

Mixed Reality Media: Integration of live video feed in 3D environments

Martin Zier

June 30, 2017
Version: Initial Drafting

Beuth University of Applied Sciences

CleanThesis

Department VI: Computer Sciences and Media

Bachelor Thesis

Mixed Reality Media: Integration of live video feed in 3D environments

Martin Zier

- | | |
|--------------------|--|
| <i>1. Reviewer</i> | Kristian Hildebrand
Department VI: Computer Sciences and Media
Beuth University of Applied Sciences |
| <i>2. Reviewer</i> | Prof. Dr.-Ing. René Görlich
Department VI: Computer Sciences and Media
Beuth University of Applied Sciences |
| <i>Supervisors</i> | Kristian Hildebrand and Joachim Quantz |

June 30, 2017

Martin Zier

Mixed Reality Media: Integration of live video feed in 3D environments

Bachelor Thesis, June 30, 2017

Reviewers: Kristian Hildebrand and Prof. Dr.-Ing. René Görlich

Supervisors: Kristian Hildebrand and Joachim Quantz

Beuth University of Applied Sciences

Department VI: Computer Sciences and Media

Luxemburger Straße 10

13353 Berlin

Abstract

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Acknowledgement

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

This is the second paragraph. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Contents

1	Introduction	1
1.1	Overview	1
1.2	Motivation	2
1.3	Problem Statement	2
1.4	Relevance & Challenges	3
1.5	Results	3
1.6	Thesis Structure	3
2	Extending Reality	5
2.1	Motion Video Production	5
2.2	CGI & Video Composition	5
2.2.1	History of Green & Blue Screen Productions	5
2.3	What's VR - Differentiation of AR, VR & MR	5
2.4	Immersion vs Communication	5
2.4.1	Evolution of Virtual Reality Footage	5
2.5	Mixed Reality and its use cases	5
3	System Setup	7
3.1	Hardware Configuration	7
3.1.1	PC Workstation	8
3.1.2	Inogeni 4K2USB3	8
3.1.3	Panasonic GH2 Systemcamera	8
3.1.4	HTC Vive with Controllers and Lighthouses	8
3.1.5	Vive Controller Tripod Mount	9
3.2	Software	9
4	From Video to Mixed Reality	11
4.1	Chroma Key	12
4.1.1	Euclidean RGB Difference	12
4.1.2	Euclidean YCgCo Difference	12
4.1.3	Euclidean Lab Difference	13
4.2	Camera Offsets	18
	Bibliography	19

Introduction

“ *If a technological feat is possible, man will do it.
Almost as if it's wired into the core of our being.*

— **Motoko Kusanagi**
(Ghost in the Shell)

Extending reality with the help of computer generated imagery is no new concept. Ever since real time 3D graphics was possible there was an attempt to extend the understanding of reality. Within the recent years there have been great successes in the industry, most notably in image augmentation was "Pokémon Go" with an estimated install base of 750 million downloads worldwide in June, 2017. [Ann17] Just before this thesis started, Apple and Google showed off their consumer-ready hard- and software for augmented reality experiences.

Virtual Reality Head Mounted Displays have had a similar push in sales with an approximate of 5.83 million sold devices, which range in a sales price between 80 - 900€ for a VR kit, ranging from the very simple Google Daydream View and the very sophisticated HTC Vive. [Erg17] And in these figures are the sales of Google Cardboards missing, which is approximated at around 80 Million.

This generation of computer systems, in which are PC workstations, game consoles and smartphones, is finally sophisticated enough in computation speed and sensor-sensitivity to allow low latency tracking, precise to just a few millimeters.

1.1 Overview

The idea of Virtual Reality (VR) and Head Mounted Displays (HMDs) stems from a cultural need to switch into roles of foreign worlds. Through the advancing development of hard- and software over the last decades emerges a medium which has unmatched immersion and creates an unique, transforming experience into any imaginable environment.

VR and HMDs are now advanced enough for consumer markets - but it stumbles at communicating the experience. Without having ever put on a VR-Headset it is nearly impossible to understand - or even imagine - what the virtual reality experience

means. Any observer of Virtual Reality, usually done by showing what the VR actor is seeing, will not be able to get an understanding of the importance and shift of reality perception without wearing the headset himself.

Showing the video output from a HMD as marketing material is contradicting with classic motion video productions. There is even only one famous example where the perspective of a First Person Shooter is reenacted, which was in the overwhelmingly negatively received Doom (2005) movie.

The VR industry, including but not limited to game developers, exhibition creators and creative studios is in need of better communication of their products that includes more than the current headset wearer and allows for a similar, adapted and immersive experience.

The currently method is called "Mixed Reality" (MR) and uses an external camera with the same tracking hardware of the headset to produce a video signal that shows the real world actor with the environment around him. There are currently three main ways of producing MR footage - where as only one variant allows for live compositing with highly accurate imaging results.

1.2 Motivation

My early teenage years started around the time where digitalization and global interconnectivity begun and broadband Internet became commercially available. Suddenly remote multiplayer games, unlimited image sharing - and yes, music sharing, too -, Java-Applets, Flash, HTML framesets and "Marquee" CSS emerged in that medium. 3D Acceleration became a de-facto standard and even simple office PCs got weak, but dedicated graphics processing units built in. The mass of pixels by increasing the resolution of displays was basically a yearly iteration in greater, better, smaller and brighter.

I am personally very interested and invested in Virtual/Mixed/Augmented Reality to succeed and liked the idea to merge multiple forms of media into one - which is, in my personal opinion, a great summary of my studies and its contents. This thesis represents my interests and the reasons why I chose these studies.

1.3 Problem Statement

Initially I will research motion video productions, computer generated imagery and color theory. This leads to the knowledge to implement basic, interactive live motion video.

The core aspect will be integrating a multitude of Hardware in a software that allows for dynamic video compositing in 3D environments at runtime while a user is interacting with the virtual reality scene. This allows that the person using the Vive HMD to be composited into the scenery and it looks like he is in that scene standing. The essential difference between classic post production is, that this system is planned to operate on runtime, allow additional observers to get an interesting composited imagery of what the VR actor is experiencing.

An additional extension is to dynamically track the camer position, allowing for dynamic camera movement and a freely moving actor.

1.4 Relevance & Challenges

1. Komplette Produktionsworkflow
2. Chroma Keying
3. Bildreproduktion
4. Latencymitigation
5. Tracking

Probably should be called Relevance & Scope.

1.5 Results

Result stuff

1.6 Thesis Structure

This thesis gets contemplated by digital, mostly motion video, material hosted on GitHub. Print is a great medium, but lacks the ability for short demonstrations of video imaging solutions, problems and edge cases. To visualize these problems properly, all video media will have an annotation for cross referencing on the website. It is strongly suggested to follow these links, they will be sorted by chapters.

Extending Reality

” *You are an aperture through which the universe is looking at and exploring itself.*

— Alan W. Watts
(Philosopher)

The well known urban legend of "L'Arrivée d'un train en gare de La Ciotat" in which a train arrives at the La Ciotat station, is, that "the audience was so overwhelmed by the moving image [...] coming directly at them that people screamed and ran to the back of the room". [Wik] With that a new medium was created, which matured into a new art form of film and movies.

This sounds more like prosa text.

2.1 Motion Video Production

2.2 CGI & Video Composition

2.2.1 History of Green & Blue Screen Productions

2.3 What's VR - Differentiation of AR, VR & MR

2.4 Immersion vs Communication

2.4.1 Evolution of Virtual Reality Footage

2.5 Mixed Reality and its use cases

Summarize.

System Setup

The following section describes the hard- and software components used for the thesis and results. All demonstrations have been performed on that environment. All dependencies have been explicitly marked to allow a similar, but not exact, setup to reproduce these results.

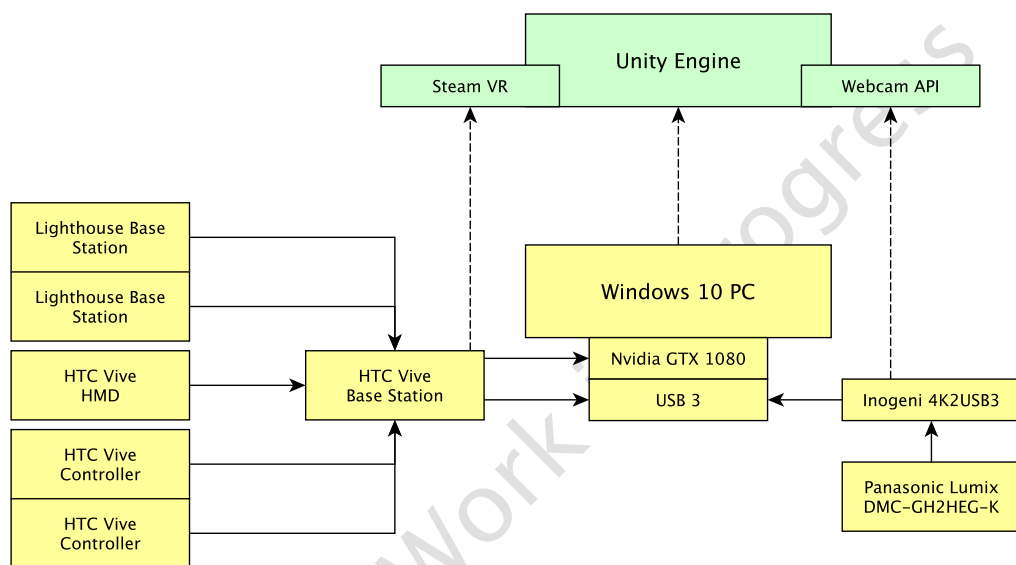


Fig. 3.1: Diagram of hard- and software components.

3.1 Hardware Configuration

The hardware configuration is split in three main parts:

1. Windows PC Workstation
2. Virtual Reality Tracking Solution
3. Motion Video Input Feed

Each individual configuration is basically interchangeable with other systems, as long as predefined conditions are met. Each condition is listed first in each subsection.

3.1.1 PC Workstation

As the software is built in the Unity Engine, the workstation is limited to either Windows or Mac OS X systems the only requirement - besides being powerful enough to render the 3D scenes - is two USB3 ports to ensure enough data throughput for the video and virtual reality solution, as well as two video outputs for a monitor and its headset.

The configuration used here is:

CPU: Intel i7-4700K @ 4.00 GHz
RAM: 16GB DDR4
GPU: Nvidia GTX 1080

This system configuration is to date a high end workstation that has an abundance of render performance, allowing it to process and keep enough framebuffer for the operation described further in.

weak text.

3.1.2 Inogeni 4K2USB3

The Inogeni 4K2USB3 converter is a standalone box that allows to receive any HDMI source and converts it as external webcam video feed. Its advantage is by the arbitrary choice of video cameras and a very simple integration with any software. With the help of the converter box it's possible to request a webcam as video resource and process that video feed as a texture on the GPU.

3.1.3 Panasonic GH2 Systemcamera

This camera provides a direct video feed via HDMI with low latency. It can directly feed into the Inogeni 4K2USB3 and produces a stable, high quality video feed with a low signal to noise ratio in well lit environments.

still unclear if this remains the target camera.

3.1.4 HTC Vive with Controllers and Lighthouses

The current best virtual reality and tracking device is the HTC Vive. It includes two infrared sending stations called "Lighthouse", two Vive Controllers and a Headset,

both systems with 6 degrees of freedom (6DOF) tracking. The tracking system is a blackbox, in which only the transformation matrices for the hand controllers and the HMD can be accessed. By default this transformation has a normalized length of 1 unit to 1 meter. Designing scenery and sense of size is therefore rather easy. The data providing is done by a library called "SteamVR for Unity", which makes the usage in engine transparent.

3.1.5 Vive Controller Tripod Mount

Most cameras have a standardized way of mounting tripods. Since the Vive controllers have no reference plane and minuscule differences in mounting angles changes the projection parameters to noticeable effects, it was necessary to build a mount for the camera to keep controller and video equipment transformation in sync, I built a mount that fits on tripod attachment points and keeps the controller locked in the same position.

add example for incorrect projection and model of mount

3.2 Software

The software of choice is Unity3D, which is free for students, non-profit organizations and small studios. It provides a

interrupted writing

From Video to Mixed Reality

Needs better sourcing.

To achieve a real time rendering environment, as previously mentioned, there are two main production cycles. The one discussed in this thesis resolves this problem by staying inside on application with multiple render cycles per frame.

The first and foremost render cycle is the stereoscopic output of the Vive HMD, which has a set framerate of either 45 or 90 frames per second. It is important to have a consistent performance, otherwise the experience for the actor with the HMD will have a terrible experience.

The secondary render cycle has to be done on the same frame, which is a virtual camera inside the virtual scene and the relative position of the real world HMD and real world camera. Since the SteamVR library for the HTC Vive already exposes a normalized, synchronized tracking, it is easily possible to position the virtual camera at an accurate location.

The following chapter describes the techniques used to transform motion video inside a greenscreen into a mixed reality image. As brief overview, the steps required are performed in referential order from the motion video from the camera feed. This is different to the render order but gives a better understanding of the techniques used to achieve mixed reality imagery.

4.1 Chroma Key

Beginning from the camera, the video signal travels through the Inogeni converter and is accessible with the system API for webcams. (See figure 3.1)

The initial step is to remove the green background from the image, which should be greenscreen. For a reference green, there has to be a color picked manually in the material editor of Unity - this was made easy by a checkbox to show raw output from the camera. Then a middle-ground green can be picked. This is an important setup step, since lightning situations can vary greatly and minor differences in light setups can have a great effect on the outcome of visible green background captured by the camera.

4.1.1 Euclidean RGB Difference

Assuming a source pixel color C_S and a reference color C_R we can calculate the euclidean distance between these colors.

$$distance = \sqrt{(C_{R,R} - C_{S,R})^2 + (C_{R,G} - C_{S,G})^2 + (C_{R,B} - C_{S,B})^2} \quad (4.1)$$

This is a computationally very low cost and works well enough for tell a difference between two separate colors. It fails to accommodate for colors that are perceived as different, but are tinted by the reference colors. Since the greenscreen will never achieve 0% reflectivity, some residue of the background color will mix with the filmed actors.

An extreme example case is used for comparing these chroma keying variants:

4.1.2 Euclidean YCgCo Difference

YCgCo stands for Luminance (Y), chrominance green (Cg) and chrominance orange (Co) and helps decorrelating color spaces. Since it is a fast, lossless color transformation it is used in example for H.264 video encoding. The two chrominance channels are then split into green to magenta and orange to blue color values and allow for a more accurate distance calculation between two colors.

Transforming any arbitrary RGB color to YCgCo can done with a single matrix multiplication:



Fig. 4.1: Comparison Image - sRGB Output

$$\begin{bmatrix} Y \\ Cg \\ Co \end{bmatrix} = \begin{bmatrix} \frac{1}{4} & \frac{1}{2} & \frac{1}{4} \\ -\frac{1}{4} & \frac{1}{2} & -\frac{1}{4} \\ \frac{1}{2} & 0 & -\frac{1}{2} \end{bmatrix} * \begin{bmatrix} R \\ G \\ B \end{bmatrix} \quad (4.2)$$

Given two colors, one from the video source C_S and a reference color C_R it is now possible to calculate the euclidean distance on the two chrominance channels:

$$distance = \sqrt{(C_R Cg - C_S Cg)^2 + (C_R Co - C_S Co)^2} \quad (4.3)$$

Since the increased decorrelation, the result is more accurate and shows less artifacting on unwanted pixels.

4.1.3 Euclidean Lab Difference

The International Color Consortium (ICC) defined 1976 *Lab* ΔE as a standard way of calculating color differences with *Lab* colors. The final distance calculation is the linear euclidean distance as with all other models, but accommodates for perceived color differences.



Fig. 4.2: Chroma Keying by using euclidean RGB distance

this is very rough and only contains equations currently used

sRGB conversion to linear RGB in respect of energy per channel:

$$v \in \{r, g, b\} \wedge V \in \{R, G, B\} \quad (4.4)$$

where:

$$v = \begin{cases} V/12.92 & \text{if } V \leq 0.0405 \\ ((V + 0.055)/1.055)^{2.4} & \text{otherwise} \end{cases} \quad (4.5)$$

from there 1

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = [M] \begin{bmatrix} R \\ G \\ B \end{bmatrix} \quad (4.6)$$



Fig. 4.3: Chroma Keying by using euclidean RGB distance

where:

$$[M] = \begin{bmatrix} RX_r & GX_g & BX_b \\ RY_r & GY_g & BY_b \\ RZ_r & GZ_g & BZ_b \end{bmatrix} \quad (4.7)$$

$$\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} [M] = \begin{pmatrix} X_r/Y_r & X_g/Y_g & X_b/Y_b \\ 1 & 1 & 1 \\ \frac{1-X_r-Y_r}{Y_r} & \frac{1-X_g-Y_g}{Y_g} & \frac{1-X_b-Y_b}{Y_b} \end{pmatrix} \quad (4.8)$$

Where $[M]$ for RGB D65 is:

$$\begin{bmatrix} 0.4124564 & 0.3575761 & 0.1804375 \\ 0.2126729 & 0.7151522 & 0.0721750 \\ 0.0193339 & 0.1191920 & 0.9503041 \end{bmatrix} \quad (4.9)$$

Based on a reference white $U_r \in \{X_r, Y_r, Z_r\}$:

$$U \in \{X, Y, Z\} \wedge W \in \{L, a, b\} \quad (4.10)$$

$$\epsilon = 0.008856 \wedge \kappa = 903.3 \quad (4.11)$$

where:

$$w_r = \frac{U}{U_r} \quad (4.12)$$

$$f(w) = \begin{cases} \sqrt[3]{w_r} & \text{if } U > \epsilon \\ \frac{\kappa w_r + 16}{116} & \text{otherwise} \end{cases} \quad (4.13)$$

$$\begin{bmatrix} L \\ a \\ b \end{bmatrix} = \begin{bmatrix} 116f_y - 16 \\ 500(f_x - f_y) \\ 200(f_y - f_z) \end{bmatrix} \quad (4.14)$$

With this conversion from sRGB to linear RGB to XYZ to Lab we can now calculate the euclidian linear distance between two colors C_1 and C_2 , which already have been converted to Lab:

$$\Delta E = \sqrt{(C_2L - C_1L)^2 + (C_2a - C_1a)^2 + (C_2b - C_1b)^2} \quad (4.15)$$

These values are rated by their perceptive difference [MW]:

0.0 ... 0.5	the difference is unnoticeable
0.5 ... 1.0	the difference is only noticed by an experienced observer
1.0 ... 2.0	the difference is also noticed by an unexperienced observer
2.0 ... 4.0	the difference is clearly noticeable
4.0 ... 5.0	fundamental color difference
> 5.0	gives the impression that these are two different colors



Fig. 4.4: Chroma Keying by using euclidean RGB distance

4.2 Camera Offsets

```
> Is this text?  
< No, this is doge.
```

Bibliography

- [Ann17] App Annie. *App Annie 2016 Retrospective*. Retrospective Report. 2017, pp. 2, 25 (cit. on p. 1).
- [MW] Tatol M. Mokrzycki W.S. *Colour difference ΔE - A survey*. Survey. Faculty of Mathematics, Informatics, University of Warmia, and Mazury, p. 20 (cit. on p. 16).

Websites

- [Erg17] Deniz Ergürel. *The latest virtual reality headset sales numbers we know so far. As of March 2017*. 2017. URL: <https://haptic.al/latest-virtual-reality-headset-sales-so-far-9553e42f60b5> (cit. on p. 1).
- [Wik] *L'Arrivée d'un train en gare de La Ciotat*. URL: https://en.wikipedia.org/wiki/L%27Arriv%C3%A9_d%27un_train_en_gare_de_La_Ciotat (cit. on p. 5).

List of Figures

3.1	Diagram of hard- and software components.	7
4.1	Comparison Image - sRGB Output	13
4.2	Chroma Keying by using euclidean RGB distance	14
4.3	Chroma Keying by using euclidean RGB distance	15
4.4	Chroma Keying by using euclidean RGB distance	17

List of Tables

Colophon

This thesis was typeset with \LaTeX 2_ε. It uses the *Clean Thesis* style developed by Ricardo Langner. The design of the *Clean Thesis* style is inspired by user guide documents from Apple Inc.

Download the *Clean Thesis* style at <http://cleanthesis.der-ric.de/>.

Declaration

You can put your declaration here, to declare that you have completed your work solely and only with the help of the references you mentioned.

Berlin, June 30, 2017

Martin Zier

