordnung der Landmarks und patch experts zum Erfassen der Variante der einzelnen Landmarks.

Auf diese Weise werden 68 Gesichts-Landmarks und weitere 28 pro Auge erfasst. Zur Brechung auf den Gesichtern sollten sie eine Mindestbreite von 100 Pixeln für eine zuverlässige Detektion Originalgröße besitzt.

Bestimmung der Gesichtsposition

Zur Bestimmung der Translation und Orientierung des Gesichtes wird ein CLNF bzw. PDM eingesetzt. Dabei wurde es mit der Kameraabbildung der 3D-Landmarks eines Kopfes in verschiedenen Positionen initialisiert. Womit auf eine Normierte Abbildung gerechnet wird, diese kann mit den passenden Kameraparameter für die Aufnahme angepasst werden um die reale Position zu bestimmen. Sind keine Parameter bekannt, so können diese anhand der Bildauflösung geschätzt werden.

Bei der Schätzung der Brennweite für ein Bild mit einer Dimension $I_x \times I_y$ wird das Standardobjektiv mit 50 mm und einer Auflösung von 640×480 Pixel angenommen, somit ergibt sich die Brennweiten f_x und f_y wie folgt:

$$f_x = 500 \cdot \frac{I_x}{640}$$

$$f_y = 500 \cdot \frac{I_y}{480}$$

Bestimmung der Blickrichtung

Durch die Landmarks der Augen werden die Augenlider, Iris und Pupille dargestellt und für jedes Auge separat bestimmt. Dabei wird der Augenbereich, basierend auf dem detektierten Gesicht, verwendet, um mit einem weiten CNN die 28 Landmarks des Auges zu bestimmen.

Zur Bestimmung der Blickrichtung wird wie folgt vorgegeben. Zuerst wird der Strahl bestimmt der, ausgehend vom Zentrum der Kamera, durch das Zentrum der Pupille verläuft. Nun wird der Schnittpunkt zwischen diesem Strahl und einer Sphäre bestimmt, die das Auge repräsentiert. Nun wird ein Strahl bestimmt der vom Zentrum der Sphäre ausgehend durch den berechneten Schnittpunkt verläuft, dies ist die resultierende Blickrichtung.

Detection der Gesichtsmerkmale

Dieser Schritt kann von OpenFace ausgeführt werden, ist aber im aktuellen Fall nicht von Relevanz.

2.4.2 Veröffentlichte Genauigkeit

Die Messung wurde auf dem Biwi Kinect head pose und BU Datensatz ausgeführt. Für die Genauigkeit der Kopfposition haben sich folgend Werte ergeben in Grad:

	Yaw	Pitch	Roll	Mean	Median
Biwi Kinect [FGG11]	7.9	5.6	4.5	6.0	2.6
BU dataset [CSA00]	2.8	3.3	2.3	2.8	2.0
ICT-3DHP [BRM12]	3.6	3.6	3.6	3.6	-

Für die Qualität zur Bestimmung der Blickrichtung ergab sich ein durchschnittlichen Fehler von 9.96 Grad.

2.5 ELSE

Um die Blickrichtung möglichst exakt zu bestimmen, sind die Landmarks der Pupille ausschlaggebend. Zu diesem Zwick kann ElSe eingesetzt werden, da dies ein Verfahren zur Detektion von Pupillen in Bildern unter realen Bedingungen ist.

2.5.1 Funktion

Das Verfahren ist in der Lage aus Bildern die Umriss einer Pupille zu ermitteln. Bei realen Aufnahmen sind Bildfehler unvermeidlich, es können Reflexionen (Brille, Kontaktlinse usw.) Make-Up oder körperliche Eigenschaften wie Augenfarbe auftreten um die Detektion erschweren.

Als Ergebnis liefert ElSe eine Ellipse, die den Umriss der Pupille beschreibt.

2.5.2 Funktionsablauf

Als Input wird im Original ein Graubild verwendet, auf dem das Infrarot beleuchtete Auge abgebildet ist. Für den Test im Vergleich zu anderen Verfahren, wurden Bilder von 384×288 Pixel Größe verwendet und ist auf diesen Echtzeit fähig.

Kantendetektion

Da die Pupille als schwarzen Fleck sichtbar ist und die Iris einen helleren Farbton aufweist, wird ein Kantendetektor verwendet, der alle Pixel markiert, bei denen eine starke Farbänderung auftritt. Bei ElSe wird ein Morphologischen Ansatz eingesetzt, von Relevanz sind nur zusammenhängende Kantenpixel, alle anderen können ignoriert werden.

Bestimmen der Ellipse

Um jene Kantenpixel zu erhalten, die die Pupille beschreiben, wird versucht fortlaufende Kanten zu finden, die eine Ellipse bilden. Jene die nicht diesen Anforderung entsprechen können recht schnell ignoriert werden. Anschließend können auch alle offenen Ellipsenverläufe und jene die am meisten vom bestimmten Verlauf abweichen, verworfen werden.

Das beste Ergebnis aller so bestimmten Ellipsen, wird als Lösung verwendet.

Grobe Bestimmung

Sollte die Bestimmung der Ellipse scheitern, so wird das Zentrum des dunkelsten Kreises ermittelt, so ein Punkt kann immer gefunden werden, ist aber nicht zwingend die Pupille.

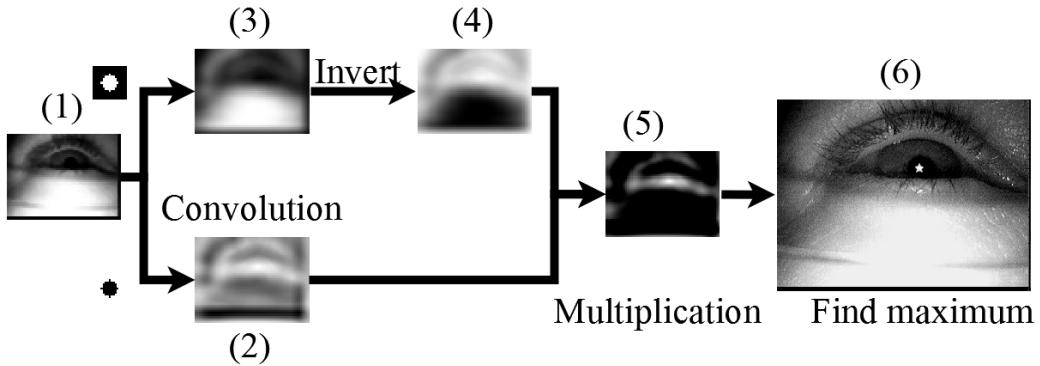


Abbildung 2.7: Ablauf der alternativen Berechnung zur Pupillen-Detektion

Auf einem verkleinerten Bild Abbildung 2.7 (1) wird ein kreisförmiger Mean-Filter eingesetzt mit Ergebnis Abbildung 2.7 (3). Zur zweiten Faltung mit Ergebnis Abbildung 2.7 (2) wird der negative Durchschnitt über ein Quadrat ohne inneren Kreis eingesetzt, wobei der Mean- und negativen Mean-Filter den selbe Radius haben.

Nun wird das Ergebnis des Mean-Filters invertiert Abbildung 2.7 (4) und mittels Punkt-Multiplikation mit dem negativen Meanfilter zusammengebracht Abbildung 2.7 (5). In diesem Bild, wird nun der Höchste Wert gesucht, da dies das Zentrum des dunkelsten kreisförmigen Ortes im Bild ist.

Ergebnis des Bsp ist als Kreuz in Abbildung 2.7 (6) markiert.

2.5.3 Ergebnisse

Im Vergleich zu en anderen Verfahren im Test, zeig sich das ElSe in den meisten Fällen als Sieger hervorgeht mit einer Verbesserung der Erkennungsrate um 14.53%.

Ein Problem entsteht wenn der Farbunterschied zwischen Iris und Pupille recht gering ist, oder durch Reflektionen der Kantenverlauf gestört wird.

To DO: Quellen für Vergleich

2.6 Berechnung der Position

Zur Bestimmung der Position $P = (X_{avg}; Y_{avg}; Z_{avg})$ des Gesichtes im Kamerakoordinaten wird die Größe, ein Skalierungsfaktor S , des Kopfes im Bild verwendet.

Da bei der Abbildung von den Koordinaten ins Bild gilt $x = f \cdot \frac{X}{Z}$ und $y = f \cdot \frac{Y}{Z}$, somit kann die Tiefe wie folgt abgeschätzt werden.

Sei $P_1 = (X_1; Y_1; Z_1)$, $P_2 = (X_2; Y_2; Z_2)$ die Beschreibung der Größe G eines Kopfes mit:

$$\begin{aligned}
a &= \frac{\sqrt{(X_1 - X_2)^2 + (Y_1 + Y_2)^2}}{\frac{Z_1 - Z_2}{2}} = \frac{G}{Z_{avg}} \\
S &= \frac{S_G}{G} \\
\Rightarrow a \cdot f &= f \cdot \frac{G}{Z_{avg}} = S_G \\
Z_{avg} &= \frac{f}{S_G} \cdot G = \frac{f}{S} \\
X_{avg} &= \frac{x \cdot Z_{avg}}{f} \\
Y_{avg} &= \frac{y \cdot Z_{avg}}{f}
\end{aligned}$$

Dies beschreibt allerdings nur eine Annäherung an die Tatsächliche Position, da mit einem Durchschnittlichen Kopfgröße gerechnet wird.

2.6.1 Zusammenhang Bildposition & Weltposition

Als Ausgangspunkt werden die Ergebnisse des CNN eingesetzt um mit deren Hilfe wie in Abschnitt 2.6 beschreiben die Position zu bestimmen. Zur Bestimmung der Orientierung R liefert auch das CNN ein Ergebnis R_{CNN} . Allerdings stimmt es nur im Zentrum des Bildes, da am Rand immer mehr die Orientierung der einzelnen Pixel mit berücksichtigt werden muss.

$$\begin{aligned}
euler_x &= \tan^{-1}\left(\frac{\sqrt{X^2 + Z^2}}{Z^2}\right) \\
euler_y &= \tan^{-1}\left(\frac{\sqrt{Y^2 + Z^2}}{Z^2}\right) \\
R_{pos} &= R(euler_x, euler_y, 0) \\
R &= R_{CNN} \cdot R_{pos}
\end{aligned}$$

Eine weitere Verbesserung kann erreicht werden, indem die gefundenen 2D-Landmarks mit Hilfe des PDM in 3D zu überführen. Um anschließend die Überführung von 2D nach 3D-Koordinaten erneut zu bestimmen um die Orientierung und Position zu ermitteln. Auch bei diesem Verfahren muss die Pixelorientierung beachtet werden.

3 Implementierung

3.1 Ablauf der Implementierung

Zur Bestimmung der Kopfposition und Orientierung wird ein mehrstufiges Verfahren eingesetzt. Am Anfang müssen alle Gesichter, die im aktuellen Frame vorhanden sind, detektiert werden. Dazu wird die MTCNN Face detection verwendet, da dieses Verfahren auch kleinste Gesichter erkennen kann. Abschnitt 3.2

Für die weiteren Berechnungen muss bekannt sein welchen Bereich das Gesicht in Frame einnimmt und um welches es sich handelt. Der Bereich wird vom MTCNN als Box geliefert, als Zuordnung zur Person wird ein Matsching zum vorigen Frame verwendet.

Damit auch eine Berechnung auf den kleineren Gesichtern stattfinden kann, werden alle zu kleinen Bildbereiche hochskaliert. Dabei muss wegen Ungenauigkeiten die gefundene Box etwas räumlich vergrößert und dann auf eine Mindestgröße gebracht werden. Abschnitt 3.3

Diese Bildbereiche werden nun mit OpenFace weiterverarbeitet, um die Position der Landmarks im Bild zu bestimmen. Durch die Berechnung auf der selben Person kann das CNN sich auf jeden einzeln einstellen, um so bessere Ergebnisse zu erreichen. Außerdem könne alle gefundenen Personen gleichzeitig (parallel) ausgewertet werden. Abschnitt 3.4

Bei großen Gesichtern wird nun ElSe auf den Augenbereich angewendet, um die Position der Pupille noch exakter zu ermitteln, damit die Blickrichtung genauer wird. Dazu muss allerdings die Differenz zwischen ElSe-Ergebnis und OpenFace-Ergebnis betrachtet werden um Fehler zu erkennen. Abschnitt 3.5

Nun wird auf Basis der Landmarks und der Kameraparameter die Position und Orientierung des jeweiligen Gesichtes bestimmt und können für weitere Anwendungen verwendet werden. Abschnitt 3.6

3.2 Detektion der Gesichter

Da nur eine einzige fest montierte Kamera ohne Zoom eingesetzt wird, muss sie eine entsprechend hohe Auflösung besitzen damit alle Personen zu erkennen sind. Allerdings machen die eigentlichen Bereiche der Gesichter nur einen sehr geringen Anteil des gesamten aus und diese müssen noch Nachbearbeitet werden. Siehe Abschnitt 3.3

Für die automatische Detektion wird Face-MTCNN Abschnitt 2.2 eingesetzt, da dieses Verfahren die meisten Gesichtern mit verschiedenen Größen im selben Bild findet, sogar recht kleine mit 20×20 Pixeln. Bei diesem Schritt müssen alle Gesichter gefunden werden, auf denen die Berechnung stattfinden soll. Dabei muss das gesamte Gesicht in der Box sein, ansonsten muss es nicht sehr exakt sein, da OpenFace einen eigenen Facedetector besitzt. Wird MTCNN-Face dedector eingesetzt hat sich eine Vergrößerung der Box um 30% als sinnvoll erwiesen, damit sichergestellt wird, dass alle Merkmale wie Nasenspitzen, Kinn, Augenbrauen usw. sicher im Bildausschnitt enthalten sind.

Ebenfalls in diesem Schritt werden die einzelnen Boxen den Personen zugeordnet, damit im späteren Verlauf das korrekte CNN verwendet wird. Für die Zuordnung reicht meist einen einfache Übereinstimmung der aktuellen Box zum vorigen Frame, da die Gesichter sich meist weder groß Bewegen noch sich die Boxen überlappen.

Damit auf allen Gesichter gerechnet werden kann, ist eine Semiautomatische Korrektur erforderlich damit Falsch-Detectionen entfernt und fehlende Boxen ergänzt werden können. Alle nicht gefundenen Gesichtern können manuelle gesetzt oder zwischen dem letzten und nächsten Frame interpoliert werden.

Die gefundenen 5 Landmarks sind für die nachfolgende Berechnung nicht relevant, da sie gerade bei kleinen Gesichtern zu ungenau sind.

3.3 Skalierung auf Mindestgröße

Da OpenFace optimiert ist auf Gesichtern von mindestens 100 Pixel zu arbeiten, werden die Bildbereiche auf diese Größe hochskaliert. Abschnitt 2.3

Die von MTCNN gelieferten und vergrößerten Boxen werden nun auf mindestens 130×180 Pixel gebracht, sollte sie kleiner sein. Neben der einfachen Skalierung, muss die Überführung von Bildkoordinaten des Bildausschnittes in die Koordinaten im Frame bekannt sein, damit dies bei späteren Berechnungen berücksichtigt werden können.

Die Skalierung ist für jeden Bildausschnitt individuell und kann sich durchaus über die Zeit ändern, wenn sich z.B. die Distanz zwischen Person und Kamera verändert.

Von einer zu starken Vergrößerung ist abzuraten, da sich dann der Rechenaufwand pro Gesicht erhöht und die Zuverlässigkeit der Berechnungen von OpenFace wieder sinkt, z.B. durch Falschdetektion des Gesichtes.

3.4 Bestimmung der Landmarks

Für die Bestimmung der Landmarks wird OpenFace eingesetzt. Dabei wird jeder Bildausschnitt unabhängig der anderen Betrachtet und da bekannt ist, um welche Person es sich im Bild handelt, kann direkt mit dem jeweiligen CNN gearbeitet werden, das auf diese Person optimiert wurde.

Durch die vorige Selektion wird nur auf jenen Bildausschnitten gerechnet auf denen auch die Person zu sehen ist, wodurch nicht unnötig gesucht werden muss und auch ein Lernen auf Personen stattfinden kann die nur selten zu sehen sind, da sie nur resettet werden, wenn sie eigentlich zu sehen sein müssten aber nicht detektiert wurden.

Für die eigentliche Bestimmung der Landmarks bietet OpenFace zwei verschiedene Methoden, die Berechnung auf Bildern und Videos. Der Hauptunterschied ist das Lernen, dass bei der Videoauswertung verwendet wird, wodurch sich die Bereiche, auf denen Ergebnisse geliefert werden, deutlich erhöht.

Dies ist interessant für die spätere Anwendung, da somit auch Einzelbilder verwendet werden können, die eine deutlich höhere Auflösung haben als ein Video. Allerdings sinkt dann der maximale Winkel relativ zur Kamera beträchtlich, zu Gunsten der maximalen Distanz. Außerdem können schon kleinsten Farbänderungen im Bild beim Hochskalieren ausschlaggebend sein, ob ein Gesicht erkannt werden kann, wodurch bei gleicher Bildqualität Gesichter im Video besser erkannt werden.

Da die gesamte Berechnung auf Grau-Bildern basiert ist auch eine Farbkorrektur, wie Verbesserung des Kontrast, Farbverlauf usw. möglich, um etwaige Einflüsse bei der Aufnahme zu korrigieren.

Dennoch kann es passieren, dass trotz allem ein Gesicht falsch detektiert wird, wie z.B. das erkennen eines Gesichtes in der Ohrmuschel, diese müssen entsprechend behandelt werden, da ansonsten das Lernen auf diese Bereiche stattfindet und im nächsten Frame erneut nach diesen Merkmalen gesucht wird.

- Verbesserung durch Farbkorrektur

3.5 Verbesserung der Augen

Zusätzlich zu den 64 Landmarks, die ein Gesicht beschreiben, kann von OpenFace weitere 28 Landmarks für ein Auge bestimmt werden, aus denen dann die Blickrichtung ermittelt wird.

Um die Position der Landmarks zu verbessern, kann auf dem Bildausschnitt der Augen der ElSe-Algorithmus eingesetzt werden. Dieser Algorithmus arbeitet auf einem Farbbild um so die Umrisse der Pupille zu berechnen.

Da unter den 28 Landmarks die Umrisse von Pupille und Iris beschrieben wird, müssen diese aus dem Ergebnis von ElSe abgeleitet werden. Dabei hat sich eine Veränderung des Radius mit ?? für Pupille und ?? für die Iris bewährt.

Allerdings muss das Auge für die Berechnung aus entsprechend vielen Pixeln bestehen, wodurch es im Originalbild mindestens mit 10 Pixeln dargestellt wird, um sinnvolle Ergebnisse zu erhalten. Da diese Berechnung unabhängig der Landmarks ausgeführt wird, empfiehlt sich das Ergebnis zu überprüfen, damit die bestimmte Ellipse auch innerhalb der Augenhöhle liegt.

Dabei wird jedes Auge unabhängig vom anderen betrachtet, wodurch sich verschiedene Blickrichtung ergeben. Ab einer Distanz von mehr als ??cm kann die Blickrichtung beider Augen als parallel angesehen und kann entsprechend behandelt werden. Eine Verbesserung ergibt sich, wenn beide Augen anhängig von einander bestimmt werden, damit sich der Fehler minimiert.

3.5.1 To Do

- Größe für Else
- Grenze für Rechnung

3.6 Bestimmung der Position & Orientierung

Für die Bestimmung der Position und Orientierung des Gesichtes wird wie in Abschnitt 2.6 beschrieben ausgeführt. Dies kann Wiederrum von OpenFace übernommen werden, dazu muss nur das Zentrum des Bildes und Brennweite f_x, f_y bekannt sein. Außerdem werden noch erweiterte Verfahren angeboten, bei dem die Position im Bild besser mit einbezogen werden, um die Winkel der Kameraabbildung zu berücksichtigen.

Der signifikanteste Parameter für die Position ist die Brennweite f_x , da mit ihm die Tiefe geschätzt wird und sollte entsprechend exakt bestimmt sein. Von Interesse ist vor allem der Punkt auf den der Blick bzw. das Gesicht ausgerichtet ist, dadurch muss neben der Position im Kamerakoordinatensystem auch die Orientierung bekannt sein.

Da nur die Position des Kopfes und seine Orientierung bestimmt werden kann, ergibt sich das Problem, den konkreten Blickpunkt zu ermitteln, da ein ganzer Kegel, wenn eine Fehlertoleranz berücksichtigt wird, als mögliche Lösungen in Frage kommen.

Außerdem liegt der Blickpunkt meist außerhalb des Bereiches der Kamera und muss entsprechend von einer Anwendung interpretiert werden.

4 Ergebnisse

4.1 Erreichte Werte

4.1.1 Auswirkung der Größe

Durch den Aufbau, muss das Verfahren zuverlässig bezüglich der Größe sein, zur Messung wurde der Datensatz von Labeled Faces in the Wild [HMLLM12] verwendet. In diesem Datensatz ergibt sich im Originalbild eine durchschnittliche Kopfbreite von 94 Pixel.

Zur Durchführung wurden die Größe der Bilder mit dem Faktor multipliziert um so kleinere Gesichter zu erhalten und anschließend mit dem Image-Detector von OpenFace zu detektieren, siehe Abbildung 4.1.

Es ist zu erkennen, dass die Wahrscheinlichkeit auf eine erfolgreiche Detektion ab 0.5, also etwa Gesichert mit 47 Pixel Breite, rapide abnimmt. Bei der verwendeten Kamera Abschnitt 1.4 entspricht dies einer Distanz von etwa 4.5m.

Bei der maximalen Distanz auf der gearbeitet werden soll (8.5m) ergibt sich eine Gesichtsgröße von etwa 22 Pixel, das einer Skalierung von 0.25 entspricht. Bei dieser Bildgröße ist keine Detektion möglich, siehe Abbildung 4.1.

4.1.2 verschiedenen Skalierungsverfahren

Um auf den gewünschten Distanzen arbeiten zu können, wird der jeweilige Bereich Hochskaliert. Dazu wird das Ursprüngliche Bild (250×250) linear um den angegebene Faktor verkleinert und anschließend mit den angegebenen Verfahren auf 300×300 wieder vergrößert. Die Wahrscheinlichkeit auf eine Detektion ist in Abbildung 4.2 abgebildet.

Es ist zu erkennen das durch die Vergrößerung, Gesichter in Bereichen die normal nicht erkennbar sind, bestimmbare werden. Als das ungeeignetste Verfahren hat sich Nearest-Neighbor herausgestellt, siehe blaue Linie Abbildung 4.2. Die anderen haben sehr ähnliche Ergebnisse, nur das Lineare Verfahren ist etwas schlechter. Dennoch werden die Anforderungen, einer Detektion auf Gesichtern von 22 Pixel (Skalierung 0.25) von allen erfüllt.

Ausgehend vom Skalierungsfaktor des Linearen-, Bicubic- und Lanczos-Verfahren wären mit der verwendeten Kamera auch Distanzen bis zu 14m möglich. Allerdings ist das Bild durch die Verkeilung deutlich besser als Originalaufnahmen, da Pixelrauschen nicht vorhanden ist.

4.1.3 Auswirkung von Pixelrauschen

Durch Aufnahme eines Schwarzbildes der Actioncam zeigt sich, dass das Pixelrauschen recht hoch ist, siehe Abbildung 4.3. Das Rauschen hat keine Normalverteilung, sondern es besteht aus kleinen Bereiche, die den selben fehlerhaften Farbwert besitzen.

4.1.4 To Do

- Patch Experts und Optimierungsfunktionen CLM

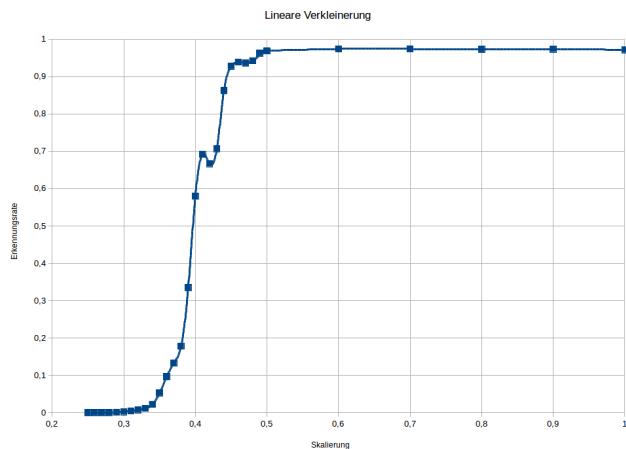


Abbildung 4.1: Die Bilder aus Labeled Faces in the Wild [HMLLM12] wurden mit den Faktor auf der X-Achse linear verkleinert und die Erkennungsrate Y-Achse abgebildet

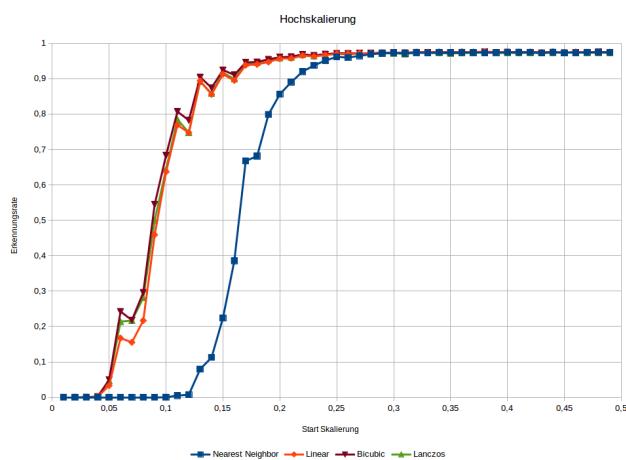


Abbildung 4.2: Die Bilder aus Labeled Faces in the Wild [HMLLM12] wurden mit den Faktor auf der X-Achse linear verkleinert und mit den verschiedenen Verfahren wieder vergrößert Abschnitt 2.3. Aufgetragen gegen die Detektionswahrscheinlichkeit. Nearest-Neighbor (blau), Linear (rot), Bicubic (braun), Lanczos (grün)



Abbildung 4.3: Aufnahme eines Schwarz-Bildes (2688×1520) der Actioncam um den Faktor 7 verstkt und invertiert.

- Auswirkung von Pixelrauschen
 - Rauschen der Actioncam bestimmen
Done
 - Simulation des Rauschens
Add Gaußverteilung auf Image
- Distanzen
- Winkel
- Auswerten der Messung
- Wann ELSE
- Mittlung Ergebnis / Landmarks
- Zuverligkeit mit Farbkorrektur

4.2 Fehleranalyse

Mit entsprechend hochauflenden Kameras knnen auch bessere Resultate auf grerer Distanz erreicht werden. Gerade die Bestimmung der Blickrichtung ist meist nicht mglich, da die Augenpartie viel zu klein fr eine Berechnung ist. So bleibt meist nur die Gesichtsorientierung.

Da Bewegung erlaubt ist passiert es immer wieder, dass Teile des Gesichtes verdeckt werden durch

Hände beim Melde, andere Schüler oder dem Lehrer, der vor der Kamera steht oder sich der Kopf zu weit wegdrehen. Aber auch die Frisuren spielen eine Rolle, da dadurch diese einige Landmarks verdeckt werden und so das Gesicht nicht erkannt wird wie z.B. die Augenbrauen.

Eine Lösungsansatz währen Landmarks in Profilbildern zu detektieren und das verwenden von weiteren Kameras aus anderen Perspektiven.

4.3 Verbesserungen

- Mehrere Kameras für 3D und weniger verdecken und wegdrehen

Literaturverzeichnis

- [App15] Johannes Appel. *Die Bedeutung der Aufgaben für das Beteiligungsverhalten der Schüler : eine Videostudie zur Wirksamkeit des Unterrichtsprozesses*. PhD thesis, 2015.
- [BK08] Gary Bradski and Adrian Kaehler. *Learning OpenCV*. O'Reilly Media Inc., 2008.
- [BRM12] Tadas Baltrušaitis, Peter Robinson, and Louis-Philippe Morency. 3d Constrained Local Model for Rigid and Non-Rigid Facial Tracking. In *Computer Vision and Pattern Recognition (CVPR 2012)*, Providence, RI, June 2012.
- [CSA00] Marco La Cascia, Stan Sclaroff, and Vassilis Athitsos. Fast, reliable head tracking under varying illumination: An approach based on registration of texture-mapped 3d models. *IEEE Trans. Pattern Anal. Mach. Intell.*, 22(4):322–336, 2000.
- [FGG11] Gabriele Fanelli, Juergen Gall, and Luc J. Van Gool. Real time head pose estimation with random regression forests. In *CVPR*, pages 617–624. IEEE Computer Society, 2011.
- [HMLLM12] Gary B. Huang, Marwan Mattar, Honglak Lee, and Erik Learned-Miller. Learning to align from scratch. In *NIPS*, 2012.
- [HR92] Andreas Helmke and Alexander Renkl. Das Muenchener Aufmerksamkeitsinventar (MAI): Ein Instrument zur systematischen Verhaltensbeobachtung der Schueleraufmerksamkeit im Unterricht. *Diagnostica*, 38(2):130–141, 1992.
- [Kyb07] Jan Kybic. Point distribution models, 2007.
- [KZ15] Zhifeng Li Yu Qiao Kaipeng Zhang, Zhanpeng Zhang. Joint face detection and alignment using multi-task cascaded convolutional networks, 2015.
- [MWM08] Louis-Philippe Morency, Jacob Whitehill, and Javier Movellan. Generalized Adaptive View-based Appearance Model: Integrated Framework for Monocular Head Pose Estimation. In *8th International Conference on Automatic Face and Gesture Recognition*, Amsterdam, The Netherlands, 2008.
- [TB16] Louis-Philippe Morency Tadas Baltrušaitis, Peter Robinson. Openface: an open source facial behavior analysis toolkit, 2016.
- [WF16] Thomas Kübler Enkelejda Kasneci Wolfgang Fuhl, Thiago C. Santini. Else: Ellipse selection for robust pupil detection in real-world environments, 2016.
- [Wik15] Wikipedia. Opencv — wikipedia, die freie enzyklopädie, 2015. [Online; Stand 16. Mai 2017].
- [Wik16a] Wikipedia. Bicubic interpolation — wikipedia, the free encyclopedia, 2016. [Online; accessed 6-May-2017].

- [Wik16b] Wikipedia. Lanczos-filter — wikipedia, die freie enzyklopädie, 2016. [Online; Stand 6. Mai 2017].
- [Wik17] Wikipedia. Point distribution model — wikipedia, the free encyclopedia, 2017. [Online; accessed 9-May-2017].