Computer Aided Medical Procedures
Prof. Dr. Nassir Navab

Dissertation

# XR For All: Closed-loop Visual Stimulation Techniques for Human and Non-Human Animals

Hemal Naik

Fakultät für Informatik
Technische Universität München

# Technische Universität München

Fakultät für Informatik

Lehrstuhl für Informatikanwendungen in der Medizin

# XR For All: Closed-loop Visual Stimulation Techniques for Human and Non-Human Animals

## Hemal Naik

Vollständiger Abdruck der von der Fakultät für Informatik der Technischen Universität München zur Erlangung des akademischen Grades eines

Doktors der Naturwissenschaften (Dr. rer. nat.)

genehmigten Dissertation.

| | |
|---|---|
| *Vorsitzende:* | Prof. Dr. Anne Brüggemann-Klein |
| *Prüfer der Dissertation:* | 1. Prof. Nassir Navab, Ph.D. |
| | 2. Prof. Dr. Iain Couzin |
| | 3. Prof. Dr. Oliver Deussen |

Die Dissertation wurde am 25.08.2020 bei der Technischen Universität München eingereicht und durch die Fakultät für Informatik am 02.02.2021 angenommen.

# Abstract

Augmented Reality (AR) and Virtual Reality (VR) are truly exciting technologies for artificial sensory stimulation. Especially because they can provide immersive and interactive experiences (closed-loop). We have observed rapid progress in the development of these technologies in the last three decades, especially in the domain of visual stimulation. AR/VR solutions have been proposed for a wide range of tasks, from military training to health care and from industrial manufacturing to animal behavior studies. Although such applications are proven to be useful and cost-effective, many do not make it beyond the scope of research labs. Lack of attention towards workflow integration is considered one of the major causes behind lower acceptance of AR/VR applications. This is especially true when users intend to use the technology to create a tool for a specific purpose, e.g. a visualization or a measurement device. In such cases, it is extremely important to understand the application domain and to work with the end-users for developing useful AR/VR solutions. In this thesis, we will discuss the importance of developing workflow specific solutions for the application of AR/VR.

Our research focuses on two very different application domains of AR/VR: industrial manufacturing and animal behavior studies. We show that closed-loop visual stimulation techniques of AR and VR are used in these domains as interactive visualization tools. However, the users of both domains have very specific requirements and expect the tools to be compatible with existing workflows. Our research shows that working closely with the users, understanding their practices, and the environment is essential to design workflow specific solutions.

Industrial users intend to use AR visualization tools for various industrial processes without losing reliability and accuracy provided by existing techniques. We developed Spatial Augmented Reality(SAR) solutions for two industrial processes: part alignment and marking. We worked closely with industrial partners to understand the working principle of existing methods and proposed an alternative AR solution while maintaining the same principles. The methods are developed with the intention of replacing the existing methods to save production costs and time.

Similarly, biologists intend to build novel AR/VR visualizations tools with live animals to study animal behavior. Biologists have been using Virtual Environments to study behavior but the existing tools have many limitations e.g. limited tracking methods, small setups, etc. We identified that technology experts and biologists have to work together to overcome these limitations. In this spirit, we prepared a literature review to introduce technology experts to the field of sensory stimulation based behavior studies. Taking a step further, we worked closely with biologists to develop a novel experimental setup suitable for conducting novel AR/VR experiments with a group of animals.

# Zusammenfassung

Augmented Reality (AR) und Virtuelle Realität (VR) sind sehr aufregende Technologien zur künstlichen Stimulation der Sinne. Immersive und interaktive Erfahrungen (Closed-Loop) ermöglichen können. Wir haben in den letzten drei Jahrzehnten rasche Fortschritte in der Entwicklung dieser Technologien beobachtet, insbesondere im Bereich der visuellen Stimulation. AR/VR-Lösungen wurden für ein breites Spektrum von Aufgaben vorgeschlagen: von der militärischen Ausbildung bis zur Gesundheitsversorgung und von der industriellen Fertigung bis zu Verhaltensstudien an Tieren. Obwohl sich diese Anwendungen als nützlich und kosteneffektiv erwiesen haben, verlassen viele nicht die Forschungslabore. Die mangelnde Aufmerksamkeit in der Integration von Arbeitsabläufen gilt als eine der Hauptursachen für die geringere Akzeptanz von AR/VR-Anwendungen. Dies gilt insbesondere dann, wenn Benutzer beabsichtigen, die Anwendung als Werkzeug für einen bestimmten Zweck, z.B. eine Visualisierung oder ein Messgerät, einzusetzen. In solchen Fällen ist es äusserst wichtig, das Anwendungsgebiet zu verstehen und mit den Endnutzern zusammenzuarbeiten, um nützliche AR/VR-Lösungen zu entwickeln.In dieser Doktorarbeit werden wir die Bedeutung der Entwicklung von arbeitsablauf-spezifischen Lösungen für die Anwendung von AR/VR diskutieren.

Unsere Forschung konzentriert sich auf zwei sehr unterschiedliche Anwendungsbereiche der AR/VR: industrielle Fertigung und Tierverhaltensstudien. Wir zeigen, dass Closed-Loop-Techniken der visuellen Stimulation von AR und VR für diese Bereiche zum Zweck der interaktiven Visualisierung nützlich sind. Die Benutzer beider Bereiche haben jedoch sehr spezifische Anforderungen. Sie wollen AR/VR-Werkzeuge als Unterstützung ihrer Aufgaben einsetzen und erwarten, dass die Werkzeuge mit bestehenden Arbeitsabläufen kompatibel sind. Wir zeigen, dass es möglich ist, massgeschneiderte Lösungen zu entwerfen, indem wir eng mit den Nutzen zusammenarbeiten, und um ihre Praktiken und die Umgebung zu verstehen.

Industrielle Nutzer möchten interaktive Visualisierungstools für industrielle Prozesse nutzen, ohne die Zuverlässigkeit und Genauigkeit der vorhandenen Techniken zu verlieren. Wir haben eng mit Industriepartnern zusammengearbeitet und alternative Loesungen mit AR-Techniken vorgeschlagen, wobei das Arbeitsprinzip der bestehenden Methoden beibehalten wurde. Wir haben auf Spatial Augmented Reality basierende Loesungen fuer den Prozess der Ausrichtung und Markierung entwickelt. Unsere Methoden sind so konzipiert, dass sie die bestehenden Methoden ersetzen und gleichzeitig haben wir deren Arbeitsprinzip replizieren, um den Anforderungen der Endnutzer gerecht zu werden.

Auf ähnliche Weise wollen Biologen interaktive Visualisierungstools bei lebenden Tieren einsetzen, um neue experimentelle Techniken zur Untersuchung von Tierverhalten zu entwerfen. Biologen verwenden bereits VR-Setups, die auf ihre Bedürfnisse zugeschnitten sind. Jedoch

haben die bestehenden Techniken viele Einschränkungen, z.B. begrenzte Trackingmethoden, kleine Setups, usw. Wir haben festgestellt, dass eine enge Zusammenarbeit zwischen Technologieexperten und Biologen erforderlich ist, um viele dieses Einschränkungen zu beheben. In diesem Sinne haben wir eine Literaturübersicht für die Technologieexperten erstellt, um ihren Beitrag zum Thema zu unterstützen. Wir entwickelten ausserdem neue Methoden zur Verfolgung mehrerer Tiere in großen Gebieten. Darüber hinaus entwickelten wir verschiedene Designkonzepte unter Verwendung unserer Tracking-Lösungen wie z.B. als Anwendungsfaelle für die Motivation der nächsten Generation von Verhaltensexperimenten mit XR-Technologien.

# Acknowledgments

It has been a long journey and I do not know how I have made it to the other side but I surely know that without the support of some really cool people it would not have been possible. First, I would like to start with my mentor and supervisor Prof. Nassir Navab. Nassir has played a great role in shaping my ideas about science and life. I will remember our discussions very fondly, especially the ones on the bus while going for the annual retreats. He always advised me to do what I wanted to do and for that, I can not thank him enough. The second person I would also like to thank is Prof. Iain Couzin. Iain introduced me to the fascinating world of animal behavior studies. He gave me absolute freedom to pursue my own directions. He has been extremely friendly and always offered room for discussing professional as well as personal life. I will always be thankful to him for giving me *the wild card* entry into the wildlife studies. I would also like to thank Prof. Oliver Deussen for joining my thesis committee.

I have been fortunate enough to spend time at different institutions during my Ph.D. First, I would like to thank all my colleagues and friends at the CAMP chair. I can not forget to mention the team from Extend3D, Peter, Nick, Bjoern, Andres, Chris, Theresa, Beatriz, and Andy the "laser man". I spent 5 years with Extend and learned a lot with them. Eventually, I also learned that industrial manufacturing was not the path for me. I took a detour and moved to the Max Planck Institute of Animal Behavior. My friends and colleagues generously took time to teach me a thing or two about life on the planet. A special thanks to Vivek, Martina, Stefano, Anja, Camilla, Heiko and Jacob for supporting me during #Coronatime thesis writing. I could not find a better place to mention Simon and Renaud, both of them helped me through many tough days, *Maximum Bamboule* for life. I learned many things through them and I am grateful for all those evenings discussing stories from different parts of the world.

I must mention three individuals who played a very special role in shaping my ideas. Michael Emenaker introduced me to the fascinating world of birds in the foothills of Himalayas. I still remember those days very fondly, undoubtedly they have changed my life for good. Mate Nagy, an exceptional friend, and the lighthouse to my stranded boat in the rough sea of the Ph.D. I know that he is always there for me with the right kind of advice and I am truly grateful that he treats me not as a colleague but as a younger brother. I have great respect and admiration for Prof. Martin Wikelski, from him I learned the importance of promoting scientific activities in wildlife at a global scale. He took interest in my ideas of working with wildlife researchers in India. I am thankful to him for providing support for visiting various research institutions in Europe and India for promoting research collaborations and science outreach activities.

# Contents

## IV    Conclusion            155

## 7   Conclusion            157

## V    Appendix            159

## A   Survey on requirements of collective behavior studies        161

## B   Tracking errors and open problems        167

## C   List of Authored and Co-authored Publications        171

## D   Abstracts of Publications not Discussed in this Thesis        173

## Bibliography        177

## List of Figures        191

## List of Tables        193

# Part I

Introduction and Background

# Introduction

> *Man is a tool-making animal*

— **Benjamin Franklin**

The art of tool making is considered to be one of the greatest achievements of humans as a species. Our journey started from primitive tools made of stones and now we build complex tools such as satellites and smartphones. Tools primarily allow us to accomplish certain tasks or to build products which make our life easier. Modern-day humans have reached the digital age in which we are creating breakthrough innovations with use of computers. Some of these novel ideas and technologies enable us to investigate the laws of nature and gain a deeper understanding of our environment.

One of the primary topics in this direction is related to our understanding of the world through our sensory experiences. Over the years we have built many tools to study humans and other animals. We learned that each animal can only sense a limited amount of information from its environment. Humans can not hear ultrasonic sounds like bats or see an ultraviolet spectrum of light like insects. Understanding of sensory organs lead to many questions such as: can we build tools to experience "non-human" sensory experiences, e.g. UV vision? Can we enhance our sensory experience, e.g. visualize microscopic objects? Or can we create novel experiences detached from our immediate environment, e.g. deep-sea diving or floating in space? These questions have captured the imagination of many. This curiosity to extend the limits of our sensory experience vision possibly contributed to inventions such as a telescope, the microscope for visual enhancement and speakers, microphones for auditory enhancement. Now in the digital age, we are trying to create interactive and immersive sensory experiences using technologies like Augmented Reality (AR) and Virtual Reality (VR).

In the past few decades, AR/VR researchers have designed tools and methods that add virtual content to the real world or replace the real world and "teleport" humans into a virtual world. AR/VR tools are considered useful for both curiosity-driven (fundamental) research and application-driven research. Fundamental research for investigating sensory-motor mechanisms involve question-related to functional aspects of the sensory organs. For example, how do the sensory systems (e.g. vision, auditory) function? And how can we design artificial stimulations to create desired motor responses? This field of research is not only limited to humans but also pursued with other animals [149]. The fundamental knowledge gained by understanding the sensory systems inherently drives the application-driven research which focuses on creating synthetic sensory experiences e.g. television or radio. Application-driven research specific to AR/VR broadly focuses on two aspects. The first aspect is to create novel tools and experimental setups to support fundamental research and to validate its findings. The second aspect is to create a useful set of applications for improving human sensory experiences for tasks in various fields such as education, entertainment,

manufacturing, etc. Our research presented in this thesis will deal with both of these aspects of application-oriented AR/VR research.

## 1.1 Motivation

The motivation for promoting AR/VR research and developing AR/VR tools are broadly defined in the previous section with their suitability for curiosity and application-driven research. These goals are broad and often not in direct perspective while researchers are working on problems in AR/VR. This is because designing virtual sensory experiences is a complex problem. The state-of-the-art AR and VR systems are the product of years of research in various fields such as vision, perception, psychology, computer vision, software engineering, display technology, computational algorithms, etc [149]. Many technological components come together and function seamlessly to design and produce artificial sensory stimulation. Designing each component presents its challenges and research questions. Often the intention to find a mathematically elegant solution or aesthetically pleasing solution drifts the inventors far from the overall objective. As a result, we may end up creating solutions for problems that do not exist or do not fit into the scope of the problem.

It is extremely important to keep an eye on the final objective because a tool is ultimately developed for a specific user who would use it in a specific scenario e.g. a teacher in a classroom or a worker in a factory. Therefore, the tool designer must have an insight into the user's working methods and the working environment to understand all aspects of the problem. It is ideal to claim that the objective must always solve the problems with the simplest possible solution that meets the requirements of the user without adding unnecessary or complex operational hurdles. AR/VR tools and applications are often developed without close cooperation between end-users and the researchers. As a consequence, many tools are rejected by their intended end-users and they are not used to solve *real world* problems.

The central argument made in this thesis is that AR/VR developers have to gain a complete understanding of the user environment and practices to create a solution that the user can appreciate. This argument is true especially where strict working practices are already in place e.g. industrial manufacturing [70, 153]. A new practice or a tool is introduced after evaluating its performance for multiple factors including its compatibility with existing workflows. Other factors may include accuracy, efficiency, cost, time or effort, etc. The solutions or practice that do not meet these criteria are rejected because it impedes the existing process.

Similar rules apply to the fields where novel AR/VR approaches are required for conducting fundamental research. For example, virtual environments are used with insects, mammals, and fish for studying animal behavior (ethology) [149, 195, 205]. The technical development is done by biologists rather than computer scientists. Biologists have designed customized setups for working with animals using existing AR/VR methods and tools designed for humans. We argue that the direct involvement of technology experts is required for improving the existing technology and to make it more suitable for animal behavior experiments. These developments are not possible without working with biologists to understand their problems.

Taking these examples as main motivation this thesis covers work done by the author in both of these fields, industrial manufacturing, and animal behavior studies, to emphasize the need for inclusive development. Keeping in mind the diverse nature of applications presented in this thesis, the author has written this thesis with the hope that readers from both industrial and biology backgrounds will appreciate our understanding of their problems and will be able to replicate the methods. Additionally, the readers from computer science (developers of AR/VR) may benefit from the insights presented in this thesis.

## 1.2 Structure

The thesis is divided into 5 parts: **Part 1** contains an introduction to the topic of sensory stimulation and brief background of sensory stimulation techniques. After that, we will introduce virtual stimulation techniques using AR and VR. This chapter includes details of technical concepts, mainly visualization and tracking, that are required to design AR and VR tools. Both, AR and VR, have thriving research communities and many active directions of research application. Therefore, we have limited our discussions to material that is relevant for understanding and replicating our research.

In **part 2**, we start with a discussion on the desired properties of industrial AR (IAR) solutions that are outlined in Chapter 3. We will present two processes in industrial manufacturing i.e. part alignment and marking process and mentioned existing practices. Further, we propose the idea of designing new solutions for these processes using Spatial Augmented Reality (SAR) and discuss relevant challenges. Chapter 4 includes the implementation details of the IAR solution developed by the author and industrial collaborators. The chapter will conclude with closing remarks related to this work and the author's learning from this project. This research is done in collaboration with Extend3D GmbH and other industrial partners. The funding for this project was supported by the German Federal Ministry of Education and Research under the ARVIDA initiative (grant no. 01IM13001N).

In **part 3**, we will discuss applications of virtual environments as interactive sensory stimulation tools in fundamental research, more specifically for studying animal behavior. These applications are currently underrepresented in AR/VR research community and therefore we dedicate chapter 5 to provide a detailed understanding of the use of artificial stimulation in fundamental research and the impact of such studies in other fields of research such as robotics, AI, medicine. We will conclude the chapter by highlighting the need for developing novel AR/VR solutions for studying animal behavior. In chapter 6, we propose a novel experimental setup that allows setting up interactive sensory stimulation experiments with different types of animals i.e. humans, dogs, birds, or insects. We will present this idea as "XR for all" and outline some novel use cases that are not yet explored in behavior studies. The author has focused on developing real-time tracking solution (marker and marker-less) for tracking the head movement of birds. Finally, the chapter is concluded with two example applications with the head tracking approach. The research proposed in this part is done in collaboration with an interdisciplinary team of scientists at the collective behavior at the Max Planck Institute of Animal Behavior. The research presented is supported by the DFG center of excellence 2117 Center of the Advanced Study of Collective Behavior (ID: 422037984).

In **part 4**, we will conclude both the projects in chapter 7. The author provides highlights of the importance of working with the domain experts for the development of AR/VR solutions.

The research presented in this thesis is conducted with the help of many collaborators from both industry and academic institutions. The author has used we-form throughout the thesis to give credit to all of the collaborators.

# Artificial Sensory Stimulation with AR and VR

<div style="text-align:right">2</div>

In this chapter, we begin with a short background on the methods and tools developed to enhance or alter sensory experiences. We will move on to digital methods for sensory stimulation (virtual environments) and explain the concept of real and virtual world continuum. At this point, we will define the scope of our research in the context of the virtual continuum. Visual stimulation is central topic of this thesis and therefore we have included discussion on vision properties.

After introducing virtual environments we will focus specifically on Augmented Reality (AR) in section 2.2 and Virtual Reality (VR) in section 2.3. AR and VR have a long history of research [39, 137], therefore we will limit the discussions to a brief introduction with some example applications. The presented examples will provide a general overview of different visualization schemes used in AR and VR applications. Finally, we will explain core technological concepts for designing of AR/VR tools in section 2.4. This section will include technical considerations required for designing the AR/VR tools. Overall this chapter will provide a reasonable background for research presented in this thesis.



**Fig. 2.1.** An Illustration of The Allegory of the Cave, from Plato's Republic [171]. The concept explains that perception of the real world for the prisoners is manufactured by shadows on the wall. The concept was originally used for making philosophical arguments on "the effect of education and lack of it on our nature". This concept is valid even today because we know that artificial sensory stimulation presented in the right manner can indeed create illusions of an alternative reality.
*Author:* 4edges, *Source:* Wikimedia commons, *License:* CC-BY

## 2.1 Background

Humans have designed many tools for the artificial stimulation of sensory organs with the objective of modifying the sensory information available to us from our environment. Such ideas and tools for manipulation of senses are invented by entertainers, artists, thinkers, and scientists for various purposes. Some ideas require the use of a particular tool or an apparatus, such as a lens or a mirror. Lenses are widely used to make simple tools for modifying visual information or creating an optical illusion. For example, lenses are used with wearable glasses to alter the appearances of real-world objects in terms of distance, shape, or color. Similarly, lenses are also used for making scientific apparatus (e.g. microscopes and telescopes) that allow us to enhance our visual abilities by several orders of magnitude in terms of scale and distances. Some ideas of sensory stimulation rather involve modification of the environment itself to create an alternative form of reality, e.g. theatrical performances. Theatrical events are generally performed in specially designed arenas, such as an opera hall or a black box theater, to enhance the auditory and visual experiences while detaching the viewer from the outside world. Viewers are subject to intense experiences with the help of carefully designed manipulation of lights and sound using techniques like shadowing, echo, music, etc. Entertainers consider sensory manipulation a vital part of the storytelling as it allows the viewer to become part of the story during the show. Artists working with static forms of art (sculptures and paintings) often chose to exploit the limitations of human sensory organs to design the illusion of motion or depth. Some of the most exciting examples of such art form can be seen in works of Rembrandt or Cezanne and M.C. Escher [65].

Scientists design sensory stimulation based experiments to understand the functional capabilities and limitations of human sensory organs [149]. The knowledge of the functional properties of human sensory organs is used in various fields of sciences e.g. biology, psychology, engineering, etc. For example, the visual illusion displayed in figure 2.2 is used to study depth perception and an illusion displayed figure 2.3 is used for understanding motion perception. The mechanism of sensory organs is replicated to design tools that mimic sensory organs. For example, recorders are designed to record sounds and speakers are designed reproduce sounds for auditory stimulation. Such tools become vital in designing the next generation of scientific experiments that depend on artificial means of sensory stimulation. For example, biologists conduct sensory stimulation based experiments with animals having completely different sensory modalities than humans (details in part 3). Artificial stimulation based animals behavior studies are useful for fundamental research (e.g. vision, neurology and medicine) and applied research for technology development (e.g. bio-inspired robotics, artificial intelligence) [7, 97, 149, 202]. Such experiments not only provide a deeper understanding of our environment but also provide means to design sensory experiences beyond the limitations of our sensory organs e.g. recorders that sense ultrasound waves.

Over the years, extensive research in sensory modalities has led to the invention of digital technologies which in turn have reshaped the idea of sensory stimulation. The invention of faster computers, high-resolution display technologies, and high-speed cameras have triggered a new wave for artificial stimulation based experiments. For example, display technology has improved dramatically in the last few decades, from CRT screen to LCD screens and more recently 4K screens, to produce realistic images. It is now possible to design sensory stimulation

**Fig. 2.2.** Fraser's Spiral is an optical illusion described by British psychologist Sir James Fraser in 1908. The picture contains a series of concentric circles which give false illusion of a spiral. This illusion is created by combining lines with misaligned color patterns [73].
*Author:* Mysid, *Source:* Wikimedia commons, *License:* CC-BY



**Fig. 2.3.** The image creates an illusion of motion even though the image is static. The circles appear to be moving due to the retinal effects of opposing edge contrasts and size variations.
*Author:* Fiestoforo, *Source:* Wikimedia commons, *License:* CC-3.0

techniques in an interactive (closed-loop) manner such that stimulation is context-dependent, personalized, and react to the behavior of the observer e.g. augmented reality and virtual reality. Such technologies have given way to a new set of tools and applications using dynamic sensory stimulation. For example, AR equipped microscopes may have virtually embedded text or color schemes that make scientific studies easier. Similarly, artificial stimulation techniques are getting closer to real-world simulations. The idea behind virtual reality applications is to design artificial stimulations that are indistinguishable from real-world stimulations. This also includes designing completely artificial experiences that may seem unnatural and yet believable e.g. teleportation. The research presented in this thesis focuses on the topic of interactive sensory stimulation of humans and other animals.

Arguably *true* virtual experiences can not be provided without stimulation of multiple senses in addition to vision i.e. auditory, tactile. Multiple attempts have been made in the direction of multi-sensory stimulation in humans [132] and other animals [83, 205]. However, it is still not possible to design rich multi-sensory experiences due to a lack of complete understanding of sensory-motor mechanisms and a lack of a technological tool to support such stimulations. As a result, most applications of artificial stimulation are focused on specific sensory organs. The research discussed this thesis will focus on the topic of visual stimulation. The readers should assume vision as default sensory modality for any text referring to generic terms such as virtual simulation or artificial stimulation. In the following text, we will first discuss important biological factors to be considered while designing visual stimulation.

## 2.1.1 Vision Properties

Many animals possess the ability to visually obtain information about their surrounding environment. Each these animals possesses specific visual properties that are quite different from another in terms of depth, color, motion, etc. These properties affects the way in which information is captured by the visual organs. It is crucial to learn the important properties of the sensory organs to create realistic simulations. The vision properties of humans and other animals have been studied extensively in biology. We have outlined some key vision properties which play a major role in designing visual stimulations for humans. The same properties in context of other animals are covered separately in chapter 6.

- Field of view: It is the area/volume of the environment observed by the eyes at any given moment. It is usually defined in terms of degrees and varies widely between different animals. The human field of vision is, for example, 210° horizontal range and 150° vertical range.

- Spectral resolution: It is defined as the ability to visualize different spectrum of light. Humans have trichromatic vision, which means human can see in combination of Red, Blue, and Green.

- Spatial resolution: It is defined as the ability to see spatial detail ans it is measured in degrees. This property is comparable to a term used to define image detail i.e. resolution.

- Depth perception: It is the ability to perceive depth. Different animals use different cues such as stereopsis, motion parallax or focusing, overlap, shadow, vertical distance to the horizon, retina to image size ratio, perspective and texture [149].

- Motion perception: It is the ability to perceive motion. There is a critical threshold value beyond which a flickering pattern appears as a continuous motion to the observer. The technical term is called Flicker Fusion Frequency and it is measured in Hz, for humans it is approximately 25 Hz.

Modern display technologies such as screens or projectors are designed to meet requirements of the human visual properties e.g. spectral or spatial resolution. Earlier we mentioned that artists used the knowledge of visual properties for creating illusions. They were able to create perception of motion (Fig. 2.3) and depth (Fig. 2.3) by using two dimensional surfaces. Nowadays, digital technology is used to display dynamic stimuli using 2D screens or projection surfaces. The stimuli are designed to be consistent with the visual properties of the observer i.e. to create the desired illusion of depth, color, and motion. It is even more challenging to design interactive stimuli using 2D display technologies. Vision properties are cleverly exploited for creating illusions of life like visual experience in three dimensions. In the next sections, we will discuss such techniques and their technological components required.

## 2.1.2  Real and Virtual World

Over the years many concepts are designed to provide artificial visual stimulations using display technologies. Some of the first approaches include the use of simple slide based projectors or cathode ray tubes for displaying static or dynamic scenes [88]. These technologies were used to display real-world scenery captured via a still camera or video camera. These were mostly open loop approaches where stimulus is displayed in fixed manner. After that the concept of "*Virtual Environments*" (VE) was introduced with the intention of providing a closed-loop method for visual stimulation [199]. The stimuli are constantly updated to respond to the viewer with the goal of mimicking the sensory experiences of the real world. The content displayed as stimuli includes artificially edited recordings of video sequences [88] or synthetically rendered images using computer graphics [11, 149]. VEs are divided into many categories and subcategories in the technical literature. The divisions are based on the content or the techniques used for presentation of the stimuli.

### Reality-Virtuality Continuum

All different categories of VE are defined as some form of reality [132, 140, 193] e.g. Augmented Reality, Virtual Reality. The nomenclature is consistent with the motivation of creating an alternative form of reality using virtual content. The categorical distinctions are mainly dictated by the extent of virtual content displayed to the observer. For easier understanding, we present these categories in a simplified version of real-virtual world continuum. Figure 2.4 presents a pictorial view of the reality-virtuality continuum where real world and virtual world are displayed at two ends. The left end represents real world where the observer does not see any artificial stimulation. The extreme right end represent a completely Virtual Reality (VR) where the observer is blocked from reality and only the virtual content is displayed. Between the two realities we have two broadly defined categories called Augmented Reality (AR) and

**Fig. 2.4.** A simple depiction of Realiy-virtuality continuum.



**Fig. 2.5.** *Reptiles* M.C.Escher, 1943. Humans are able to infer three dimensional information from two dimensional scenes using their experience and imagination. Escher has carefully introduced depth cues (e.g. perspective, shadow, scale) that lead the viewers to a conclusion that the image is a pictorial representation (2D) of a real world (3D) scene. The reptile creates a conflict for the viewer as it appears both as a three dimensional entity and a two dimensional drawing. (*Source:* [61] This way Esher is successfully able to create a paradox for the viewer by using imagination of the viewer. *License:* CC-BY)

Augmented Virtuality (AV) where the observer is able to see both real world and virtual world embedded in same space. In the technical literature often the term Mixed Reality (MR) is used to represent mixed environments. Recently, a more generic term **XR** (extended reality) is being used as an umbrella term to represent AR, AV, MR and VR.

It is important to point out that the research in VE is ongoing and there are many variations of virtual continuum presented in the literature. Milgram and Kishino presented [140] the virtual continuum with focus on Mixed Reality. Stapleton [193] highlighted the importance of having imagination added to the real-virtual continuum. The role of imagination in the context of creating believable artificial stimulations is discussed later in this section. More recently, Mann et al. [132] discussed a need for creating new terminology and presented a time line of introduction of different terminologies. The authors have presented a new concept of **All Reality**, which may serve as an umbrella term to represent all different forms of realities in the future. The readers should refer to these publications for detailed understanding of the terminologies and ideas behind offering such distinctions.

Imagination performs a big role in the acceptance of sensory stimulations, here we limit the discussion about sensory stimulation of humans. Stapleton and Davies [193] have discussed this topic in length in context of reality-virtuality continuum. The observer is aware of the artificial nature of the stimulus and yet imagination of the observer, if triggered correctly, is able provide a sense of immersion through artificial stimulation [193]. This is especially true for applications of VE that require a sense of immersion e.g. story telling. It allows the observer to ignore conflicting sensory experiences and especially when multi-sensory experiences are not available. It should be noted that imagination plays an important role in acceptance and application of the tool. As an example see figure 2.5 made by renowned artist M.C.Esher [61, 65]. Esher has brilliantly introduced a sense of wonder into a two dimensional image using imagination of the viewer. Even today simple sensory augmentation techniques used in cinemas and theaters are able to provide the viewers a good sense of immersion without being interactive. Similarly, videos games although are interactive medium of stimulation where the sense of immersion for the player is enhanced by provoking imagination with rich interactive visuals. In conclusion, we can say that interactive stimulation technique for humans must consider the aspect of capturing user's imagination to enhance the acceptance of the application.

## 2.2  Augmented Reality

Augmented Reality is a variation of VE in which virtual content is augmented over the real world. The goal is to present the virtual content such that it appears and behaves as objects present in the real world. All Augmented Reality applications must have three inherent properties,

- Real and virtual worlds are combined

- The virtual content is interactive in real-time

- Virtual content is registered in 3D

which were first described by Azuma [11] in one of the most influential survey paper on AR applications.

The idea of augmenting virtual content in the real world is mainly to provide enhanced information about the real world (Fig. 2.7 and Fig. 2.6). The augmentations are designed to assist the user to perform a specific task such as surgery [11, 81], manufacturing [175, 184] or to provide entertaining interactive experiences [23, 27]. Augmented reality is considered very useful because it can provide both shared and personalized experiences. This means that multiple users can look at the same scene and get customized visualizations [26] or many people have displayed the same visualizations simultaneously for shared experiences e.g. projection mapping [173]. The augmented content can be simple geometrical shapes and texts or complex three-dimensional objects [11, 74]. The users can interact with the virtual content while the augmentations remain registered to the real world. Accurate registration means

**Fig. 2.6.** Left: Magic mirror is an augmented reality system displaying augmented content on the user's body for educational purpose [30] ©2012 IEEE
Right: The image shows another example of medical AR application where 3D model is augmented with informative text for educational purpose. *Author:* zedinteractive, *Source:* Wikimedia commons, *License:* CC0



**Fig. 2.7.** Left: Image shows 3D model of industrial part augmented in front of the presenter *Author:* Eawentling *Source:* Wikimedia commons *License:* CC-SA-3.0
Right: Application of AR in museum where additional information about the art piece is augmented in the tablet screen. *Author:* Kippelboy, *Source:* Wikimedia commons, *License:* CC-BY-3.0

that the augmented content "sticks" to the real world and new images are rendered to match the perspective of the viewer. In AR applications, it is more important to present accurately registered virtual content rather than realistic-looking virtual content [11]. The realistic appearance of the virtual content is an added advantage but not an absolute necessity.

In terms of implementation, there are various ways of implementing Augmented Reality (AR). The implementations vary based on the technology used to achieve interactivity and the methods used to display the augmented content. Design variations are motivated by the requirements of the application domain. Complete information on different technologies and methods used in AR is beyond the scope of this thesis. In this section, we will present major design concepts with suitable example applications. Our goal is to provide a necessary foundation for understanding the AR research presented in this thesis. More ideas on AR applications can be obtained from literature reviews specific to AR applications [11, 21, 27, 70].

**Fig. 2.8.** (left) A view-master model G, introduced in 1962. It is a special type of stereoscope where two images are displayed simultaneously using transparent color film. *Author:* ThePassenger, *Source:* Wikimedia commons, *License:* CC-BY-SA 3.0.
(right) Ivan Sutherland wearing the head mounted augmented reality display. © Ivan Sutherland, reprinted with permission.

## 2.2.1  First designs

The first application of the concept of Augmented Reality was demonstrated by Ivan Sutherland in 1968 [199] (Fig. 2.8). The goal of the project was to put virtual information around the user such that it would appear as other real world objects place in the room. Sutherland designed a head mounted display mechanism where the user can see the virtual content in different parts of the room. The virtual information appears three dimensional because separate images are displayed to the eyes with slightly different perspective. This idea is similar to using a stereoscopes (see Fig. 2.8). The system is programmed to update the images constantly to match perspective of the user. On basis of psychological studies Sutherland argued that change of perspective imagery with movement of the head (i.e. kinetic depth effect) is more important than disparity is important for perception of depth.

Figure 2.8 shows Sutherland wearing the head mounted display (HMD), also referred as optical system in the paper [199]. It consists of a specially designed spectacles with two small cathode ray tube displays (one for each eye) to render the virtual content. The displays are reflected through a half silvered mirror arrangement which allows viewing the real world and virtual objects simultaneously. This concept in modern designs is known as see-through head mounted displays. The direction of the user's perspective is inferred from the position and orientation of the optical system. The position of the optical system is tracked in the coordinate system of the room (i.e. reference space) using mechanical and ultrasonic sensors. Position of the virtual objects are defined in the room coordinate system and they are transferred to the coordinate system of the optical system using a transformation matrix. Finally, the coordinates of the virtual content are projected from the optical system (3D) to the display screen (2D) using principles of projective geometry [85]. The virtual content is always rendered from perspective of the user and the appearance of virtual objects change when user goes towards them or away from them, just like the real world objects.

Sutherland demonstrated that virtual objects can be displayed as 3D object hovering in the air or they can be registered to the objects in the room such as a table or a wall. In the original application the observers are displayed a virtual room around the user with numbered

walls, ceiling and floor. In another example observers are introduced to the 3D chemical bond structure of cyclo-hexane. The user is able to interact with the objects, move them in 3D space and even create new objects. The virtual objects are registered correctly to the real world objects provided that positions of real world objects are known in the reference space. Pose computation of real world object is important in cases where virtual objects are only partially rendered for certain perspectives as they are occluded by real world objects.

## 2.2.2 Design concepts

Sutherland's ideas opened a new line of inquiry in the field of Augmented Reality. To date, several design concepts have been proposed where augmented content is displayed using different tracking and display technologies. We will describe them in two broad categories based on the visualization methods i.e. head-mounted displays and spatial displays. For the sake of simplicity, we will present spatial displays in two different categories: display based and projection-based. In principle, the application domain (e.g. medicine, industry, etc.) determines the selection of the design concept in terms of visualization and tracking approach. It is common to select the visualization scheme first and then implement the best possible tracking approach for the use case. Tracking involves the computation of the user's viewing direction, the position of the display screen, or the position of the objects in the real world. Tracking methods are often interchangeable and therefore we will discuss them separately discussed in more detail in Sec 2.4. The reader should note that the primary research contributions of the author are related to the development of tracking approaches for AR/VR applications.

We would like reiterate that the information given here is by no means an attempt to cover all modes AR visualizations. The readers may refer to review papers in orders to get complete grasp on the topic [11, 21].

### AR with Head Mounted Displays (HMD)

Visualization approach with HMDs requires the user to wear a device on the head such that the optical system is placed in front of the eyes (see sec. 2.2.1). There are two types of HMD based approaches, optical see-through HMDs and video see-through HMDs [11, 180]. The working principle of all HMD based methods is more or less similar to the concept demonstrate by Sutherland. This means 3D information is transferred to the user's perspective by tracking the user's perspective in 3D.

**Optical see through HMDs** use optical combiners to display merged view of the real and the virtual world e.g. transparent display or mirror reflection (Fig. 2.9). Optical see-through HMD based designs are proposed for a range of applications [11] in medicine, manufacturing, and entertainment, etc. The main advantage of this technique is that the user does not have to hold the device and augmentation is presented in a naturalistic way i.e. along the viewing direction. This design allows the user to perform regular tasks while having an augmented view. It is also possible to control or manipulate the augmentation using audio commands or gestures [129]. In medical applications, a typical use case is for displaying medical scans (CT, x-ray) or vital stats augmented over the patient's body [19]. The augmentations are useful while performing surgical procedures [19] or more generally for education and training

**Fig. 2.9.** Two designs of head mounted displays. Left: The image shows an optical see-through HMD where a screen or see through mirror is placed in front of the user's eves. *Author:* Shyuan1977, *Source:* Wikimedia commons, *License:* CC-BY 4.0
Right: Visette45 SXGA Video See through head mounted display shows cameras mounted in front of the device. *Source:* Cine Optics



**Fig. 2.10.** A worker receiving visual instructions for repairing industrial part through HMD

purposes [78]. It can be argued that the doctors are used to wearing headgears for other procedures and therefore they accept such designs. In manufacturing applications, a typical example is training assistance for maintenance and repair as depicted in Fig. 2.10 [153, 175]. The user can see the inner working mechanism of the machine and internal components are displayed as if the machine is transparent. The visualizations are configured to provide step with guidance to the user such that visualizations update as the user progresses through the maintenance job.

Optical see-through methods are considered computationally cheap because only virtual content is rendered. However, one of the main challenges is to display virtual content without losing the aesthetic appearance of the real world i.e. light, colors, etc. In some applications, the special techniques are used to block a certain wavelength of light for improved vision. This approach remains popular to this day and many commercial companies are constantly coming up with new designs of optical see-through HMDs. Nowadays it is possible to produce see/through HMDs with transparent displays using "holographic projection films" that show projections diffused in to the screen e.g. Microsoft Hololens or Magic Leap. We have proposed

an AR application to visualize the 3D movement of birds using the concept of optical-see through HMD. The details of the application are discussed in detail in Chapter 6.

**Video see-through HMDs** are designed to block the real world view of the user and the augmented view is seen through the display screens placed in the HMD (Fig. 2.9). The real world view is captured by cameras attached in front of the HMD. The cameras are positioned in such a way that they produce a stereoscopic view. Virtual content is added to the videos in real-time before displaying the augmented view to the user. It is possible to remove details of the real world using image processing techniques, also known as diminished reality. Video see-through HMD devices are also suitable for other variations of VE i.e. Augmented Virtuality.

This approach is computationally expensive as it requires very fast tracking and image processing operations. Fast computational techniques and computers are required to update the augmented view as soon as the user changes perspective. This is often challenging to achieve and slow processing introduces a lag in the system which may lead to motion sickness [180]. Another disadvantage with such an approach is that the real world view is completely blocked and failure of equipment may cause discomfort for users or even harm patients in case of surgery. Such designs are considered less suitable for medical or manufacturing tasks where unwanted accidents may occur due to technical failure.

A major limitation with HMDs is that users have to wear the equipment. Wearing HMD is cumbersome and many users do not prefer to use it. Therefore, the weight of the HMD is an important consideration for application designers. It is challenging to produce small high-resolution displays that fit into the HMD. Mobility of HMDs is limited to an area if tracking is done using external sensors i.e. outside-in approach (details in 2.4.4). To avoid this problem modern HMDs are mounted with sensors on the HMD itself to track real-world objects in 3D [129]. This approach requires on-board computation which comes as a trade-off for mobility since every additional sensor increases the weight of the system. One of the most important challenges with HMDs is to provide a realistic field of view. The peripheral vision of the user is blocked which may make the user uncomfortable.

On a positive note, display technologies are improving rapidly and tracking sensors are becoming smaller which is a promising sign for designing smaller and lighter hardware. Technology companies like facebook and google are investing in the development of a new designs for see through HMDs that resemble wearable glasses. Their vision is to bring AR in daily life interactions for applications in entertainment and communication. It is likely that users of specialized domains such as manufacturing or medicine will only accept novel AR solutions if the applications solves the inherent need of the field without overheads.

### AR with display screen

This concepts involves viewing the augmented view through a display screen placed in front of the user e.g. computer monitor or cell phone display. This configuration is used for applications where users can access the information on demand and do not always need the augmentations. Display screen based AR approaches are designed to be in fixed or mobile configurations as per needs of the application. We have defined four configurations based

A. Moving display - Moving camera

B. Moving display - Fixed camera

C. Moving display - Fixed camera

D. Fixed display - Fixed camera

E. Fixed display - Fixed camera

F. Fixed display - Moving camera

**Fig. 2.11.** Different configurations of display screen based AR solutions.

on the position of display and camera 2.11. Their advantages and disadvantages are also explained with example applications.

**Fixed display and Fixed camera** approach is one of the easiest to setup for AR applications. The augmented content is displayed from the camera's perspective. The viewpoint the camera is fixed and this means knowledge of static objects in the real world is exploited to create rich visualizations as some visuals are per-configured. The use of a camera can be avoided if the real world scene is completely static [26].

A good example venue is a museum where a display can be used to provide augmented view of an artifacts placed in front of the camera e.g. JURASCOPE [1]. If the camera is mounted directly behind the display it can also be called a video see through configuration (see Fig. 2.11). A similar application is also valid for archaeological sites where augmented views can be displayed to digitally replace the broken or incomplete places [214]. A slight variation to this approach is the augmented desk configuration where the camera is placed at a specific location [3]. For this configuration the camera may not be in the immediate vicinity of the viewer. For example, watching augmented views of a remote sport telecast or CCTV footage. Virtual mirror configuration is used when interactive visuals are displayed on the user. The camera is facing the user and display acts as a mirror. For example, magic mirror system [30] is used to display medical information is augmented over the body of the user.

The main advantage of these methods is that the user does not have to carry any hardware and can exploit wide range of augmentation options based on personal preferences. One common problem is that stereopscopic cues can not be displayed as rendering two separate images for two eyes is not possible. One idea to circumvent this issue is to use 3D glasses with screens [11]. Fixed configurations do not offer mobility as the user is require to be at the location of installation. Such installations are cost effective option for shared spaces where many people can avail the facility (e.g. museums or work bench).

---

[1]https://artcom.de/en/project/jurascope/

**Fixed display and moving camera** approach is used to display augmented view of the real world seen from perspective of a moving camera. In a typical use case, the user operates a device mounted with the camera and augmented information is available to the user via screen. This approach is chosen in cases where viewer requires an additional view for assistance with a specific task. For example, endoscopes mounted with cameras are used for navigation in the body while performing surgeries [117]. Similarly, AR enabled displays are used for structural inspection using aerial robots [166]. The augmented view is vital for decision making in such cases. In chapter 3, we will describe the use of such configuration in manufacturing assistance application. We have used a mobile stereo camera setup connected a display screen to show augmented information about the task (Fig. 3.9).

The moving camera configuration is extremely useful when direct access to the real world scene is difficult e.g. inside human body, remote locations. Unlike HMDs the displays screens are placed at designated locations and the user has to look at the screen which can be distracting. However, this can also be an advantage when the user can choose to perform the main task and only use augmented views for verifications. Such designs can be supported with powerful computers attached to the display. Fast computation and data transmission is crucial to minimize the lag and provide real time performance.

**Moving display - Fixed camera** approach is an extension of the fixed display - fixed camera approach. It is used when display locations is required to be flexible as per needs of the user for additional mobility or collaborative tasks. For example, industrial scenarios where parts on a conveyor belt are scanned by rigidly mounted cameras for quality check or counting purposes. Different users are able to view the workpiece using a fixed screen or a mobile screens e.g. tablet. The augmented display shows information about parts such as identity, number etc. The worker is able to control the augmentations by interacting with the display. Such configurations are helpful in collaborative environments when multiple users are consulted for a task.

The configuration is useful where the computational requirements are really low. However, data transmission limitations may introduce a lag. The camera perspective does not change due to rigid design and therefore perspective based rendering is not required. In most cases, tracking object position in the image is sufficient for 2D augmentations. Additionally, the augmented content is independent of movement of the display. The interactivity is limited to user interactions on the screen.

**Moving camera - Moving display** approach is one of the most popular spatial display based AR configuration. The development a this approach has increased substantially due to availability of mobile devices with good cameras and fast computing power i.e. smart phone, tablets. A typical device has a display on one side and video camera on the other side, also known as video-see through configuration (Fig. 2.11). The user points the camera at the object of interest and the augmentations are displayed in the screen facing the user. This configuration is more intuitive because it provides the user "a virtual window in the real world" (see Fig.2.6 and Fig. 2.7). Remove viewing configurations are also possible with camera mounted on drones. The user is able to move with the display and visualize the scene from aerial view. Such configuration is used in sports entertainment [159], rescue operations [118].

This video see though concept is popular as it provides augmentation on demand for the users. Nowadays, low-cost smart phones and tablets are ubiquitous and therefore very suitable for developing AR applications for masses. Such devices already come optimized software framework to support AR applications e.g. ARCore [2] for android and ARKit [3] for iOS. The user do not have to carry or purchase a special device such as a HMD or wearable glasses. Most importantly mobile configurations allows use of AR in both indoor and outdoor environments. Applications of this concept is shown in multiple domains ranging from entertainment and gaming to manufacturing and medicine [107, 183].

Real-time augmentations with devices having limited in size and computational power is challenging. The choice of tracking approach is crucial for achieving real-time augmentations (discussed in 2.4). However, video-see through configurations do not require video transmission which reduces the delay. External sensors are deployed to reduce computational needs and to track both display and the real world objects. Tracking with external sensors often limit mobility to a specific area (details in sec. 2.4.3).

### AR with projectors

Augmented reality with projectors is also referred as Spatial Augmented Reality (SAR) techniques [173]. In this approach, one or more projectors are used to directly augmented the objects in the real world. The augmented objects can be three dimensional object [15, 174] or a flat surfaces such as screen or a wall [24]. 3D information of the real world object is used to accurately register the augmentations. The orientation of real world objects w.r.t projectors are computed using tracking sensors. These setups require a one time calibration procedure in order to transfer 3D coordinates from sensor coordinate system to projector coordinate system [27, 173]. The projected content can be altered based on perspective of the user if user's movement is tracked, otherwise it is rendered from perspective of the projector [26]. The projector type is selected based on the type of application i.e. color or laser projection. AR with projectors can be implemented in two configurations based on mobility of the projector i.e. static and dynamic. The principle behind both approaches is the same but the application of dynamic configurations are more complex.

**Static configuration** is useful when the projector is expected to provide augmentation in a fixed region. The projector is fixed to a ceiling or a wall for indoor applications. This

---

[2]https://developers.google.com/ar
[3]https://developer.apple.com/documentation/arkit

**Fig. 2.13.** Spatial AR application for entertainment and manufacturing assistance. Left: sandbox installation for teaching geoscience education. A camera-projector setup is attached rigidly to the ceiling above a table with box of sand. The projection is limited to boundaries for the sandbox. Depth camera is used to reconstruct surface of the sand in real-time and topographical information is projected with different colors. Users are able to modify the surface by moving sand and the projections change accordingly. *Author:* Karlbrix, *Source:* Wiki Commons *License:* CC3.0.
Right: Laser projector projecting a circle on a industrial part. This is used by the worker for manufacturing guidance or quality inspector for verification of job © *Extend3D*.

configuration is relatively easier to set up in a static scene i.e. objects do not move [15]. This concept is popularly used for projection mapping applications where one or more projectors are used to provide vivid visual displays for entertainment purposes [27]. Visuals from the projectors create anamorphic illusions where two dimensional surfaces appear three dimensional (or vice-versa) [25] and static objects appear dynamic [174].

Static projector configurations are commonly used for large scale entertainment events such as projection on large buildings or castles (see Fig. 2.12). If target object is also static, a one time calibration is performed using the 3D model of the real-world object and real-time tracking is not required. The augmentation remains registered as long as the target object or the projector does not move. A calibrated camera-projector setup is used for projecting on movable objects [15]. For projecting on 3D objects it is vital to know the geometry of the objects for accurate registration. The advantage of such setup is that 3D reconstruction of the scene can be done by projecting a known sequence of images on the scene [177]. In case of dynamically changing scene or objects real-time camera tracking methods are deployed [25, 99, 100, 115]. The perspective of the observer is tracked for some SAR applications for personalized experiences e.g. museums [26].

Static configurations are less flexible but extremely useful in providing rich visual experiences due to possibilities of pre-configuration. The computational power is not a limitation for such configurations. It is cost effective approach for shared spaces e.g. museums or schools, however multiple installations may increase overall cost for setup.

**Mobile configuration** is used when the application requires projectors to be used as mobile augmentation devices. Portability increased the range of applications. There are several mobile projection based AR applications proposed in medicine [63] and entertainment [35]. Often projections are used as input devices e.g. virtual keyboard, where user interaction is captured using cameras [35, 142]. Often small projectors are preferred for their compact size and less weight which makes them easier to carry. Smaller designs do not always provide enough contrast and there are technical challenges in manufacturing such devices. Another variation of mobile configuration uses off the shelf projectors combined with cameras in a

portable housing [67, 188]. We have used this configuration for our research in Industrial Augmented Reality (IAR) applications (details in chapter 4).

Mobile setups are useful in industrial scenarios because they can be moved around the in the workshop for different applications. A major advantage of using spatial augmentation is that multiple people can view the augmentation simultaneously as virtual content is directly projected on the object. This is particularly useful in industrial scenarios where multiple users are required for decision making. Another strong point of projection based methods is the possibility of hands free operation. Hands-free operations are preferred in industrial settings because users have manufacturing tasks to carry out. Augmented reality serves as perfect companion in this scenario as the users can focus on the task at hand with help of projections on the work part (See figure 2.13). DLP (Digital Light Processing) projectors are used for applications which benefit from colorful visualizations e.g. rapid prototyping or data visualization [67]. Often industrial environment is are very well lit and DLP projections do not provide sufficient contrast. In such cases, laser projectors are used instead of DLP projectors. Laser projectors are suitable for manufacturing guidance tasks where only limited information is sufficient e.g. annotation or marking task as depicted in Fig. 2.13. It is useful to have the camera-projector setup in same housing for the position between the two devices does not change with respect to each other. This way positions remain rigid the system requires one time calibration which hold for longer time.

In general, projections (static or dynamic) are limited to a specific area. One way to alleviate this problem is to use multiple camera-projector setups. However, they require more complex calibration procedures for alignment of projections and powerful computers for processing and rendering [173]. Research in this field is extensive [23, 173] and automatic calibrations protocols are developed for such setups [100, 222]. Line of sight is another problem and augmentation are easily distorted when other objects or humans can occlude the projections. This problem is more prominent in static configurations and can be partially solved to some extent by using mobile configurations. The operating frequency of the projections must be fast enough to provide smooth transitions between different visualizations. This limitation is more relevant to the laser projectors as the lines drawn per second are limited. Similarly, fast tracking is necessary if target surface or object is moving.

### Special cases

There are some special cases which are used with both screen and projection based spatial AR categories. Therefore, we are mentioning them separately.

**Transparent displays** are used to display augmentations in optically see-through manner. These displays are used with fixed display and mobile display configurations. Transparent LCD [4] and OLED[5] displays using active display components are being developed from some years [21] for hand-held and mobile applications. Another method of creating similar affect is by projecting information on transparent screen using a screen or a projector by using a beam-splitter mechanism [26]. Such implementations do not require the user to wear a head mounted gear but the user is often required to operate in a specific area. In future, small and mobile transparent displays could be extremely useful for using AR with handheld devices.

---

[4]https://crystal-display.com/products/transparent-lcd/
[5]https://otilumionics.com/transparent-displays/

**Retina projection** is a technique where the images are projected directly on the retina of the user. This is a specialized technique where small projector is integrated into head mounted device [186]. As the image must be projected correctly on the retina. This approached requires retina tracking method in addition to tracking of real world objects. Retina projection based augmentations are clear and spatial resolution of such augmentation is high. Technology required for this approach is specialized and expensive for mass production. In the near future, such designs could become reality but as of now there are technical limitation for designing and using such devices.

**Fig. 2.14.** Alternative reality experiences with VR. (left) An advert for a VR cinema experience in the Netherlands. *Author:* Tero Koistinen *Source:* Wikimedia Commons *License:* CC-BY-SA-4.0
(right) An engineer navigation through four million stars in the milky way in a VR application (Point-CloudsVR) designed for astronauts to classify star groupings. The controllers are used to navigate *Author:* NASA/Chris Gunn *License:* CC-BY-SA-2.0

## 2.3 Virtual Reality

One basic property of Virtual Reality (VR) is that the user is disconnected from the real world. Visualizations displayed to the user may reflect a fantasy world, abstract patterns or mimic the real world. The visual cues such as depth and motion perception must be consistent along with the aesthetic appearance of objects. Stimulus presentation plays an important role for acceptance of the artificial world. Interactivity is another key property of VR, which directly affects the sense of immersion. The user must be able interact with virtual world which includes manipulation of objects and navigation in the virtual space [88, 199]. Different mechanisms are used to provide interactive experiences in the virtual space (see Fig. 2.15) e.g. joy stick, virtual keypad, haptic gloves or audio commands. These devices are used to get input from user about interaction and simultaneously provide the user higher degree of immersion by engaging multiple sensory organs [132]. Ideally, the means of navigation and interaction in the virtual world should be as naturalistic as possible i.e. navigation by walking and object interaction with hands. Fully mobile VR setups with gesture based interaction are not possible yet. However, navigation by walking is partially possible with treadmill based approaches, e.g. Virtualizer Elite [6].

The responsiveness of the visual stimuli is extremely important in the virtual space. The user's perception of self-motion and balance must be considered when designing the response time of the stimuli otherwise it may cause motion sickness for the users [102, 158]. This may discourage the users from using the technology on regular basis. Sensory mismatch is caused by poorly designed stimuli or technological limitations of displaying the stimuli, e.g. refresh rate of screen or rendering jitters. This problem is addressed promptly in last few years by improving real-time tracking and graphical rendering methods. It should be noted that VR is used extensively for studying sensory perception in humans and other animals such as insects, fish and small mammals [149]. Through interactive experiences we are able to understand the functional properties of sensory organs (mentioned in 2.1). Knowledge of functional

---

[6]https://www.cyberith.com/virtualizer-elite/

**Fig. 2.15.** VR design concepts for navigation and interaction. (left) Two strategies of interaction in VR using gloves with sensors. Image on top-left shows mobile HMD and image in bottom-left shows VR headset connected to robotic arm.*Author:* NASA National Aeronautics and Space Administration, *Source:* Wikimedia Commons.
(center) VR suit A VPL Research DataSuit: a full-body outfit with sensors for measuring the movement of arms, legs, and trunk. Developed circa 1989. Displayed at the Nissho Iwai showroom in Tokyo *Author:* Dave Pape, *Source:* Wikimedia Commons, *License:* Public domain.
(right) Virtualizer Elite: VR concept for simulating naturalistic movement for navigation and handheld controllers for interaction in the virtual world.

properties is directly useful for improving the existing methods for virtual stimulation [149]. We have introduced this topic as fundamental research earlier and more discussion on this particular application of VR is addressed in Part 3 as one our core research contribution.

In terms of implementation, different approaches have been introduced in terms of design (Fig. 2.15). Design configurations in VR, like in AR, have two main components i.e. tracking and visualization. Tracking involves the mechanisms to provide interactive (closed loop) feedback to the user. This mainly includes methods to track change of user's perspective (in 3D) and methods to get input from user for manipulation the virtual stimulus. Tracking approaches for VR and AR methods are quite similar, therefore we will discuss tracking approaches separately in Sec 2.4. In this section, first we will describe one of the earliest demonstration of the concept of multi-sensory stimulation in human users. Following that we will describe two popularly used VR design concepts based on the approach chosen for visualization of the stimuli. Each will include a short discussions on advantages of using each approach with relevant application examples.

### 2.3.1 First Designs

In sec 2.1, we mentioned that some of the earliest ideas of creating alternative reality using sensory manipulations came from the art world. The act of providing immersive experiences

**Chapter 2** Artificial Sensory Stimulation with AR and VR

**Fig. 2.16.** Illustration of Sensorama from U.S patent no.3050870. one of the first setup designed to provide multisensory stimulation to the user. Morton Heiig designed the setup in 1962 and called it a "Experience Theater". He made 5 short films which included stereoscopic visuals, stereo sound effects, olfactory cues and haptic cues in form of vibrations. *Source:* Wikimedia commons *License:* Public domain.

the entertainment world does capture the true spirit of Virtual Reality. The first use of the term "Virtual Reality" was mentioned in the context of theatrical performances in 1938 [9, 132].

One of the first design concepts of virtual environments was demonstrated by Morton Heilig in 1960s (Fig. 2.16). As a filmmaker Morton's idea was to design a rich and entertaining cinematic experience for users. Cinema experiences are limited to two dimensional visuals and covered only 18% of the observers visual field of view. To overcome this problem Morton invented "Sensorama", one of the first prototype of a multi-sensory (visual, auditory, olfactory and haptic stimulation) stimulation system for providing immersive experiences [8]. As illustrated in the figure, the user is seated in front of the machine and head is placed in the canopy consisting the optical system. It is designed to block the lights coming from the surrounding environment. The viewer is only able to see the pictures presented through the eye piece. Heilig used a custom built multiple camera setup to capture panoramic visuals to cover the user's field of view. Additionally, stereo sound effects are added for auditory stimulation, fans and chemicals are used for olfactory stimulus and slight vibrations are used for tactile stimulation. One of the films was about recreating experience a motorcycle trip through a city. The visuals of the road journey were accompanied by other sounds on the road, a gentle breeze bringing odors from the street and vibrations simulating movements of motorbike on a road with occasional potholes.

**Fig. 2.17.** One of the first VR headset (1985), made by Jim Humphries and Mike McGreevy at NASA's Ames Research Center in California. The system created a computer-generated image of what a pilot might see during an actual flight. It was intended to test concepts of presenting visual information to pilots or astronauts. Sensors tracked the movement of the wearer's head, and the displayed images moved accordingly. James Humphries donated this to the Museum in 1997. Picture taken at the National Air and Space Museum's Steven F. Udvar-Hazy Center in Chantilly, Virginia, USA. *Author:* Sanjay Acharya *Source:* Wikimedia commons, *License:* CC-BY-SA 4.0.

It can be argued that *Sensorama* does not allow user's to interact with stimuli and the user can not control the stimuli. This being said, the ideas presented with *Sensorama* were very influential in inspiring the development of interactive virtual experiences. Morton proposed that the concept was capable of providing much more than just entertaining cinematic experiences and stimulations can be used for military, education or industrial applications. Heilig predicted that such installations will be a cost effective method for demonstrating conceptual ideas to students or workers instead of using a text based approach book or a manual. For example, displaying structure of molecules or working principles of complex machines. These predictions proved to be valid and today we have several off-the shelf products available for education and training applications in VR [87, 126, 210].

Morton also introduced the idea of creating shared virtual spaces for multiple users. The patent includes a concept for displaying visuals to multiple users (four) at the same time. The idea of personalized shared experiences is that each viewer is present in a shared virtual environment and yet it is possible to manipulate visuals for each viewer individually. Taking this concept as inspiration the modern VR systems provide unlimited opportunities to study group behavior. Virtual spaces provide means to manipulate sensory stimuli presented to each individual which is not possible in real world.

It is worth noting that the concept of Sensorama failed to make impression on the entertainment industry as it was too expensive to produce such films. It became clear that the ideas had great potential if developed further with special focus on interactivity. The research in VR was pursued strongly by military, industrial users because of expensive hardware requirements.

**Fig. 2.18.** Modern VR headsets designs. (left) Sony play station 4 PSVR dedicated for playing games in VR *Author:* Evan-Amos, *Source:* Wikimedia commons, *License:* CC-BY-SA 4.0. (right) HTC vive pro headset for multipurpose VR applications. *Author:* KKPCW, *Source:* Wikimedia commons, *License:* CC-BY-SA 4.0.

Nowadays, research in VR is pushed for entertainment application, especially by the video game and cinema industry [8].

## 2.3.2  Design concepts

There are two design concepts that are popularly used for visualization in VR i.e. head mounted display and CAVE. Head mounted displays are wearable headsets which block the outside view for the user (Fig. 2.17 and Fig. 2.15). HMD based approach for VR was popularized soon after Sutherland demonstrated idea of head mounted displays for Augmented Reality (covered in sec 2.2). CAVE systems are more similar to spatial AR methods where user in surrounded by displays and large part of the field of view is covered by the display screen or projections. The user is surrounded by the stimuli without having to wear any device on the head. In both methods the visuals are always generated from the perspective of the user. In the following text we will cover both approaches, relevant tracking approaches are covered in the next section.

### VR with HMD

HMDs for VR are closely related to video-see through HMDs explained in the previous section with AR approaches (sec. 2.2). One major difference is that HMDs are not designed to display the outside world. HMDs provide portable and mobile means of deploying VR solutions. While the mobility of the user in the virtual world is unlimited, the mobility of user is still restricted in the real world due to physical constraints in the real world e.g. walls, furniture. This is a major design challenge because ideally the interactions and navigation in the virtual spaces should be as close as possible to real world. Most VR applications are designed to work around this challenge and still provide meaningful experiences.

Often dedicated VR stations are designed to simulate realistic conditions. For example, roller coaster rides in VR involve users sitting in a real roller coaster while wearing HMDs and the seats moves in sync with the visuals. This includes dedicated spaces where user's are provided unique experiences such as flying [191] or underwater swimming with fishes [87]. The user is given limited mobility and the need for complete freedom is eliminated psychologically by creating supporting narrative. In other words, users are less likely to think about walking while sitting in driving seat of a car or a plane. Another alternative is to use treadmill like devices,

which allow users to walk in the virtual world while keeping them at the same location in the real world (see Fig. 2.15). It is shown that 3D tracking methods are used effectively to extend the mobility defined by range of tracking [87, 191]. It is typical to provide users with various types of controllers to interact with virtual spaces.

The working principle of the HMDs based visualization is as follows. The 3D motion of the headset is tracked with respect to a virtual reference coordinate space. All virtual objects are represented in the same coordinate system. The location of each eye is inferred from headset tracking, and a virtual view is computed for each eye. The eye is modeled as a pinhole camera and therefore images are computed by projecting 3D world onto the 2D image plane using projective geometry[85]. The complete image is rendered virtually hence VR images are computationally expensive. Real-time rendering is challenging for highly dynamic scenarios such as flight simulators or games. Lag in tracking or rendering can easily disrupt the illusion.

The stimuli must always correspond with the user's perception of motion. Each animal learns motion cues through sensory experiences obtained while navigating in the real world. In other words, sensory organs are calibrated to continuously sample information from the environment such that we always have a sense of self-awareness and balance, for example, sleeping or standing. We can focus on a moving object with the eyes and even grab it with movement of arms i.e. hand-eye calibration. Providing accurate motion cues is a complex problem. It is easier to design applications where the user's perspective can be predicted partially. This way pre-rendered images can make computational requirements cheaper [164]. New computer graphics techniques and GPU development has eliminated many problems in real-time rendering [22, 128].

Our knowledge regarding vision bases perception is still developing. Therefore, one major use case of HMDs based VR is in the field of psychology studies [149, 158] and the studies related to sensory organs itself [182]. Such studies provide us insights about human perception and affect of different types of stimulation on our perception. In this research direction VR application are developed for non-humans as well [149]. The particular application of VR is covered extensively in chapter 6.

Video game developers are adopting the VR technology (Fig. 2.19) for providing new experiences. Gamers are accustomed to assume the role of virtual avatar and navigate in the virtual 3D world with the help of controllers (with limited buttons) and two dimensional screens. HMD based Virtual Reality has extended this experience by providing 3D experience and more realistic immersion as users have to physically more their perspective to inspect the scene. Nowadays, it is also possible to detect gestures of users with motion capture methods, depth cameras and inertial sensor based controllers (see Fig. 2.15). For example, a game of tennis or golf can be played virtually where the users are required to really move the arms as if they were hitting a real ball. VR presents a unique opportunity to expose the user to experiences which are impossible to provide in the real world. The users are virtually transported to desired locations one the planet or even visit another planets, for example space travel (see Fig.2.20). For training applications VR is relatively cost effective and the user can train for extreme conditions without any risk. Figure 2.20 shows one such example where pilots are trained for using parachutes in simulated conditions.

**Fig. 2.19.** Virtual reality for gaming applications. The pictures are from a recent gaming convention in Spain where gamers got together to try play station gamers with VR headsets. (left) The gamer is playing first person shooting game using a mock gun as a controller in virtual world. (right) The gamer appears to be flying a plane in the virtual world. The visuals of the gamer are displayed in the screen in front. Note that gamer is not able to see the screen, it is for visualization of spectators in the real world. *Source:* Wikimedia commons, *License:* CC-BY-SA 4.0.



**Fig. 2.20.** Virtual reality applications for space travel and combat training. (left) Apollo 11 astronaut Buzz Aldrin and Erisa Hines of NASA's Jet Propulsion Laboratory (JPL) in Pasadena, California, interacting with Destination: Mars experience at the Kennedy Space Center Visitor Complex. Destination: Mars allows virtual visits to Mars using real imagery from NASA's Curiosity Mars Rover. Photo credit: NASA/Charles Babir NASA image use policy. (right) A military personnel receiving training for operating parachute using VR. The simulator prepares the students for rough weather conditions and equipment malfunctions. U.S. Navy photo by Mate Chris Desmond. *License:* Public domain, *Source:* Wikimedia commons

VR experiences are equally popular for designing cultural tourism and educational applications. It is also possible to host virtual events where multiple users from different physical locations are able to come together and collaborate in virtual spaces. For example, recently IEEE VR 2020 conference was organized virtually on the Mozilla hubs platform [7]. The participants from different parts of the globe attended the conference in a virtual space. The participants using HMDs and controllers had much better control of their avatar.

Mobility offered by HMD based VR is limited due to multiple reasons. One of the major reason is real world obstacles such as walls or furniture. This can often lead to unwanted accidents and users are scared to walk into a wall or fall. Another limitation for mobility is requirement of cable connectors for data transmission. This limitation is currently being solved with high-speed wireless transmission protocols, for instance HTC Vive uses bluetooth. Tracking approaches also add limitation to movement, especially when HMD is tracked by external sensors which have limited range. Development of high-resolution video cameras,

---

[7]https://hubs.mozilla.com/

depth cameras and other inertial sensors are solving this problem rapidly. However, additional sensors on the HMD require compromise in terms of weight and performance. Most systems need a dedicated computer for rendering virtual images using the 3D tracking information. Even with the existing limitations, VR is promoted strongly for commercial uses at home and other indoor environments. The cost of a HMD is still relatively high, applications are specific and therefore it is not as popular as other gadgets for personal use such as cell phones, laptops.

It is believed that compact computers with integrated tracking sensors may solve the problem of mobility in future and increase acceptance of VR dramatically. Many companies are have proposed smartphone mountable HMD designs as a low cost alternative to dedicated VR headsets like Google cardboard or Samsung gear VR. Smart phones are inexpensive, have reasonable computation power and are versatile in use. Most are equipped with inertial sensors and cameras which can be used for tracking headset movement. Additionally, on board sensors are useful for tracking real world obstacles which may enhance mobility of the user in VR. Moreover, smart phones are also useful for AR applications. So far this option has not gained real success due to laggy performance and limited resolution. However, the research trend is in favor of developing compact computational devices with multiple sensors which can be used for a wide range applications of interactive visualizations.

Another major design challenge for HMDs is to match the visual properties of the human eye, e.g. field of view and spatial resolution. Existing displays fail to match the vision capability of the human eye. On the positive note modern displays sufficiently meet flicker fusion frequency (motion perception) and spectral resolution (color representation) requirements. Depth perception is better in HMDs because of stereoscopic perspective corrected images. This technique is better than display techniques used in 3D cinema where observer wears polarizing goggles to observe the 3D effect. At the moment, the development cost of miniature sized high resolution displays is high. Improving field of view in HMDs is a difficult design problem. It does add significant value to user's experience but it is not absolutely necessary for many applications.

## CAVE

The concept of Cave Automatic Virtual Environment (CAVE) was introduced by Cruz Niera [45]. The idea is reminiscent of Plato's allegory of the cave [171] mentioned earlier in the introduction (see Fig. 2.1). It is a digital space where the user completely is surrounded with virtual stimuli. The stimulation is achieved by projecting walls, ceiling and floors of a room using multiple projectors [45].

The working principle of the system is as follows. The users position and head movement are tracked in 3D using tracking sensors. Projectors are rigidly placed and their location is registered to the tracking sensors. It is also possible to use an array of large display panels as in CAVE2 [68]. The user's head location and distance from the wall is computed to render a perspective corrected view. In most cases users are given polarizes glasses or shutter glasses for stereoscopic view. Often 3D glasses are integrated with tracking sensors and user's perspective is inferred from position of the glasses. The user is able to move freely in the designated area and can see their own body too. Additionally, controllers are provided to interact with the

**Fig. 2.21.** CAVE setup at Center of Advanced Energy Studies. The image shows a user operating in a CAVE VR setup where visual stimulations are projected in .*Author:* Idaho National Laboratory,*Source:* Wikimedia commons,*License:* CC-BY-2.0

visuals (see Fig. 2.21). Nowadays, motion capture technology is also used to allow gesture based interactions.

This approach is less intrusive as the user does not wear cumbersome hardware. CAVE systems fare better in terms of field of view and spatial resolution [68]. The visuals are surrounding the user and therefore always cover the complete field of view. High resolution screens or projectors are used and pixels are not visible because displays are not very close to the eyes. Multiple users can access the same space and visualize 3D views, however, perspective correction is not possible for all users at the same time. Often it is not necessary to use sterescopic views and 2D visuals are sufficient. CAVE is used as a 3D and 2D visualization tool depending on the need of the application.



**Fig. 2.22.** The image displays students interacting with the CAVE setup at Max Planck Institute of Animal Behavior. The students are displayed migration routes of various bird species on a global map. Such setups are extremely valuable for displaying global events and corresponding changes at various time scales. This method creates more intuitive learning experiences. ©Maxcine, Max Planck Institute of Animal Behavior, Radolfzell.

This approach is used for many different applications in military, manufacturing and tourism [68, 131]. It is useful for designing realistic simulators or improved visualization of geographical maps. Workers or engineers can take virtual training for operating machines in the factory (Fig. 2.21). It is an ideal setup for students can take a virtual walk in world heritage cites to learn about history and art with realistic visualizations(Fig. 2.22). One major advantage is that users can navigate naturally in the virtual space.

There are certain disadvantages in terms of mobility and personalized experiences. The user is restricted to the area and building large environments is expensive. The user's own body or other users can disrupt the visualizations. Personalized experience for multiple users are difficult as all viewers see the same content. The projectors must be calibrated for accurate registration and they must be synchronized to avoid projection lag. Overall 3D tracking and real-time rendering of 3D views using multiple projectors is computationally expensive. Taking overall costs in to consideration, CAVE setups are better suited for applications in shared spaces e.g. museums, companies or schools.

The concept of CAVE was adopted in the early 2000 by biologists to make virtual environments for studying animal behavior in insects. It inspired VR designs for insect, fish and mammals are used for fundamental research in neurology, medicine and bio-inspired robotics [149, 195]. These applications of VR are rarely discussed in the technical community. In this thesis, we have discussed this topic in more detail in Part 3. Our contribution includes first review of the state of the art research in the field of animal behavior using VR 6. Further, we have introduced a novel setup design capable of providing different sensory stimulation to a group of animals (including humans).

## 2.4 Technical concepts

In this section, we will cover technical concepts required to understand functional aspects of XR applications. 3D tracking and visualization (stimulus display and rendering) techniques are the two crucial components for designing the closed loop visual stimulation techniques. In the previous sections, we discussed vision properties (sec 2.1) that play an important role in the selection of display technology and after that we discussed different visualization concepts of AR/VR application in section 2.2 and 2.3. We learned that all XR methods require some type of 3D tracking technology for manipulation of stimuli in interactive manner. In this section we will focus on 3D tracking technologies.

The tracking requirements of XR applications are defined after the application designer selects a probable visualization concept. This includes defining which objects to track, where to place tracking sensors, which tracking sensors to use, etc. It is essential to understand these requirements to select the best possible tracking strategy for the application. All tracking methods have their limitations in terms of speed, range, or resolution. These limitations must have minimal impact on the implementation as poor tracking may influence the performance and consequently affect acceptance of the application. Eventually, both tracking and visualization approaches must be compatible and the combined trade-offs must be acceptable to the user.

**Fig. 2.23.** Pictorial illustration of a cube, with origin at $O_{obj}$, placed in a reference space at defined as $O_{ref}$. Four points $P_{1..4}$ (in red) are defined at four corners of the cube. The Position of the points ($P_i$) can be defined in two coordinate systems i.e. $P_{ref}$ and $P_{obj}$. $Rt$ defines the relationship between corresponding points in two coordinate spaces $O_{obj}$ and $O_{ref}$ in terms of rotation and translation parameters (6-DOF pose). Minimum 3 corresponding points are required to compute the relationship between the two coordinate spaces. In this case, the pose $Rt$ can be computed from $P_{ref}$ and $P_{obj}$ using eq 2.1.

Tracking methods are developed to represent the three-dimensional information of the world (real or virtual) in mathematical terms. This information is used to compute two-dimensional images that are displayed to the user as visual stimuli. The position of objects in the real world keeps changing with respect to the user's viewpoint either due to the movement of the object itself or movement of the user. These changes are measured and reflected in the stimuli for a continuous experience. In the following text, we will discuss the details of tracking approaches and their advantages and disadvantages. 3D tracking is a well-studied topic and explaining all the details are beyond the scope of this thesis. We will focus more on concepts and methods that support understanding the research ideas presented in Part 2 and Part 3.

## 2.4.1  Position and Pose

Real world or virtual world, all XR applications are represented in a Cartesian coordinate space i.e. reference space. All tracking systems maintain location of each point or object of interest in the reference space. Mathematical representation of points and objects is done using position and pose respectively.

Position of a point is defined as 3D coordinates (x,y,z) along the 3 reference axis in the reference space (see Fig. 2.23 for illustration). This is true for any point in the space and therefore each point in space has 3 degrees of freedom. An object can be seen as a group of rigidly connected points. If the object is relocated in reference space, the 3D coordinates of each point on the object will change. However, the relative position between the points on the object always remains the same unless it is a deformable object. Similarly, a local coordinate system, i.e. object space, can be defined for each object where each point on the object always maintains a fixed position irrespective of how the object moves in the reference space. A pose is used to express the relationship between the same points defined in different coordinate

systems i.e. object space and reference space. A pose has 6 degrees of freedom, 3 for the position (x-y-z coordinates), and 3 for orientation (e.g. Euler angles). Pose parameters define how to transfer any point from the reference space to the object space and vice versa. The tracking methods are typically designed to compute the pose of objects in reference space. In this way, the pose is a simple and effective way to represent information about all points belonging to the same object.

### Mathematical representation

Let us define a reference space with the origin at $O_{ref}$ and place an object $A$ with origin at $O_{obj}$ in the same space. Each point $P$ on the object is defined as $P_{ref}$ in $O_{ref}$. The same point $P$ is represented as $P_{obj}$ in $O_{obj}$.

$$P_{ref} = \begin{bmatrix} X_{ref} \\ Y_{ref} \\ Z_{ref} \end{bmatrix}, P_{obj} = \begin{bmatrix} X_{obj} \\ Y_{obj} \\ Z_{obj} \end{bmatrix}$$

The pose relationship between them can be represented as $Rt_{obj}^{ref}$. It is mathematically represented with the following expression,

$$P_{ref} = Rt_{obj}^{ref} \cdot P_{obj} \tag{2.1}$$

$$P_{ref} = \begin{bmatrix} R_{11} & R_{12} & R_{13} \\ R_{21} & R_{22} & R_{23} \\ R_{31} & R_{32} & R_{33} \end{bmatrix} \cdot P_{obj} + \begin{bmatrix} T_1 \\ T_2 \\ T_3 \end{bmatrix} \tag{2.2}$$

### Application in XR

3 point correspondences are sufficient for the computation of a pose. As a general working principle, tracking sensors compute the pose of all objects by computing the position of some features, $P_{ref}$, in reference space. The position of those features in object space, $P_{obj}$, is used to compute object pose. Once the pose is determined the position of all points on the same object can be transferred to reference space. For example, the pose of the user's head is sufficient to infer the location of the eye because the eyes are rigidly connected to the head. Finally, the 3D representation of each object w.r.t. the user's perspective is created using the position of the eyes in the reference frame. This principle used in most tracking technologies for a simplified representation of the 3D world. The strategies for computing 6-DOF (Degrees Of Freedom) pose differ based on technology and sensory type (explained later).

In AR applications, the pose of real objects is computed in reference space, and virtual objects are place accordingly. In VR applications, virtual objects are maintained in the reference space and the pose of user's head is tracked in reference space for rendering perspective correct views. The 6-DOF pose of other objects such as controllers, tools or displays are also tracked in the same space. In summary, pose and position parameters allow mathematical representation of the 3 dimensional world that is used to compute interactive visual stimuli.

## 2.4.2 Image Projections

The visual stimuli displayed to the user are in the form of two-dimensional images. These images are rendered by projecting the 3D representation of the world to a 2D image plane. Mathematical concepts for computing image projections from the 3D world are defined in the field of projective geometry [85]. In the field of metrology [127], concepts of image formation and image projection are also used for extracting 3D information from 2D images. In principle, human eyes perform the same function as cameras where each eye captures an image of the real world that is further interpreted in the brain. For XR applications, camera models provide a mathematical foundation for generating visual stimuli.

A pin hole camera model is used to convert 3D coordinates into pixels (eq. 2.4). Each pixel, $p$, in the image represents a point $P$ projected from the 3D world. $P$ is represented in the coordinate system of the camera with origin as $O_{cam}$. A 3x3 matrix (upper-triangular), $K$, is used to transfer 3D space into pixel space (see eq. 2.4). The dimension Z (depth) is lost after this transformation. $K$ is known as the intrinsic matrix, it has 4 non-zero parameters which define the properties of the camera i.e. focal length, sensor size etc. These parameters are determined through a process of camera calibration (see sec. 2.4.6). The size of the image is limited by the sensor size and field of view of the camera.

In XR applications, the views are generated from user's perspective when HMDs or CAVE systems are used. In such cases, the position of the human eye is treated as the origin of a virtual camera. Tracking information is used to define the pose of objects in virtual camera space. The display surface coincides with the image plane and 3D information is projected onto this plane using the camera calibration parameters. In AR applications, images are often projected using the perspective of the camera used to capture real world view. In such cases, the cameras are calibrated separately for accurate registration of virtual stimulus with real world projection.

$$p = K \cdot P \tag{2.3}$$

$$\begin{bmatrix} x/z \\ y/z \\ 1 \end{bmatrix} = \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} \tag{2.4}$$

## 2.4.3 Tracking configuration

Tracking configuration refers to the arrangement of the sensors for 3D tracking. We present two types configurations for sensor placement: $a.$ inside-out and $b.$ outside-in. Sensor placement is decided based on the features of the application e.g. mobility, multi-user tracking or multi-object tracking etc. The selection of sensor type and configuration are crucial for designing closed loop systems. Sensory modalities and their features are explained later in this section.

**Fig. 2.24.** Pictorial examples of inside-out tracking configurations with optical cameras. This configuration is selected for *on the move* applications. It is suitable for indoor and outdoor environments. As depicted in the images, the tracking range is often limited. (left) Wearable sensor: The sensor is mounted on the body of the user, a typical use case for HMD based tracking. (center) Handheld sensor: The sensor is operated by the user with mobile device such as a smart phone or tablet. (right) Semi-mobile sensors: The sensors are mounted on a mobile setup. A separate view is also presented for better understanding of sensor arrangement.

### Inside-out configuration

In this configuration, the sensor is placed on an object moving in space and it measures sensor pose w.r.t the reference coordinate system. The sensor is carried by the user as a wearable device such as HMD, or it is operated by the user as a mobile tracking device e.g. smartphone. The inside-out approach is compatible with all design concepts for AR applications i.e. HMD and spatial displays. However, for VR applications it is mainly used for HMD based applications. In Figure 2.24 few examples of the inside-out configurations are illustrated. The final objective is to compute the pose of objects w.r.t sensor space ($O_{sensor}$) and the reference space $O_{ref}$. The tracking approach can be chosen in two ways i.e. temporal referencing and direct referencing.

**Temporal Referencing:** This strategy uses temporal sampling of sensor movement for computing it's 6DOF pose in reference space. As a first step, the sensor pose is initialized in the reference frame as ($Rt_{t0}$) while starting the application. After initialization, the sensor only measures relative change in pose w.r.t it's previous location i.e. $Rt_{t0}^{t1}$. This measurement is combined with the initial pose to compute the relation with reference frame i.e. $Rt_{t1} = Rt_{t0} \cdot Rt_{t0}^{t1}$. The initial position may serve as reference frame itself, in that case initial pose of sensor would be identity i.e. $O_{ref}$ is same as $O_{sensor}$ and $Rt_{t0} = I$.

This strategy is typically used for VR applications where sensors measure the movement of the head and controllers in the virtual environment. In this case, it is easier to define the initial starting pose in the virtual world. This approach is useful for AR applications when augmentations are not assigned to specific real-world objects. For example, rendering virtual objects like dinosaurs or fish in a living room. For such cases, the temporal referencing is used to create an ad hoc 3D representation of the reference space (i.e. the living room in our example). The sensor pose is computed in this space and this allows virtual elements to be registered to the scene.

Both inertial sensors and camera-based sensors are capable of using this technique. It is known as a dead reckoning problem in literature of inertial sensor-based tracking [40]. Optical sensor-

Marker Registration = Rt$_{cube \leftrightarrow Marker}$
Calibration = Rt$_{sensor \leftrightarrow ref}$

**Fig. 2.25.** Pictorial illustration of marker tracking based pose computation approach with optical sensors. This concept is used in both inside-out or outside-in tracking approach. Markers are used to compute pose of target object. The picture shows an example marker pattern which is designed to be unique and easily detectable in the image. The marker is placed on the object and it's position is registered in the local space of the object $O_{obj}$. The pattern location is not changed after registration. The position of the sensor is registered in the reference space using a calibration procedure. This relationship is stable until the camera position is moved. Fixed camera based tracking is used in outside-in approaches. For inside-out approaches, the reference space is aligned with sensor space i.e. $O_{sensor} = O_{ref}$ or object space i.e. $O_{obj} = O_{ref}$. The pose between sensor and object is unknown quantity which is computed in real-time. All other relationships are determined using registration and calibration information.

based methods often use static features in the scene ( edges, corners, etc.) and use them to infer the 6-DOF pose of the camera. This approach is useful in unknown environments and provides mobility. This idea suffers from the problem of drifting because each measurement has a small error and with time the error is accumulated. Sensor measurement is limited to its sampling rate and therefore recording pose of very fast movements is challenging, e.g. controller motion in VR gaming scenarios. Optical sensors may fail due to fast motion or get disoriented if the environment has symmetrical patterns. This approach is rather popular in both AR and VR applications, especially with HMDs. It is possible to use multiple sensors independently or to fuse information for optimized performance. Sensor fusion is computationally expensive and increases processing time. This approach is one of the most popular approaches with modern VR and AR applications, especially with compact and inexpensive optical sensors, for example, HTC Vive or Microsoft Hololens.

**Direct referencing:** In this approach, the sensor pose w.r.t reference space is computed using referencing objects placed in the scene. Referencing objects are registered in the common reference space and tracking sensors compute pose of reference objects to infer it's own location in the reference space. Therefore, the accuracy of the placement of referencing objects affects the overall accuracy of the application. To avoid this error, reference space is often defined by the local coordinate space of the reference object. The objects used for referencing are of different types depending on the sensor technology. For example, optical sensors use known 2D or 3D geometric patterns e.g. fiducial markers or patterns (see Fig. 2.25).

**Fig. 2.26.** Pictorial example of outside-in tracking configurations with optical cameras. The illustration displays multiple objects being tracking in reference space using multiple sensors. The sensors define a specific tracking area where object pose came be tracked. All sensors are registered in reference space using a one-time calibration approach. The sensor position does not move and therefore the pose of sensors in reference space remains static. In a typical application, pose moving objects are tracking in sensor space and transferred to reference space. This information is further used to generate desired stimulations. The illustration shows that user's perspective is tracked along with display and cube.

Direct referencing is used in AR applications when augmentations are designed for specific target objects (with known geometry or patterns) such as industrial AR applications. It is possible to use multiple target objects for extending the range of application for example in museums or factories. The local coordinate space of objects is used as a reference space to handle moving objects. This way augmentations always remain registered to the reference objects irrespective of movement. For VR applications, this approach is used by strategically placing multiple referencing objects in the real world to compute the orientation of sensors in the VR space. One such example is HTC Vive headsets which use lighthouse technology with base stations. They are placed in the corners of a room to define reference space and sensors on the HMD use detect the lasers emitted from the base stations to recover pose HMDs in the reference space. The range of tracking with base stations is limited and but the accuracy is very high.

Inside-out approaches definitely provide more mobility to the user. Mobility comes at expense the of wearing additional sensors and carrying computational devices. However, sensors are becoming cheaper and approaches supporting both indoor and outdoor applications are preferred by users. The performance of inside-out tracking is boosted by sensor fusion approaches [29]. Such applications require a calibration process to synchronize the sensor measurements. It is also possible to combine temporal and direct referencing concepts to provide more flexible tracking results [29]. Inside-out methods are more suitable for personalized experience as sensors are typically operated by the user. Shared experiences are possible in known areas where referencing is rather easier. The general trend of AR/VR applications is moving towards use of inside-out approach [21].

Outside-in configurations are used when sensors are placed at a stationary location to track moving objects in the scene such as HMD, displays, etc. This tracking strategy and range are defined by the sensor type e.g. magnetic, optical, mechanical, etc. Outside-in configurations are compatible with most design concepts of AR and VR applications. It is particularly useful when perspective tracking is important for displaying the stimuli e.g. Spatial AR, CAVE. HMD based methods may or may not use this configuration depending on the requirements.

Unlike the inside-out approaches, the reference coordinate system is defined by the sensors. The configuration almost always requires multiple sensors to be used and thus the calibration process is performed to register measurements of all sensors in the same space i.e. $Rt_{sensor}^{ref}$. The pose of moving objects is computed by one or more sensors in sensor space and transferred to reference space using calibration information i.e. $Rt_{obj}^{ref} = Rt_{sensor}^{ref} \cdot Rt_{obj}^{sensor}$. Ultimately, tracking information is used to create visual stimuli.

The tracking method in outside-in configurations can be contact-based on contact less. Some methods require physical contact e.g. mechanical tracking with robotic arms (Fig. 2.15). Physical constraints restrict the movement and therefore contact less methods, such as magnetic or optical tracking, are preferred when mobility is an important criterion. Contact less approaches require special sensors or markers to track the object of interest. The concept of marker based tracking is illustrated in figure 2.25. Sensors tracks the markers and that is used to infer the pose of the object. There are two types of markers i.e. active and passive. Active markers communicate back with the tracking sensors with a signal such as light emitters or acoustic sound emitters. Passive markers are placed on the object and their position is detected by the scanning within the tracking range, for example visually scanned AR Tags or RFID tags used in the super markers. Active markers are expensive and have power requirements whereas passive markers are cheaper to produce. Both may suffer from issues like a line of sight or signal interference.

In most XR applications, optical sensors are preferred for pose computation over other type of sensors. Outside-in configuration is ideal for scenarios where multiple objects tracking is required. Motion capture systems are popular for VR approaches where gesture tracking is important for interaction. Outside-in setups are usually placed at a fixed location and computational support for fast processing and tracking is possible for highly dynamic XR applications. Outside-in configurations with large network of sensors is less mobile and accessible for user's only at the installation sites. This is suitable for facilities where fixed installations are not hindrance of installation e.g. museums, class rooms etc. It is possible to design mobile setups for outside-in configuration using a limited number sensors e.g. VR base stations. This approach is largely limited to indoor environments and therefore less preferred for modern AR applications.

## 2.4.4  Sensor technologies

The types of sensors mainly differ in terms of the information that they acquire from the scene i.e. optical, acoustic. The selection of sensors depends on the need of the applications. Most XR applications require 6-DOF pose computation of one or more objects. Not all sensors can

Video Image

Depth Image

Infrared Image

**Fig. 2.27.** Image of a dog obtained from Intel Real-sense D-435 sensor in an outdoor environment. The sensor provides three type of images i.e. video, infrared and depth image. Video image has color information stored as an RGB value for each pixel. Infrared images capture infrared reflection from scene and they are represented as gray scale values. Each pixel in depth image represents depth value of the imaged point. Depth information is displayed as color coded image to give a sense of depth i.e. blue to red.

. ©Hemal Naik.

measure the 6-DOF pose of objects but it is common to combine multiple tracking technologies to achieve the best possible results. The sensor measures different types of information e.g. microphones record audio information, cameras record visual information. This raw information collected via sensors is further processed to extract the 6-DOF pose of target objects. Over the last few years, different types of sensors have been demonstrated for XR applications. Optical sensors are one of the most commonly used sensors in modern XR applications. We have chosen to discuss optical sensors in detail because our research mainly focuses on tracking approaches with optical sensors.

### Optical Sensors

Optical sensors represent a family of sensors that capture light signals reflected from the scene in the form of images. The core idea is to measure visual signals from the scene and filter the information required for computing object pose. A wide variety of optical sensors are designed to capture selective information from the scene. For example, video cameras capture visible light, infrared cameras capture only IR signals. We will only discuss commonly used as optical sensors.

The raw information obtained by the optical sensor is represented as two-dimensional images. Earlier we explained that images formed by projecting 3D information from sensor space to a 2D image plane using different camera models (sec 2.4.2). Similar mathematical techniques are used to recover the 3D information in metric units in the sensor coordinate system [127]. These techniques depend on sensor configuration, the number of sensors used, and the type of sensors used in the application. First, we will describe different types of sensors and the information that they gather and then we will shortly discuss commonly used approaches for

3d measurements. Figure 2.27 shows different types images obtained from different sensors used in XR applications.

**Infrared sensors** measure the infrared light reflected from the scene. This is one of the most commonly used sensors for indoor XR applications e.g. Outside-in configuration with Vicon motion capture [8] or Mobile inside-out configuration with HTC Vive lighthouse [9]. Artificial light sources commonly used in indoor environments have negligible infrared frequencies and do not interfere with sensor's IR light source. This fact is used as an advantage for designing tracking approaches with IR sensors. Active markers that emit IR light or passive markers that reflect IR light are attached to the object of interest. In the case of passive markers, IR light sources are combined with cameras to illuminate the environment with IR light. Generally, markers appear as bright oversaturated spots the images captured by IR sensors. Pixel locations of these markers from one or more images are used to determine the 3D position of markers and from that pose of the object is computed [127]. This is a general working principle with most IR based methods. However, special techniques also exist where IR lights are used differently to obtain 3D information e.g. steamVR lighthouse or depth cameras. Lighhouse technology uses IR sensors to emit phase-shifted infrared lights which are detected by sensors on the headset. The sensors on the headset can compute pose based on the difference between the signals emitted from each base station. The choice of using active or passive markers depends on parameters such as cost, size, and weight of markers. This approach is not suitable for outdoor approaches as natural light contains significant IR components (see Fig. 2.27) and it is reflected from all sources present in the scene, which makes spotting objects of interest difficult.

IR tracking with markers is considered a very reliable and accurate method for object tracking. The complexity of image processing and computer vision operations is low. This makes the method extremely responsive which is ideal for real time applications. Marker based methods are limited to tracking objects with markers. However, good detection is guaranteed which in turn provides the highest form of accuracy and therefore such methods are preferred in industrial settings [127]. There sensors are used for both AR and VR applications, especially when fixed number of objects are being tracking for highly interactive XR applications.

**Video sensors** measure the visible light reflected from the scene. Each pixel of the image is represented by a color value in RGB format. Video cameras are used in both inside-out and outside-in configurations. For XR methods, both single and multiple camera methods are used to compute the pose of the target object. The scope of obtaining 3D information of all objects in visual field is very high. All information in line of sight and withing the FOV is captured. Therefore, it is possible to reconstruct almost all the visible points in 3D space [85, 127].

Video cameras based tracking solutions have wide range of variations [85, 127]. These problems are studied widely in computer vision and close range photogrammetry. Single or multiple sensors are used to compute 3D informations of the scene. Typically, markers are used to save time and compute pose of objects faster for AR/VR applications. However, markerless techniques are currently being developed for many 2D or 3D tracking applications.

---

[8] https://www.vicon.com/
[9] https://www.vive.com/eu/accessory/base-station/

Video camera-based tracking is preferred in many mobile AR applications (all design concepts). Implementation wise video cameras are widely available, cost-effective, and low maintenance sensors. Earlier we discussed the advantages of offering smartphone-based XR solutions. For applications with a requirement of high accuracy, multi-camera based methods are more suitable. Multiple cameras need more processing time and additional preparation i.e. calibration, synchronization. They may not be able to support real-time applications in highly dynamic environments. Modern machine learning methods have made real-time object detection, posture computation, and 3d reconstruction significantly faster for both single and multi-camera approaches [5, 157]. In summary, video sensors are useful for mobile and lightweight indoor or outdoor application of XR.

**Depth sensors** measure 3D information directly from the scene and provide depth images where 3D depth values are encoded with 2D pixels (Fig. 2.27). Depth sensors are often designed using multiple optical sensors together and their working principle is depends on the sensors.

Time of Flight sensors used temporal difference between sending and receiving time of the signal to estimate distances of objects in the scene like ultrasound or LiDAR technology. Another type of depth sensors use IR projectors to project known patterns and synchronized IR cameras that capture image of these patterns. The sensors are temporally synchronized and 3d geometry of the scene is computed by observing deformation in the projected patterns. Stereo triangulation (explained in next section) is used to compute 3D features from identical features detected in the two images. The known patterns act as unique features that are easy to identify and detect in the images. These depth sensors are also known as RGB-D sensors and they come as a single sensor with predefined calibration, synchronization and inbuilt processors for fast data fusion.

The same principle is used with a projector-camera (Pro-cam) setup for computing detailed 3D reconstruction of the scene. It is commonly used to compute detailed 3D reconstructions with contact less approach for industrial scenarios e.g. structured light, whitelight [127]. Structured light is commonly used approach where a sequence of striped pattern is projected on the scene. This allows spatial and temporal encoding which is recovered from corresponding images. This approach is commonly used in metrology [127] for detailed 3D reconstructions with very high accuracy.g. GOM[10]. A faster approach is to project a single image with prominent features [177]. This approach allows fast but sparse reconstruction because only distinct features can be robustly detected and matched from the captured images.

The performance of these depth sensors vary in terms of accuracy and time. RGB-D cameras are becoming popular for all XR applications with both HMD [129] and spatial display based designs [30]. RGB-D sensors also have an additional video camera that is calibrated to the infrared sensors. The overall information obtained is 2D color images, infrared images, and a depth image where depth value is mapped to each pixel of a video camera as displayed in Fig. 2.27. Pose computation using such information can be much faster since 3d features can be used in addition to 2d image-based methods. Commercially available RGB-D are less accurate but provide real-time results e.g. Microsoft Kinect or Intel RealSense. Depth sensors have become popular for real-time gesture tracking and offer natural interactivity for XR

---

[10]https://www.gom.com/3d-software/gom-inspect.html

applications. Research in this direction is pushed by the gaming industry because these devices support human posture tracking [189]. It is possible to manufacture compact depth sensors which are an advantage for HMD and smartphone-based applications. It must be noted that such advantages come at the cost of resolution and accuracy. Industrial pro-cam sensors are accurate (up to mm) but not suitable when extremely fast response time is required. They are used with AR scenarios in manufacturing or projection mapping when target objects are static and or slow moving [177]. Moreover, the accuracy of depth measurements decreased at larger distances. Similarly, infrared-based methods fail to provide accurate representation in outdoor environments.

**Summary:** Overall optical sensors are favored in XR application because they are less intrusive and provide a contactless method for object tracking. The development of machine learning methods is pushing the field in the direction of markerless tracking. Line of sight is a common problem with optical sensors. The technical specification of optical sensors affects the overall performance of the tracking application. The tracking range is limited by the field of view and sampling rate (frame rate) of sensors determines the limitation in terms of tracking moving objects. Sensor size and lens quality are affected by image quality in terms of resolution and clarity. Resolution and capture rate has an inversely proportional relationship which introduces a trade-off in terms of accuracy and real-time performance. The selection of optical sensors requires a good understanding of application requirements because sensors ultimately define many limitations of the system in terms of performance and experience. Their discussions will appear throughout the thesis, especially in part 3 where we present a new setup.

## 2.4.5 Tracking strategies with optical sensors

Optical sensors are a widely used tracking modality for XR applications. The sensor type and configuration are selected based on the application requirements. Different computer vision techniques are used in different scenarios e.g. static camera or moving camera etc. Tracking strategies are often specialized and chosen to be compatible with the application at hand. The implementation varies based on the number of sensors used and each have their advantages and limitations. We will explain some approaches in generic form based on the type of sensor and number of sensors.

### Monocular tracking

Monocular tracking methods use a single camera for tracking the 3D pose of the desired objects. A camera image is formed by projecting the 3D scene in front of the camera to a 2D image plane. The depth information is lost in this process. Recovering the depth of a single point is not possible from a single view. However, recovering pose of a 3D object or a 2D pattern is possible using several mathematical techniques that are developed in the field of computer vision. As a general workflow, first features of interest are detected in the images and then 3D pose of the object is determined [85].

It is always important to consider the tracking needs of the application. Some AR applications only require 2D tracking and do not need 3D pose computation. Consider a simple AR application of face detection where the objective is to draw a circle around the faces that appear in the image. It is sufficient to detect the faces in 2D space using image processing

**Fig. 2.28.** Image shows 2D and 3D marker patterns used by ZapBox for mixed reality applications.

draw circles around the region of interest. Tracking 3D orientation of a face is advantageous but not necessary. The augmentation changes if the location of face moves in the image but not if orientation changes. This approach is very fast and suitable for simple real-time applications where the 2D location of the target object in enough for augmentation.

Depth estimation using a single image is a challenging problem. The complexity of the approach depends on the geometry of the tracking object. Tracking the 6-DOF pose of planar objects is easier compared to 3D objects because all points belong to the same plane. Therefore, the geometric invariants can be extracted from the image features and the problem is reduced to computing 3D orientation of the plane ([85]). The easiest method for monocular 6-DOF tracking is to use custom-designed marker patterns that are easier to detect in-camera images (see Fig. 2.28). The position of these features ($P_{obj}$) must be known in the local coordinate system of the object. Once these features are identified in the image space ($p_{img}$), the pose of an object can be computed in a reliable manner using the following equation.

$$p_{img} = k \cdot Rt \cdot P_{obj} \tag{2.5}$$

where $Rt$ is a pose of the object and $k$ is the calibration matrix. 4 point correspondences are sufficient to compute a unique pose [124]. Markers patterns can be 2D or 3D (Fig. 2.28), generally 2D markers are preferred as they are easy to produce and measure. 3D markers need to be precise to compute an accurate pose. If markers are attached to the objects, their position is registered in the object coordinate system (explained with Fig. 2.25). This holds true for all marker based camera tracking approaches irrespective of the tracking strategy. Many AR applications require the users to use a predefined pattern which serves as a marker and does not require the user to perform any registration. Markerless methods exist, however reliably

finding features can be difficult due to various factors e.g. low light, occlusion. Machine learning based methods have made remarkable progress in this particular direction [122, 134, 138, 157, 226]. Convolutional Neural Networks (CNN) are used to predict the pose of an object directly from images. Marker less tracking is considered extremely useful for developing *plug and play* XR applications [36, 138, 189].

Monocular approaches are most suitable for inside-out configurations with either temporal referencing or direct referencing techniques (see 2.4.3). Marker-based methods explained earlier are used for direct referencing. Temporal referencing techniques are often used in combination with inertial sensors for additional accuracy. For indoor applications, the camera is moved slowly to scan reliable features (corners, floors, etc.) in the room and a 3D feature map is created in-camera coordinate system. These approaches are called SLAM (Simultaneous localization and Mapping) [109, 146] where feature maps are created and simultaneously used to track the motion of moving the camera in the space. Such approaches are useful for AR applications in unknown environments.

Monocular tracking methods are suitable for mobile XR applications. Good quality high-resolution cameras are compact and relatively inexpensive to produce. AR applications with smartphones are one of the most popular use cases for monocular tracking based XR applications. The tracking range is limited to the field of view of the cameras. Limited computation power of mobile devices can introduce lag in the application and to circumvent this problem inertial sensors are used in combination with this approach. Markerless methods with supervised machine learning techniques are a promising step in this direction. Most importantly, single camera-based applications require one-time calibration procedure to compute intrinsic parameters and lens distortion.

### Stereo-Multi camera tracking

Stereo and Multi-camera tracking strategy involve pose computation by fusing data from multiple sensors. Stereo approaches with two cameras are preferred for inside-out tracking and multi-sensor installations are commonly used for outside-in sensor configurations. The cameras are placed in such a way that they have an overlapping field of view. This way objects moving in the area are seen and detected by both cameras. Calibration is performed to compute the camera positions relative to each other and define a reference space. The position of cameras is not changed during the application.

Determining the pose and orientation of objects using the 3D data is easier than with the mocular strategy. 2D features of the target object are detected in each image and then triangulated to obtain 3D features via triangulation [85]. 3D features are directly computed in common reference space from 2D information. Further, the pose is computed by using the 3D correspondence between object space and reference space[92]. The accuracy of triangulation depends on image processing and calibration accuracy.

Optical systems for XR applications consist of both IR and video sensor based methods. Multi-camera motion capture systems prefer infrared cameras as they can perform at very high speed with very high accuracy rates e.g. VICON[11]. Infrared sensors with retro reflective

---

[11]https://www.vicon.com/

markers require less time for image processing and pose computation (discussed earlier) [127]. Reliable commercial solutions available in marker and used for many XR applications where speed and accuracy is paramount e.g. medical, industry. We have discussed this approach for tracking multiple animals in large environments (discussed in part 3). However, such methods always require markers and therefore video cameras are used as alternative. Video camera based multi-camera applications are becoming popular with both marker and markerless approach. Processing time for video images is a huge bottleneck especially in case of the markerless methods. Real-time performance for dynamic applications may not be possible. However, multi-camera tracking with video cameras is possible for semi-static environments where the position of objects does not change rapidly. We have used such a stereo system for developing industrial AR solutions (discussed in part 2).

Multiple cameras provide larger tracking range which may be crucial for applications with multiple users e.g. CAVE, multi-person VR. Larger area also allows more mobility for the users. Line of sight problem can be reduced by strategic placement of sensors, however occlusion remains a problem. Technical challenges of using optical sensors are already mentioned in monocular approaches. Multi-camera approaches have several additional technical limitations in terms of installation. Cameras must be synchronized i.e. capture images at the same time. Unsynchronized data capturing may lead to error in triangulation of features and as consequence wrong depth 3D point computation. Camera calibration process must be performed every time the camera position is changed. Calibration routine is well established and relatively easier to perform (see sec. 2.4.6).

### Depth tracking

This strategy involves object tracking based on a 3D point cloud. It is used mostly for inside-out configurations with static or moving sensor. Generally, depth sensors or projector-camera setups are used to obtain 3D point cloud for depth based tracking.

The working principle of depth sensors is discussed earlier in section 2.4.4. The outcome of each frame is a sparse or dense 3D point cloud based on projected patterns. Unlike stereo methods, 3D information from the whole field of view is available as a point cloud. Therefore, filtering and matching techniques are used to identify the 3D points belonging to the object of interest. After that the 3D point cloud is compared with known geometry of 3D object and the pose is computed.

Pose computation from 3D reconstruction is rather a challenging task. Classic 3D feature matching technique [62] were used but now most are replaced by machine learning techniques with arrival of depth sensors [189]. The accuracy of the pose depends on the quality of the point cloud. Pro-cam systems provide high accuracy because resolution of point cloud is high. Nowadays, it is possible to design customized pro-cam setups to perform 3D reconstructions in real-time. using projectors and cameras with higher operation frequency [143, 152]. These setups turn out to be extremely expensive and cumbersome to set up in terms of technical requirements. As a result, the solutions requires longer time for to converge. RGB-D sensor provide a suitable compromise when real-time performance is necessary [30]. They are a good compromise in terms of accuracy, speed, size and cost. For VR applications, depth sensors are used for maintaining orientation of the headset. For indoor environments with unknown geometry static features of the room are used with temporal referencing approach.

Spatial AR applications with projectors are a special case. In this approach, projector is used for augmentation and therefore it is possible to use the same projector for pose computation (discussed in sec. 2.4.4). The spatial augmentation is performed using the same pose and the registration remains valid as long as the projector or the target object do not move. Recently, it is also shown that augmented patterns themselves can be used for real-time referencing [177].

## 2.4.6 Calibration

Calibration is the first step required for setting up any equipment for optical tracking. A calibration process a method to compute the parameters that establish a relationship between points in 3D space and the corresponding 2D projections. In the case of multiple sensors, calibration also includes computing relationships between coordinate systems of different sensors. These can be multi-camera systems or camera-projector systems. The calibration of multiple sensors is required to combine their information in a common reference space. Calibration process and relevant mathematical modeling is extensively studied in the field of close range photogrammetry (metrology) [127] and computer vision [85]. We will only provide a brief overview of the topic. There are two types of calibration: intrinsic and extrinsic. Intrinsic parameters are sensor-specific and extrinsic parameters establish the relationship between the camera space and the reference space.

Intrinsic parameters are used to project 3D points from camera space ($O_{cam}$) to 2D pixel space (eq. 2.4). A commonly used mathematical model for both cameras and projector is the pin hole model [85] (sec. 2.4.5). Intrinsic parameters depend on the hardware i.e. sensor type, size, etc. In practice, while taking an image the rays of light pass through the lens before they are projected on to the sensor. Due to the physical property of the lens, the light gets distorted which impacts the world to image (3D to 2D) mapping. Distortion effects depend on the physical shape of the lens e.g. wide-angle lens or fisheye lens. The distortion effect must be removed from images to obtain the direct world to image mapping. Lens distortion parameters (e.g. radial, tangential) are fixed and also computed with intrinsic parameters during the calibration process. In summary, 4 parameters of camera intrinsics ($k$) and 5 parameters (3 radial and 2 tangential) of lens distortion are computed. It should be noted that the distortion model must be selected based on the type of lens. For the highest form of accuracy, all parameters must be considered, refer work of Luhmann et al. [127]. Intrinsic parameters for the projectors are the same as cameras in the case of video projectors. In the case of laser projectors, a different process is followed since the rotating mirrors principle is used for projection. Generally, such a process is performed by manufacturers and need not be computed specifically. More details on intrinsic calibration of laser calibration can be obtained from relevant literature [64, 114, 221].

For single camera calibration, extrinsic calibration is not performed. Extrinsic calibration for the stereo camera system is expressed in the 6-DOF relationship between the coordinate system of two cameras i.e. $Rt_{cam2}^{cam1}$. For projector-camera systems, extrinsic computation involves finding the 6-DOF pose of the projector w.r.t camera i.e. $Rt_{proj}^{cam1}$. In summary, 6 extrinsic parameters have to be computed for each additional sensor.

**Single camera calibration** involves taking images of a known calibration object from multiple perspectives. The calibration object can be a 2D pattern (e.g. a checker board) [209] or 3D geometry [178]. The position of each feature ($P_{obj}^i$) in 3D space, ($O_{obj}$), of the calibration object is known. The corresponding 2D location of each features ($p_{img}^i$) is detected in the images. Mathematical equation system is created using object points $P_{obj}^i$ and image projection $p_{img}^i$.

$$
\begin{aligned}
P_{cam}^i &= Rt_{obj}^{cam} \cdot P_{obj}^i \\
p_{img}^i &= K \cdot P_{cam}^i \\
p_{img}^i &= K \cdot Rt_{obj}^{cam} \cdot P_{obj}^i
\end{aligned}
\tag{2.6}
$$

Where $Rt_{obj}^{cam}$ are known as extrinsic parameters which transfer points $P_obj$ from the object space to points $P_cam$ in the camera space . The image points mentioned above, i.e. $p_{img}^i$, are without considering the lens distortion. Lens distortion parameters are independent of camera sensor parameters. For most computer vision applications images captured by the camera are undistorted and then feature detection is performed. The equations system showing the relationship between distorted and undistorted points is shown below,

$$
\begin{aligned}
x_{corrected} &= x(1 + k_1 r^2 + k_2 r^4 + k_3 r^6) \\
y_{corrected} &= y(1 + k_1 r^2 + k_2 r^4 + k_3 r^6)
\end{aligned}
\tag{2.7}
$$

$$
\begin{aligned}
x_{corrected} &= x + [2p_1 xy + p_2(r^2 + 2x^2)] \\
y_{corrected} &= y + [p_1(r^2 + 2y^2) + 2p_2 xy]
\end{aligned}
\tag{2.8}
$$

where, $k_1, k_2$ and $k_3$ are parameters for radial distortion and $p_1, p_2$ are parameters for tangential distortion.

The equation system presented in eq. 2.6 holds true for all 2D-3D point correspondences. If 2D calibration objects are used the Z component is always 0 and therefore each 2D-3D correspondence yields only 2 equations which are not the case with 3D calibration objects. For $n$ point correspondences in each image $n * 2$ or $n * 3$ equations are used to solve 6 parameters for pose ($Rt_{obj}^{cam}$), 4 for $k$ and 5 for distortion. Theoretically, 8 (for 2D) and 5 (for 3D) correspondences from a single image are enough for solving all unknown parameters. To achieve mathematical stability the larger equation system is created using many point correspondences. For optimal lens distortion computation, it is recommended to capture images from different perspectives such that 2D features are well distributed across the image. This is done to compute unbiased distortion parameters since the effect of distortion is uneven across the image (more in corners). In general, 3D calibration objects can provide the most accurate calibration results if calibration objects are produced with very high accuracy. However, 2D calibration objects are preferred because producing making very high accurate 3D objects is expensive.

**Multi-sensor calibration** the same method is used but the equations system is much larger. The process involves taking multiple images where both the images have visible calibration patterns. The process for computation of intrinsic parameters remains the same and the

extrinsic parameters are computed using the pose of the calibration object. The extrinsic parameters for stereo calibration i.e. $Rt^{cam1}_{cam2}$ are computed using the following equation system.

$$P^i_{cam1} = Rt^{cam1}_{obj} \cdot P^i_{obj}$$
$$P^i_{cam2} = Rt^{cam2}_{obj} \cdot P^i_{obj}$$
$$P^i_{cam1} = Rt^{cam1}_{obj} \cdot Rt^{obj}_{cam2} \cdot P^i_{cam2} \tag{2.9}$$
$$Rt^{cam1}_{cam2} = Rt^{cam1}_{obj} \cdot Rt^{obj}_{cam2} \tag{2.10}$$

This relationship is important for triangulating the features detected in the two camera images [85]. Projector-camera calibration is performed slightly differently but the equation system remains the same. A sequence of known patterns is projected by the projector and the features are detected by the camera. The corresponding points are used to create an equation system where 3D points are also added as unknown parameters. To reduce complexity the projections are made on planar surfaces.

The calibration process is extremely important because this determines the overall accuracy of the system. The tracking results of the stereo or multi-camera systems depend a lot on the accuracy of 3D triangulation which in turn depends on the accuracy of image processing and the extrinsic calibration parameters. Error in calibration may cause pose computation and augmentation errors, especially for spatial AR systems. It is common to use planar board with circular markers and iterative refinement techniques for parameter optimization [48, 77, 208].

# Part II

Augmented Reality Applications For
Industrial Manufacturing

# AR for manufacturing assistance <span style="float:right">3</span>

> *When a better tool or idea or approach comes along, what could be better than to swap it for your old, less useful tool?*
>
> — **Charlie Munger**

Industrial methods and practices are highly dependent on the availability of the right set of tools. For decades, Augmented Reality is being touted as the next best tool for visualization of data. Industrial researchers have been one of the strongest support groups for the development of AR tools for various applications [70, 153] such as manufacturing, quality check, training, rapid prototyping, etc. The industrial community has backed the development of AR solutions by investing time, effort, and money. There is a major incentive in realizing AR as an interactive visualization tool because it has shown great potential to save time, effort, and money in the future. Moreover, interactive 3D visualizations can be used to create a completely new type of workflow and replace classical approaches of displaying information i.e. 2D screen, paper plans. It is firmly believed that tools equipped with AR will be used in the industrial practices of the future i.e. Industry 4.0 [184].

Before going into the details of our research, I would like to highlight that research presented in this part was part of ARVIDA initiative [184]. ARVIDA project was supported by the German Federal Ministry of Education and Research (grant no. 01IM13001N) to promote the development of virtual techniques for industrial scenarios. This project builds upon a successful series of projects from the past ( ARVIKA, ARTESAS, AVILUS). As we mentioned earlier, such projects are supported to promote innovation and research by forging collaboration between industrial partners and academic institutions within Germany. The author worked with Extend3D GmbH as a researcher and interacted with other industrial partners (Volkswagen group, Thyssenkrupp, etc.) during the course of the project. [1] The readers must note that information and state of the art research provided in this thesis are in the context of the mentioned time period. There have been several developments in the field after the conclusion of the project which is not covered as part of this thesis. This being said, we think that the discussions presented in this chapter are very much relevant for future IAR applications. Especially because the solutions presented in the next chapter are successfully implemented in the commercially available IAR products of Extend3D GmbH [67].

Even after decades of effort Industrial AR (IAR) solutions are not substantially integrated in modern industrial practices [70, 153]. The popularity of AR is rising in other domains such as entertainment, education, and tourism (see 2.2). There are many reasons for not having the desired effect in the industrial domain. In this chapter, we will argue that technological

---

[1]The project was conducted from 2014 to 2017.

superiority is not enough and the implementation method is also crucial for successful acceptance of AR solutions. IAR applications cannot be designed without the contribution of all stakeholders i.e. workers, decision-makers, workflow planners and technology experts. In the next section, we will elaborate on this point, and in the following sections we will outline the concrete research objective undertaken by the author and collaborators.

## 3.1 Requirements from IAR applications

We interacted closely with our industrial partners to understand the process of adopting new work practices. We learned that an alternative solution is always welcome, but before developing any type of solution it is vital that we understand the acceptance criteria of the industrial environment. All new manufacturing practices (or workflows), methods or tools, are evaluated rigorously before being adopted into regular work practices. These evaluations are based on performance criteria such as accuracy, reliability, repeatability, etc. and usability criteria like scalability, cost-effectiveness, and user-friendliness. New methods or tools are accepted only if they have an advantage over existing methods, otherwise they are rejected irrespective of the elegance of approach or the technology being used. While these criteria remain true for almost all types of practices, we will discuss these more in the context of AR applications.

IAR applications, like any other technology, have to go through the stringent evaluation before being accepted. IAR applications have been around for long but haven't made a big impact in the industry because many applications failed to pass the evaluation criterion for acceptance. We draw this conclusion using the support of arguments from two research articles specifically written to address the requirements from IAR applications. In early 2000s, Navab started discussion about developing "**killer apps**" for industry. He pointed out that many IAR applications end up being *over kill*. Navab stressed that AR solutions must be scalable, user friendly, and reliable. IAR developers must study the industrial processes, it's requirements, end-users, user environment and consider all aspects of the problem before developing solutions. Navab also mentioned several technological challenges such as real-time tracking, need for on-line calibration, lack of mobile computational units, etc.. Many of these are solved after his report but acceptance of IAR did not improve significantly. Georgel [70] picked up this discussion again and conducted a survey of IAR applications. Georgel uses a rubric using Navab's comments to evaluate the success of IAR applications in the industrial domain. This rubric considers several factors like workflow integration, scalability, cost-benefit, out of the lab, user-tested, out of developer's hands, and involvement of the industry. The eventual conclusion was that very few IAR applications were successful outside the research labs and to get a steady footing in the industry the applications needed better workflow integration.

Our goal was to design spatial AR solutions for manufacturing assistance in the automotive and aviation industry. We took inspiration from conclusions of both Navab [153] and Georgel [70]. We created a list of requirements and followed them while designing our work packages in the ARVIDA project. In our opinion, most of these requirements are generic and should be modified as needed to suit the scope of work for other IAR applications. For this discussion we refer the IAR application as a tool that performs measurement and visualization.

### Accuracy

Manufacturing tasks in the automotive or aviation industry require accurate tools for performing measurements. In manufacturing, a range of tolerances is provided for each part along with its design. Tolerance is defined as the threshold of error that can be accepted between the produced part and its design specifications. Accuracy requirement is directly related to tolerance values. If this requirement is not met the parts may not fit together or the overall system may malfunction. As a consequence, the parts are rejected which directly increases the production costs in terms of effort, time, material, etc. The range of tolerances can vary from mm to cm range depending on the importance of part and the field of application. For example, turbines for the aviation industry or engine design for automotive manufacturing require high accuracy and therefore allow low tolerance range. The tools used for manufacturing assistance must be accurate enough to detect deviations specified in the tolerance range. Accuracy requirements must be considered before deploying IAR solution in manufacturing.

### Reliability

The reliability of the method is directly related to the precision requirements of the tool. Tools that produce reliable and repeatable results are specifically necessary for manufacturing scenarios because each part is measured multiple times at different departments for different purposes. For example, a worker measures parts to perform drilling or welding operations, and a quality inspector measures the same parts to evaluate the accuracy of the drilling location. Deviation in measurements may lead to confusion and rejection of parts. For IAR applications, it is extremely important that measurements are reliable and repeatable which guarantees augmentations with high precision.

### Scalability

It is one of the most important success factors for any manufacturing practice. If a method or a tool is useful it must be easily replicated to achieve standardization i.e. same processes are used between different departments and locations. In modern manufacturing practices, all solutions need to be scalable because they save time and remove the need to define case specific workflows. For technology based tools scalability involves ease of deployment, maintenance and replacement. This has been a major limitation in acceptance of many IAR solutions [70].

### Workflow Integration

This requirement is specifically stressed by both Navab [153] and Georgel [70] in their report. All manufacturing practices must be well defined and each process must function seamlessly to establish efficient manufacturing pipeline. This is equally valid for IAR solutions that are designed to replace an existing method or practice. Which means that the input and output of the AR solution must be compatible with other practices and should not break the work flow. Often, workflow incompatibility means forcing workflow change for other non-AR unrelated practices. Industrial partners are reluctant to make such unnecessary changes. Changes to work flow require reevaluation of the revised workflow to determine efficiency, which is a time consuming process and must be avoided. In summary, IAR solution must be designed to do exactly what is required without disturbing other processes.

### Cost benefit

This is an important requirement for decision-makers. The investment done for the application must be rewarding enough and add more value to the existing practices. Any IAR application that is not able to save time and money will not be accepted. AR solutions require expensive hardware such as cameras, projectors, or HMDs. Therefore, the cost of deployment must meet the returns. Other requirements such as workflow integration, robustness and scalability also affect cost benefit analysis.

### User Friendliness

This requirement is important from the user's perspective. Any new solution that is closed to the previous solution in terms of working principle minimizes the need for training. In manufacturing the same tools are used by workers, engineers, and designers. Each user group must be able to use and navigate through the user interface with minimum efforts. IAR solutions often end up being technically confusing or operate as a black box (i.e. the user does not exactly understand working principles). This factor leads to trust issues with end-users. Moreover, the user prefers interfaces that indicate system failures and have well-defined methods to fix the problems easily. For example, IAR applications must be able to suggest operational accuracy or suggest a possible error in calibration. In such cases, users must be able to calibrate the systems fast and verify its accuracy before performing the required tasks. Such provision increases the user's trust in the system which leads to faster acceptance.

## 3.2 Alignment and marking process

In this section, we will introduce two widely used industrial practices i.e. Alignment and Marking. These processes are the main focus of our research and they are chosen because they are used for almost all manufacturing practices. In our opinion, finding AR solutions for such fundamental practices will cement the position of AR in the industrial environment. In the following text, we will explain both processes and discuss existing solutions with certain advantages and limitations.

### 3.2.1 Alignment Process

The alignment process is used to position a part or a workpiece in a unique orientation for working with it. A unique method for part orientation is important for standardization, accuracy, tolerance management, etc. Let us understand the importance of the alignment process with a simple example of part measurement for the task of drilling. Typically, the worker is provided 3D coordinates of target locations in the design documents. The workers have to find these points to perform the operation. In such cases, it is easier for the worker if the measurement system performs measurements directly in the coordinate system of the part. To do that the coordinate systems of the part and the measurement machine must be aligned. Besides measurement tasks, alignment processes are also used for guiding cutting or milling tools. The process is essential for assembly tasks, where each part having its own local coordinate system are placed together in a global coordinate system. It is extremely useful to have a fixed alignment strategy to simplify the process and avoid the wrong assembly.

There are multiple methods for part alignment and we will discuss two methods relevant to our work: *a.* jigs and fixtures [101] *b.* Reference point system (RPS) [4, 144]. Both are concepts for mechanical methods for aligning a part in a unique orientation.

### Jigs and Fixtures

Jigs and fixtures are used extensively in the industrial environments, especially for mass production units. They are particularly designed parts purely for alignment purpose. Jigs and fixtures [101] come in various types depending on application such as drilling, marking. Design of the part also influence choice of jig, e.g. circular, spherical, planar.

One of the simplest types is a template jig (See Fig. 3.1). The working principle of template jig similar to stencils which are commonly used for drawing patterns on paper. Jigs and fixtures for three-dimensional objects can be compared to a mold or a form with an exception that very few contact points are used to place the part. These contact points are designed with the 3-2-1 principle for unique placement (explained later). Simple criteria followed for part placement is that all contact points must touch the part. The 3-2-1 principle ensures that there is only one way to achieve this result. This design is *foolproof* as no special knowledge or skill is required to perform the alignment. Fixtures are generally designed for placement while jigs are also designed to avoid multiple measurements or guide tools. For example, template jigs displayed in Fig. 3.1 has holes at a fixed distance, and drilling through the holes guarantees equidistant holes without the need for performing measurements.

Some jigs and fixtures are standard designs while most have to be produced specifically for a part. Jigs may be used for tolerance management but it is not their primary purpose. It is a costly investment but proves to be cost-effective when used multiple times for mass production. Additionally, this method saves time and ensures the accuracy of tasks at hand, which reduces failure. The overall process is user friendly and relatively easy to execute. The user does have to perform cumbersome steps such as physically move the jig and the part at the same location.

**Fig. 3.2.** The figure shows a work piece fixed in a unique locations using the 3-2-1 principle: **3** contact points in the x-y plane, fix translation along the z-axis and roation around the x-axis and y-axis. **2** contact points in the y-z plane, fix translation along y-axis and rotation around the z-axis. Finally, **1** contact point in the x-z plane, blocks translation along the x-axis. Image ©Extend3D GmbH, published with permission.



**Fig. 3.3.** An example of RPS points mentioned with their tolerance values. RPS1 is a main mounting point. RPS3, RPS5 and RPS5 are all blocking motion along z-axis. This arrangement is over constrained and therefore tolerance values are distributed. Image ©:GOM GmbH, published with permission.

RPS is a method for alignment process which also follows the 3-2-1 principle mentioned earlier. It is a concept that is integrated with the design process to align multiple parts without using external components such as jigs or measurement tools. RPS features are defined on the parts such that RPS features of different parts coincide with each other during the assembly process. It is a clever solution to avoid complicated measurements while putting parts together. The concept is easier to understand as the worker has to match defined features. If all features come together successfully it implies that the assembly is correct.

Every rigid object has 6 degrees of freedom (3 rotation - 3 translation) in a three-dimensional reference space. 3-2-1 principle is a simple method for methodically removing degrees of freedom by adding physical constraints on movement in specific directions and thus fixing the object in a unique orientation. Figure 3.2 shows a simple implementation of the 3-2-1 principle for part placement, where 6 contact points are used to fix the motion of the part.

RPS is a well-defined protocol for integrating the 3-2-1 alignment principle at the design stage. As a common practice, existing features such as holes, slots, or edges (Fig.3.3) are selected as RPS features. A fixed naming convention is followed for RPS features to define the type of feature and fixing directions. The name also conveys priority of point i.e. main mountings point are represented with a capital letter and supporting points with small letters. For example, RPS1 HxyFz is the first point (RPS1), capital letters H and F convey that it is a main mounting point, and hole (H) is fixing motion along xy direction and the surface with the hole is fixing motion in the z-direction.

Manufactured parts often do not meet the design specifications and some surfaces may deviate from their original CAD designs. We define this as **geometric modeling problem**, it is discussed in detail in the next chapter. This problem may affect the alignment process as well. However, the RPS positions are manufactured with high accuracy and their positions are trusted for accurate alignment. Due to the accuracy of positions, the RPS alignment concept also functions as a protocol for tolerance management. RPS features are strategically defined which allow part assembly only if the manufactured parts are within the tolerance range. Features are selected based on tolerance requirements. Features such as holes and points are suitable for low tolerance scheme because they are defined with fixed 3D coordinates. On the contrary, features like surfaces and long holes allow higher flexibility while fixing motion along a defined axis. Parts that do not align as per RPS guidelines are considered faulty and not used. In some cases, more than necessary RPS points are defined for blocking motion along a specific direction to distribute the errors in deviations (Fig.3.3).

The readers should note that RPS is a design concept. It is also used for designing jigs to maintain consistency in alignment protocols of the industry. RPS positions are also used for aligning parts with measurement systems (tactile or optical) to the object coordinate system. It is very useful for scenarios where multiple parts are assembled to make one single product such as a car or plane. This allows easy comparisons between measurements done at a different stage of the assembly.

In general, alignment practices with jigs and RPS are deployed successfully in most industries. The alignment processes are manual and the working principle (3-2-1) is very easy to

**Fig. 3.4.** The image displays a typical use case of optical probing technique for 3D measurements. The markers on the probe are used to track motion of probe and compute its pose. The dimensions of the probe are known and different points are measured by touching the tip to the part. Image ©:GOM GmbH, published with permission.

understand for workers and designers. They save time and effort because unwanted measurements and subjective measurement biases are avoided. Well defined alignment practices accommodate for manufacturing deviations which provide a certain guarantee that overall assembly will work and parts will function as planned. Low error rates instill a high degree of confidence in workers and engineers. The methods are repeatable and scalable for large scale manufacturing operations. In a way, existing alignment methods fulfill all parameters defined in the earlier section.

## 3.2.2  Marking Process

The marking process is used for marking locations of interest for manufacturing tasks such as stud welding, drilling, etc. This process is almost always preceded by the alignment process. In a typical workflow, workers are given a list of desired marking locations of the points with the design specifications of the part in digital format as a CAD model or in paper format as isometric-orthographic views. The worker uses measurement tools or specially designed jigs (explained earlier) to mark these locations on to the real parts. Marking with specialized jigs is not always possible because it is not feasible to design jigs for all cases. Especially in heavy engineering (aviation, shipbuilding, etc.) or rapid prototyping, where parts are produced only once and the task of producing special jigs is not cost-effective. Alternatively, the user operates a measurement device and finds the required coordinates on the part for marking. We will focus our discussions on measurement device based marking.

Mainly two type of measurement systems are used for marking process: tactile and optical. Tactile systems operate on a simple working principle: mechanical motion of a probe is used

to determine 3D location of the tip of the probe. The motion of the probe is often limited in rotation and translation based on its design. This method is very accurate and precise but less suitable when multiple measurements are required simultaneously. Parts have to be brought in close vicinity of the system to perform these measurements. Tactile methods are considered extremely reliable but they can be also cumbersome to use, due to limited mobility. The alignment is performed by using jigs or by touching the probe at several predefined locations, for example RPS. Once the measurement device's coordinate system is aligned with the coordinate system of the part, all measurements are transferred displayed in the coordinate system of the part. Alignment concept is similar to the point correspondence based 3D pose computation explained in the previous chapter (see 2.4.1). Alignment is not necessary for measuring relative quantities such as distances between points, diameter of a circle or length of a part.

Optical probe-based measurement methods are more versatile in terms of use (see Fig. 3.4). The probe is tracked in a contactless manner using a marker-based multi-camera optical tracking system. The pose of the probe is tracked which means that the position of the tip is always known in the sensor coordinate system. The alignment process is touch-based similar to a tactile probe-based system. It should be noted that alignment is only necessary for measuring points in the part's coordinate system. Relative distances are measured without performing alignment, for example distance between drilling holes or length of a part . Optical probes are cost-effective and mobile alternative to tactile measurement systems. The accuracy is very high (< mm), however not as high as the tactile methods. Both tactile and optical marking methods are extensively used in industrial scenarios. The selection of measurement device purely depends on the requirements of the job.

The marking process has very high accuracy requirements to handle manufacturing tolerances. It is difficult to manufacture parts that are exactly similar to the CAD designs. Parts are often deformed or damaged during other processes such as heat treatment or vibration on shop floor. This results in deviation from CAD designs. Earlier, we refer to this problem as the geometric modeling problem [151]. Parts with deviations above the tolerance values are rejected during alignment process. However, deviated surfaced within the tolerance range remain in the process. It is likely marking locations are affected by deviations and the desired points cannot be marked on the surface. In such cases new marking points are defined by transferring the old location to the surface. It is typically done by moving the point along one of the principle axes. New points are marked using the measurement systems. It should be noted that error in alignment can affect measurements and therefore RPS positions are used for alignment, which guarantees that pose computation is free from deviations. Geometric modeling problem is omnipresent in manufacturing. Marking practices with jigs are often useful for circumventing this problem. For example, the template jig displayed earlier (Fig.3.1) guides the drill machine to produce equidistant holes even if the part's surface is deformed.

## 3.3 Spatial AR for alignment and marking

In this section, we will discuss the idea of using Spatial Augmented Reality as an IAR application for alignment and marking process. First, we will discuss the technology, the workflow, and compare it to existing methods. Finally, we will discuss the limitations of performing

alignment and marking with SAR, and build up the background for our research contributions covered in the next chapter.

Schwerdtfeger et al.[188] explored the idea of using laser projectors as an alternative to head mounted displays for AR applications. They explored different possibilities of using laser projectors for manufacturing guidance applications. They used a rigid projector-camera (video) setup on a tripod and laser projector mounted on helmets. This design of such equipment is comparable to optical tracking based measurement systems. The study concludes that tripod mounted systems are better suited for industrial environments as they are accurate and robust. In a follow up work, Schwerdtfeger et al. [187] demonstrates versatility of the tripod mounted SAR setup by performing augmentations without using markers.

The author promotes use of a laser projector, as the augmentations can be viewed with high contrast on the shiny industrial parts in well lit environments. Therefore, laser projectors do not require modification in the light conditions to highlight augmentations. Laser projectors are also explored in other AR applications with industrial robots [176, 217] and surgery [155]. It is shown that accuracy of laser projectors is close to sub-mm range. Laser projectors work with the rotating mirror principle and the errors in projection are due to incorrect angular position of mirrors. Due to angular nature the error increases with distance of projection surface [17, 188] Schwerdtfeger concludes that laser projection based SAR systems have good potential of becoming a dynamic and portable visualization tools for industrial processes. SAR tool that can project instructions directly on the part have many advantages to offer for industrial applications, we will discuss them later in this section.

Schwerdtfeger [186] and Keitler [106] worked on the topic of IAR applications for their dissertation. Based on their experience they concluded that SAR applications with projector-camera can reliably meet requirements of the industrial scenario (Sec. 3.1) and decided to optimize the product design for commercial IAR solutions. The concept is now turned in to a industrial product which we used for our experiments (Fig. 3.5). The research presented in this thesis is done entirely using this device. As of now, it is integrated in manufacturing practices of several manufacturers across the globe [2]. It is a versatile product that caters wide range of applications in multiple industries such as automotive, aerospace, transportation and marine [67].

### 3.3.1 Device specifications

The hardware consists of a stereo camera system and a laser projector rigidly mounted in one single housing (Fig.3.5). The setup is mobile and it can be mounted on a tripod or connected to a robotic arm. The device is always connected to a computer for computational and hardware communication. The cameras and projectors mounted on the devices are commercially available products. Ideally, both can be replaced and customized based on the need of the user. Camera resolutions are up to 10 Megapixels and the working range is approx 1.3 to 3 meters. The projection area is 60 degrees with high precision and can be extended up to 80 degrees. Additional, hardware-specific details can be found in the datasheet of the

---

[2]https://www.extend3d.de/en/

**Fig. 3.5.** Werklicht Pro-L, consists of two stereo cameras and a laser projector. Image ©Extend3D GmbH, published with permission.

product [67]. A software interface is designed to import the CAD design of the industrial part. The software interface is used for planning augmentations and pose computation strategies.

## 3.3.2 Working principle and workflow

Here, we will explain the general working principle of the device along with the operational workflow in industrial settings. In principle, it is a projection-based SAR device with two optical sensors (video) for tracking. It is operated in inside-out configuration with a fixed referencing approach (see sec. 2.4). The device is also compatible with display-based design concepts mentioned in sec. 2.2 i.e. fixed display - moving camera and moving camera - moving display. First, we will explain the working principle very briefly and then will explain the different stages of the workflow. A typical workflow is divided into two stages: planning and Execution. Additionally calibration is performed for setting up the device during manufacturing and, therefore, it is not considered part of the typical industrial workflow. General information on calibration is covered in the introduction.

### Working principle

The working principle of the device can be explained in two steps i.e. *tracking* and *projection*. First step is to compute object pose to align object coordinate systems with device coordinate system. Tracking strategy is based on stereo-camera strategy (sec. 2.4) similar to optical measurement systems discussed earlier. The second step involves using the pose information to highlight desired locations on the object.

The 3D position of some features on the object are triangulated in the device space. These features on tracking object are natural features such as lines, edges or holes or artificial markers (coded or uncoded) as displayed in fig 3.6. Coded markers are attached to the part as part of workflow for reliable enhanced tracking results. Uncoded markers are often already attached to the part as part of photogrammetric measurement processes [150] The positions of features or markers in object space are predefined in the CAD model during planning stage. Using this information the object pose (6-DOF) is computed in device space. Pose is nothing but a transformation matrix used to transfer points from object space to device space and vice versa. This transformation matrix remains valid until object or device position is changed.

**Fig. 3.6.** The image shows two types of markers used in the industrial scenario. Markers with encoded patterns are uniquely identifiable which makes image processing and feature matching easy. The image also shows specialized mounts used to place encoded markers on the industrial parts. These mounts are magnetic and their design is customized to be used in a hole, surface, corner and edge. The uncoded markers are circular rings that are easily detectable in images but do not have unique identification code. Uncoded markers are popularly used in industrial scenario for applications of metrology [127]. Both are often used for camera calibration procedures and for object tracking [150]. ©Extedn3D GmbH, published with permission.

After this step, position of any desired feature in object space is transferred to the projector space. The pose of projector w.r.t device space is determined during the calibration stage. Calibration parameters are computed as part of initial equipment set-up (covered in sec.2.4.6). They remain consistent as long as position of cameras and projectors do not change.

### Planning

This stage of the workflow involves planning the work steps for the end-users. In a typical industrial workflow, the work steps are designed by the engineer, and a list is provided to the worker via a paper-based or a digital medium. For example, a stud welding task has a list of work steps with details of job number, welding positions, part number, part dimensions, reference CAD design, etc. The work steps are generally prepared using a software interface with CAD models (Fig. 3.8). We use a similar software interface to plan the desired augmentation patterns at desired locations using the CAD model. Figure 3.7 displays an example of the planning stage with the software interface and the projected contours with special text depicting job instructions. After planning is complete the work steps are exported and work packages are created.

The tracking strategy is also defined during the planning stage i.e. marker or marker-less. This information is also exported with the work packages. It is useful to verify the tracking strategy at the design stage by testing it with a prototype as shown in figure 3.7. Testing helps to avoid unnecessary problems such as the line of sight issue for marker-based tracking strategies. Ideally, the placement of markers must not impede the worker from working on the

**Fig. 3.7.** Software interface can be used for both planning augmentations (left) and executing the augmentations (right). Image ©Extend3D GmbH, published with permission.



**Fig. 3.8.** Example of a CAD model of an industrial part. The holes, lines and arcs can be selected in the interface for planning projections. Image ©Extend3D GmbH, published with permission.

task object. The positioning of markers w.r.t the device is crucial for error-free augmentation from multiple perspectives.

### Execution

During the execution phase, the workers used the exported work packages using a web interface [184]. As displayed in Fig. 3.9, each work step contains an image of the CAD model depicting a recommended position of the part placement while facing the device. The screenshot also contains information about the position of markers. Figure 3.6 shows different types of markers and mounts available for the task. The worker places markers on the part and selects the desired work step. The device computes the pose and the device starts making augmentations on the part as shown in figure 3.9.

The tracking strategy changes based on the requirement of the application. Marker-based methods provide the highest form of accuracy and robustness. For the following discussions, we will assume this method as a default tracking strategy. The placement of markers is an additional step that a worker has to perform. In this case, a display based AR concept is used to avoid marker placement related errors. Uncoded markers are an exception because they are added on the part for other measurement purposes [127]. The position of these markers (in object space) are directly imported from the computed for AR applications [150].

**Fig. 3.9.** The image shows a work step being selected on a tablet and real-time projection on the industrial part based on the selection. Image ©Extend3D GmbH, published with permission.

All work packages contain the position of markers ($P_{obj}$) defined in object space with the CAD model. The 3D locations triangulated ($P_{device}$) by the system are used for pose computation ($Rt_{obj}^{device}$).

$$P_{device} = Rt_{obj}^{device} \cdot P_{obj} \qquad (3.1)$$

$$p'_{img} = k \cdot Rt \; P_{device} \qquad (3.2)$$

$$\text{Reprojection error} = \text{dist}\,(p_{img}, p'_{img}) \qquad (3.3)$$

Three markers are sufficient to compute the 6-DOF pose of the objects using the equations given in 3.1. Additional markers are used to compute the error in pose computation and to verify the validity of the pose. The pose is used to project the marker locations on the image space using equation 3.2. In an ideal scenario, the position of the markers detected in the image, $p_{img}$, and the position of the markers projected in the image $p'_{img}$ must be identical. The distance between the projected and the detected points is used to measure the error in pose computation, also known as reprojection error. Overall reprojection error is computed by taking root means square of all sum of all distances i.e. dist $(p_{img}, p'_{img})$ [85]. The reprojection error is computed in pixel space and it is further converted in metric space using the calibration data. Error in metric space provides a quantitative estimation of the augmentation error for the users. A minimum error threshold is defined with each work package to avoid faulty guidance. The reprojection error for each marker is computed individually and markers contributing to higher errors are removed for pose computation.

Earlier we mentioned that the device offers two types of visualization strategies with SAR i.e. projection and display based. The laser projector is used to project work step related instructions directly on the part. This is a primary source of augmentation for worker assistance. The screen-based visualizations are used to provide further information about the accuracy of marker placement, tracking error, etc. The indication of error is done by augmenting the marker positions in the image with different colors i.e. green for good and red for bad. The

**Fig. 3.10.** The image shows an example of two different industrial parts augmented using the laser projector. Image ©Extend3D GmbH, published with permission.

user is able to use the screen-based augmentations to position the part such that all markers are detected by the system. After the job is done, the worker uses the same web interface to update the status of the job and export a report of the job. In this way, the software interface works as an assistant to the worker by guiding through the process in a step by step manner.

### 3.3.3 Use cases

SAR devices like **Werklicht Pro** are used for multiple applications in the industrial domain [67]. Primarily, the device is proposed as an interactive data visualization tool for the industrial scenario. The major advantage is that it can be used at various departments while keeping more or less the same workflow. It is proposed as a digital alternative to the existing manual solutions, for example alignment and marking with jigs. The uniqueness of the data visualization technique is also exploited for completely new applications e.g. rapid prototyping, painting. This method does not make existing tools or methods completely redundant, rather it is designed to as a multipurpose tool. We have listed some example use cases below.

#### Alignment and Marking

We have explained alignment and marking processes is earlier with existing methods. SAR presents a digital alternative to the classical manual techniques for alignment and marking processes. The alignment is automatically done by placing the part in front of the system. The system is mobile and can be brought to the part if required. SAR device combines the step of manual alignment and finding the target points using a measurement device. The worker can skip the steps and directly perform jobs at marked locations. The marking is done by highlighting the target locations on the surface of the part using the projector. In a way, the SAR system also acts as a measurement device.

There are additional advantages of using the SAR system for alignment and marking task. It is also possible to mark multiple points or complicated contours which increases the speed of operation. Direct augmentation on the object is extremely intuitive and workers do not have to go through complicated design plans or CAD models. Jigs and fixtures use mechanical alignment principles and they need to be produced specifically for different parts. SAR devices work with all types of parts and thus work more like a *digital jig*.

**Fig. 3.11.** The image shows example of the verification process where projections highlight the desired locations. This is a simple and intuitive way of verifying quality of work. Image ©Extend3D GmbH, published with permission.

## Quality control

Quality control is extremely important in manufacturing to avoid production failures. As part of this process, each work step is verified after the part is assembled. The personnel at quality checking have to follow the same work steps as the workers. This requires repeating the process of alignment and measurements to check for missing steps or erroneous markings. SAR devices are extremely useful in speeding up this process. All locations of work steps are highlighted automatically and the verification officer can perform verification without performing any measurements.



**Fig. 3.12.** The image displays a mock car model with several uncoded markers. The position of these markers are computed using optical measurement systems. The image displays a laser pointer projected at two different markers. This is one example use case of our SAR device. The measurements are directly presented on the part which makes decision making more intuitive. The image is part of test dataset Image ©Extend3D GmbH, published with permission.

### Visualization of metrology data

This is a special use case that is applicable only when uncoded marker-based metrology tools are used in the process. Uncoded circular markers are applied to the surface of part and their 3D positions are reconstructed using photogrammetric tools [127] (see Fig.3.6). The reconstructed positions are registered in the part coordinate system and used for comparative measurements [150] such as results of vibration test, surface deviations, etc. Our device with video cameras is capable of computing pose using the same uncoded markers. The system works in sync with existing industrial practices and alignment is performed without adding special markers. Additionally, the information about surface deviations is obtained from metrology processes is displayed directly using the SAR system.

### Other use cases

Other common use cases with SAR devices include the use of projection for painting, bracket assembly, and designing jigs. Assistance for painting is rather a unique application which is possible only because augmentation is directly done on the parts. In shipbuilding or aeronautical industries, large vessels and planes are painted with colorful patterns for improved aesthetic look or marketing. Painting large curved surfaces requires specialized templates that compensate for the curvature. The job is performed by skilled painters because drawing complex textures on very large objects is a challenging task. SAR devices are used to project contours of patterns directly on the object. The patterns are planned using the software interface with the 3D geometry of the part. This is a novel use case that is extremely tedious if SAR systems are not used.

## 3.3.4 Advantages and Limitations

Earlier we discussed requirements from IAR applications and now we will discuss the advantages and limitations of the SAR system w.r.t these requirements. There are some limitations in terms of accuracy and reliability. We addressed these limitations with our research in the ARVIDA project (next chapter).

### Accuracy

Optical tracking systems with a video camera used in our applications can achieve pose computation accuracy in the sub-millimeter range. It is comparable with other optical tracking solutions (IR-based) used as probe-based measurement systems. Projection accuracy is an additional consideration with SAR systems. There are three causes of errors in our SAR system i.e. calibration error, pose computation error, and geometric modeling error.

- **Calibration errors** are static errors and therefore affect all measurements done by the system. Error in the calibration of cameras leads to triangulation errors which affect pose computation. Further, the camera-projector calibration error leads to errors in augmentation. The calibration errors can increase over time if the device experiences vibrations or physical impacts which results in the shift of camera-projector locations. These errors are minimized by performing recalibration but cannot be eliminated.

- **Pose computation** errors are dynamic in nature because they depend on accuracy of other operations such as image processing, feature matching, etc. Pose errors also include human errors, for example wrong placement of markers.

- **Geometric modeling** errors are introduced by the manufacturing processes. The deviation between the CAD model and the actual manufactured part contributes to pose errors in marker positions that are affected or projection errors if target surfaces deviate. This error is a dynamic error and may be introduced at any stage of production.

Overall, it is very difficult to separate the real source errors from these three errors [12]. Calibration and pose computation are well-studied problems but geometric modeling is less studied in the context of SAR. We worked on this problem to increase the accuracy of our SAR system using a closed-loop AR solution [151]. The implementation details and results are detailed in the next chapter.

### Realibility

Reliability in the context of SAR systems can be defined as the ability of the device to augment the same location repeatedly from different perspectives. Assuming that calibration and geometric modeling errors are minimum, the precision of the device varies within pose computation error. A typical range of variation is millimeters i.e. pose computation accuracy.

The results of marker-based tracking solutions are generally considered reliable and repeatable since similar methods are used for probing and referencing with other optical tracking based technologies. It is important to know the accurate position of markers in object space to provide reliable and repeatable results. Coded marker-based solutions with our device require the markers to be placed at locations planned by the designer. This is a possible source of error because the actual locations may vary from the planned locations. Thus, the existing solution is less reliable for high accuracy tasks. To increase the trust of the users it is vital to adopt ideas from existing practices. Tracking solutions based on uncoded markers are a good example of a reliable solution. In this case marker positions are measured after they are attached to the part removing the possibility of placement error. Markerless solutions do not deal with the issue of marker placement. However, they face challenges in industrial environments due to varied lighting conditions and property of the material i.e. shiny or reflective surfaces. Rigorous testing protocols are needed to prove the reliability of markerless methods and to gain the trust of end-users.

We learned that the idea of using RPS features as possible locations for coded markers is not yet explored for SAR systems. We concluded that integrating RPS features in the process can increase the reliability of our solutions. Therefore, we developed a strategy for marker and markerless alignment using RPS features [148]. The implementation details and results are explained in the next chapter.

### Scalability

Earlier, we mentioned several use cases of SAR systems in various departments, which proves that the SAR approach is extremely scalable. The system requires engineers and workers to be trained for operating the equipment. However, the process of setting up the task remains more

or less the same, irrespective of the final application. The working principle of the system is easier to explain and therefore does not face objections from engineers or workers. The system assists the workers to achieve the task in an effective and error-free manner. The software interface via tablet or PC is designed to be as close as possible to other workflows. The parts of the device are easily replaceable and system downtime is extremely low when replacement is required.

### Workflow integration

The planning stage is designed to be compatible with existing software practices with CAD-models. Engineers are used to operate with 3D models of parts and the additional steps of planning projections and marker positions are not cumbersome. While it is an additional step for the designer it replaces the need for creating work step specific instruction lists for workers. In the case of uncoded markers, the step of planning marker positions is completely avoided.

From the worker's perspective, they have to place the markers on the part for tracking and use the software interface to select the work step before starting the job. This step replaces multiple steps from the existing practices i.e. setting up the part on a jig or setting up measurement device, read part instructions, measure point of interest, and mark them. The workers have to perform less number of steps for SAR solutions. For example, the task of part alignment and measure of target locations is reduced to placing the markers on the part. This feature not only saved time but also minimizes the possibility of human errors. Most importantly, the SAR solution does disrupt any existing process.

### Cost benefit

In manufacturing, any practice that reduces effort, material, or time is considered cost-effective. The biggest advantage of the SAR system is the range of applications it can handle with the same set of tools i.e. equipment, markers, and computer. The product is rightfully promoted as a digital jig for manufacturing purposes. It reduces the need for creating special jig based solutions i.e. saving material and time. The SAR solution reduces work steps (explained above) and provides better means of data visualization for job execution or verification, which saves time and effort. The system indicates errors in marker placement or alignment reducing production errors. Overall, the system is cost-effective, as it saves material, effort, and time.

### User friendliness

The existing workflow is configured to be user friendly. The engineers have complete access to CAD-based models and the possibility to verify projection and tracking strategy while planning. The workers are given simplified work steps via a web browser. The task can be started by selecting the correct work step and placing the part in front of the machine. Interactivity is the biggest advantage of the system. Tracking is done in real-time and projections readjust in real-time when part or system itself is relocated. The tracking and projection accuracy can be observed and error is color-coded based on the engineer's specifications. Moreover, marker placement is also verified in most workflows, and error in marker placement is indicated via screen-based AR interface.

The system does face limitations on the line of sight and tracking (also projection) is lost when markers are not detected or occluded. Stereo tracking has its limitations in terms of speed and accuracy of tracking (covered in the introduction). Laser projectors provide required contrast however the content projected by these projectors is limited. From our perspective, the step involving marker placement and planning is a trade-off for having a reliable and accurate data visualizations tool.

# IAR solutions for Alignment and Marking

<div style="text-align: right">4</div>

In this chapter, we will discuss the solutions developed by the author as part of the ARVIDA initiative. The research contributions are focused on the process of alignment and marking with SAR devices. More specifically, we have developed methods to improve accuracy and reliability to enhance the acceptance of SAR systems in industrial environments. The research is presented in two main parts. In the first part, we propose two methods (marker-based and markerless) for pose computation using RPS features with a focus on increasing accuracy of part alignment. In the second part, we propose a closed-loop solution for marking process to reduce errors stemming from the problem of geometric modeling. Each method is explained with relevant examples and the experiments are conducted with real industrial parts. The final section of this chapter includes a combined conclusion and impact of our research.

We would like to remind the readers that this research was conducted between 2014-2017. The problems discussed in this chapter are very relevant to the acceptance of IAR applications. These specific problems were chosen to demonstrate that novel solutions can replace existing practices while maintaining relevance to the industrial workflow. The marker-based implementation of the RPS system (in sec. 4.1.1) is presented in the final report of ARVIDA project [184]. The work on markerless alignment (in sec. 4.1.2) was presented in the ISMAR conference in 2016 [148] with the title "Frustration Free Pose Computation For Spatial AR Devices in Industrial Scenario" ( ©[2016] IEEE). Our research on geometric modeling (in sec. 4.2) was first presented at the International Symposium on Mixed and Augmented Reality (ISMAR) conference in 2015 [151] with the title "A Step Closer to Reality: Closed Loop Dynamic Registration Correction in SAR" ( ©[2015] IEEE). These works are used as basis for the text presented in this chapter.

## 4.1 Solving alignment problem for SAR using RPS features

In this section, we will describe two different strategies for improving the accuracy and reliability of part alignment. We have proposed the integration of RPS protocol in the current workflow with SAR systems. In principle, RPS features are defined to physically block 6-DOF movements of the part by fixing movement along with specific directions. This could be explained as a physical way of defining a 6-DOF pose of the object (see chapter 3) where the object is manually rotated and translated to be placed in a unique location. We propose replication of the same idea using an optical tracking system. We use 3D locations of RPS features and computes the pose using the directional constraints defined in the RPS protocol

for building an equation system. Unlike the manual method the part does not have to be placed in a specific manner but the system adapts to the existing position and location. This way we offer a contactless alignment solution without compromising the reliability offered by the existing practices.

### 4.1.1  RPS alignment with markers

The first strategy is implemented using markers based tracking solution. We used coded markers and place them at the locations of defined RPS features. Figure 4.1 shows a mock example of the implementation of RPS based referencing of parts with markers. The RPS protocol for the part is defined by the engineers. The workers place markers based on the RPS strategy defined for the part and the 6-DOF pose of the part is computed according to the defined strategy. RPS strategy is nothing but a selection of features defined for alignment. This strategy depends on the features present on the part i.e. holes, edges, or surfaces.

#### RPS configurations

Different features of the parts are used as RPS features because each part is designed uniquely and does not have same type of features. To accommodate the variations in part design different strategies are used to define RPS features. Circular features or corners act as fixed points and they are used to lock motion in all 3 directions. Edges (lines) are used to block one degree of freedom and usually defined by two markers. An edge can be combined with another marker in the same plane to fix more degrees of freedom. Similarly, 3 markers on a same plane are used to define a plane. These three features are used in different combinations for defining RPS configuration based on the 3-2-1 alignment principle. The total number of markers required can vary between 3 to 6 [184].

**Fig. 4.2.** Example of different RPS strategies displayed on the same model. The image shows typical use cases where markers mounted on edges, holes and surfaces. ©Extend3D GmbH

- 3 point configuration: Three features defined with holes or corners. One feature can be primary which blocks motions in all three directions. Second feature blocks motion in two directions and finally third feature blocks motion in one direction (top center 4.2).

- 3 plane configuration : The planes are defines where primary plane is defined by 3 markers. Secondary plane is defined by 2 markers and tertiary plane is defined by 1 marker (top right 4.2).

- Plane-line-point configuration : 3 for surface, 2 for line and 1 for point

- Plane-line-point configuration : Two edge , one point on corner, all points on one plane define the plane (bottom left 4.2)

- Plane-point-line configuration : Three markers on the same plane, two of them placed at fixed locations and third marker along a line in a slot hole. (bottom right 4.2)

- Plane-line-line configuration : Three markers define a plane, two markers define one edge and one marker defines another edge.

- Plane-line-line configuration : The markers on the same plane, two of them define one edge and third one defines another edge. (bottom center 4.2)

### Pose computation strategy

The strategy for solving the pose equation is determined by the RPS configuration. RPS configuration for the part is defined by the design engineers while designing the parts. The engineer defines the marker placement positions while planning the work steps. During the execution stage, 3D marker positions are computed and feature specific RPS guidance is used to build the direction-specific equation system.

A simple case is a three-point configuration. It is similar to classic method for pose computation using absolute orientation [92]). The position of fixed points (holes and corners) in known in object space. The same points are reconstructed with the stereo points and used as direct correspondence for pose computation i.e. $P_{obj} = Rt \cdot P_{device}$. Each 3D correspondence offers

three equations (one for each coordinate as shown in eq. 4.1) and therefore three 3D-3D point correspondences (9 equations) are sufficient to solve the 6-DOF pose. There are several cases when more than 3 RPS features are defined and the solution in such cases is overdetermined i.e. more equations than unknowns.

$$
\begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}_{obj} = \begin{bmatrix} R_{11} & R_{12} & R_{13} & T_1 \\ R_{21} & R_{22} & R_{23} & T_2 \\ R_{31} & R_{32} & R_{33} & T_3 \\ 0 & 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}_{device}
\tag{4.1}
$$

Let us discuss other special cases when only partial information is used as per the specifications of RPS. A point feature may be specified for blocking movement in only one or two directions. For example, a feature defined as Hxy is a hole (point) that blocks movement in x and y-direction. Therefore only two equations (x and y correspondences) are used for pose computation and equation with z is discarded.

For maintaining design simplicity RPS features are usually aligned with the coordinate system i.e. features are defined on principal planes (XY - YZ and XZ) [4]. This design makes sure that motion is always blocked along the principle direction. This design constraint is used while computing the pose for features where fixed-point correspondences are not defined e.g. lines or planes. In such cases, features aligned with the coordinate axis simplify the construction of the equation system. For example, let us assume that two markers are placed to define an edge and it is parallel to the x-axis. This implies that movement is blocked along with the y-axis and y coordinates of the measured points can be used. The part is still free to move along the x-axis and markers can be placed anywhere along the edge the equations with y coordinate will always result in the same values because motion is along the principle direction.

Similarly, RPS features are also defined as planes and therefore markers can be placed anywhere along with the specified areas. In such cases, point correspondence is not possible for building an equation system. 3 points are sufficient for defining a plane and therefore their positions are used to define a plane in the device coordinate system. Let us assume that the plane is defined in the XY plane and blocks the movement of the part in direction of the normal i.e. along the z-axis. Therefore z-coordinate all the points on the plane is useful for building the equation system irrespective of their positions. In summary, constraints on each point are used to create a set of minimum 6 equations which is sufficient for pose computation in a unique manner. It should be noted the same method for building an equation system with 6 parameters holds for any alignment strategy working with 3-2-1.

### Advantages

The main advantage of using the RPS system is to win the trust of the user by providing improved workflow integration. We replicated the same method with optical tracking which not only improves accuracy but also increases reliability. RPS based maker placement allows us to avoid the geometric modeling problem while computing object pose. Mathematically, six

equations for the 6-DOF pose is a minimum condition and therefore results are bound to be consistent.

In terms of workflow integration, our method has an additional advantage that marker positions do not have to be specifically planned or tested by the engineer. At the same time, workers already know the concept of RPS and therefore marker placement at defined RPS features is rather intuitive for the workers. The previous approach required the markers to be placed at locations specified by the engineer. This approach was prone to placement errors and the potential cause of production failures. In summary, the new proposed approach reduces training time for workers as the workflow is simplified. Using trusted practices such as RPS increases the reliability of the system.

### Limitations

Our approach still has some limitations in terms of usability which are mainly due to the use of optical tracking and coded markers. The line of sight problem is omnipresent in optical tracking systems. The visibility of each marker is also important as the detection of each marker is important for the minimum solution based equation system (except the over-determined solution). Our solution does not handle cases when RPS features are not defined in principle planes. It is ideal to have automatic feature detection using the markerless method. We have proposed a solution in this direction in the next section.

## 4.1.2  Markerless alignment using Spatial AR

In this part of the chapter, we have proposed a resourceful solution to achieve end-to-end workflow integration with the minimum effort from the user end. In particular, we have focused on improving our existing solutions for alignment by adopting a markerless tracking strategy.

We have already shown that SAR solutions exhibit clear benefits for industrial users in terms of saving time, material, and enhancing productivity. However, the needs of the industrial environment in terms of usability and practicality remain partially unmet. The existing SAR solutions require additional work steps e.g. planning and placing markers. This additional effort is required from the user for accurate pose computation. We argue that markerless methods are a suitable alternative to marker-based solutions. However, the accuracy and reliability of markerless solutions are not proven in the industrial domain. The tracking results may vary between two departments, which is not appreciated in industrial workflows. Mostly because markerless methods are less customized for dealing with challenges of industrial environments such as complex part geometry, reflective or textureless metal surfaces. Additionally, feature detection varies due to light conditions which pose a big challenge for markerless tracking methods. Many markerless methods include an initialization step which requires moving the object in a specific perspective in front of the camera. The initialization is cumbersome and impractical for SAR systems having a bulky projector-camera (SAR) systems. In summary, we can conclude that markerless solutions for part alignment must perform with consistent accuracy and provide repeatable results without adding complicated and time-consuming work steps.

A car part projected with laser projector. The accuracy of the projection can be visually verified. A small error in pose, calibration or modeling error can lead to imperfect projections. It can be seen in the image that contours are lower part of the curve do not exactly project on the part. ©2016 IEEE [148]

Most markerless methods require prior knowledge of the part model (CAD-model) for pose computation. One major issue with this approach is the problem of geometric modeling is omnipresent in industrial manufacturing. Unknown deviations in part geometry affect both marker and markerless methods. Earlier we already proposed a strategy to use RPS features with a marker-based solution to eliminate errors stemming from geometric modeling. Now, we propose a new tracking strategy for achieving a markerless part alignment. We use the concept of laser scanning [127] in two parts: initialization and refinement. During initialization, we use the SAR device to project laser patterns on the part. A sparse point cloud is computed by triangulating the projected patterns. Reconstructed points are compared with the CAD model to compute an initial pose. Using the initial pose it is possible to project the laser patterns are specific locations on the part. Ideally, the new locations of projected patterns should be free from geometric modeling errors to ensure reliable results. One idea is to follow this strategy for pose refinement stage. Alternatively, we used an ICP based refinement strategy to improve initial pose estimation. We worked with industrial parts that are measured very accurately and do not suffer from geometric modeling. Therefore, we argue in favor of using the ICP algorithm for refinement in cases where the part is computed with relatively high accuracy. This approach is designed as a proof of concept of the proposed workflow.

In the future, we propose to use RPS locations for rescanning the object in the refinement stage. Markerless tracking combined with RPS features is a reliable strategy for seamless workflow integration. It will ensure both accuracy and reliability without imposing novel workflows on the user.

### Tracking workflow

In this part, we discuss the overall workflow and the details of the tracking method. In the initialization stage, the device projects a crosshair by default. The user has to point the device towards the object such that a projected cross-hair falls on the object. The position of the

**Fig. 4.4.** A car model with example projection on the wind shield. The model is covered with uncoded markers used in metrology. ©2016 IEEE [148]

crosshair is triangulated and this provides an estimate of the object's distance. Based on the distance and the size of the part a rough scanning volume is computed. The scan volume is limited to the object of interest to avoid generating unnecessary noise from the background. A diagonal grid pattern is used as the scanning patters as shown in Fig.4.5. The projected pattern of laser lines is detected in the cameras and reconstructed in three dimensions. Normals are estimated from the sparse point cloud [181] and feature vectors are created using the reconstructed points using PPF descriptors [62].

The CAD model of the part is used to compute feature vectors using the same strategy. This step can be performed as an off-line operation at the planning stage for saving time. The feature vectors from the CAD model and real-time construction are matched during the execution stage and a set of poses are suggested by the matching algorithm. The pose with the highest confidence level is selected as the initial pose.

The proposed workflow is free from any type of user introduced errors as marker placement is not required. The user only has to place the object in front of the device for initialization. However, the method relies a lot on the quality of the initial pose. The PPF descriptor-based matching has proven to be a robust method even if the scanned volume is smaller compared to whole the object size.

The success of descriptor matching relies on the quality of the normal estimation. The quality of the normal estimation technique depends on the density of the point cloud and its distribution. A dense point cloud results in better normal estimation but preparing dense point cloud requires time. The laser projectors are limited in terms of projection capability and therefore can not project many patterns simultaneously. The idea is to select a pattern that is dense enough to produce good normals within a reasonable time frame. We used a diagonal grid pattern as opposed to the classical vertical or horizontal scan pattern to get well distributed points in the volume. In our tests, diagonal patterns reduced the scan time by 30% without dropping the accuracy of the pose. Additionally, we can vary the intensity of the projections to suit the material of the part. This ensures that laser detection is robust in varying lighting conditions.

### Evaluation

We performed an evaluation using two industrial objects free from geometric modeling problems. The refinement stage uses the ICP algorithm and the initial pose is computed using the laser scanning approach. We use the same point cloud for both refinement and initialization stage. The refined pose is compared with the results from the classical marker tracking approach. The marker positions are computed with high precision metrology systems hence marker-based pose is used ground truth for comparison of results.

We used a car mock-up, Model 1, and a metal part, Model 2. For evaluation, we record the final pose and the overall time taken by the algorithm to converge to a solution. The Table 4.1 shows the results of our experiments. As we mentioned before the density of points is important for initialization and time. *Scan Density* column in the table represents the metric distance between the lines in the grid. The *Accuracy* column states the deviation in the translation and rotation parameters.

### Discussion

Our experiments show that a dense grid pattern significantly increases the time of scanning (as expected) but does not influence the accuracy significantly. We believe that the complexity of the part design also had to play a role, parts with more planar surfaces may require larger scanning volume. An option is to keep scanning density as a variable feature that the user can control. We observed that scanning consumes 80% of the overall workflow and pose

computation only consumes 20% of the time. The scanning time of our device is slower than commercial laser scanners because the hardware synchronization between the laser and the cameras is a limitation. This can be improved in the future with synchronization. It should be noted that the existing strategy with ICP refinement saves time because rescanning the object is not necessary. This strategy is only useful for objects produced with high confidence. In order to use RPS features, the object has to be rescanned and this would increase the scanning time.

Earlier we discussed the requirements of the IAR applications. It should be noted that the scanning based approach does not require is more or less a *plug and play* method of alignment. We give priority to the accuracy, reliability, and repeatability of the solution. In the future we expect the introduction of the RPS concept to add credibility to the results. It should be mentioned that the method performs poorly for purely planar surfaces and cylindrical symmetric bodies. Descriptors for such objects are not unique which leads to false initial pose and the ICP algorithm itself requires a variation of surface angles to perform a better alignment.

| Object | Scan Density *(mm)* | Time *(seconds)* | Accuracy (T) *(mm)* | Accuracy (R) *(deg)* |
|---|---|---|---|---|
| **Model 1** | 100 | 31 | 1.5 | 2.28 |
| | 60 | 48.8 | 0.67 | 2.45 |
| **Model 2** | 100 | 25 | 31.77 | 2.95 |
| | 60 | 37.7 | 9.84 | 2.14 |

**Tab. 4.1.** Time performance and evaluation against marker tracking.

## 4.2 Solving geometric modeling using closed Loop SAR

In this sections, we will tackle the problem of geometric modeling for the marking process. It is significant to solve the problem for the industrial partners willing to use SAR systems. We have discussed the problem earlier but we will discuss it again for broader understanding of the problem in context of SAR devices. Further, we will propose an iterative solution matching the existing industrial practices.

### 4.2.1 Background

Accurate 3D registration is a priority for IAR applications. Augmentation of trivial information such as texts or logos may not always require a high degree of registration accuracy. However, some industrial AR applications, such as point marking, do demand extremely high registration accuracy [17]. In SAR applications, usually a video projector [15] or a laser projector [188] is

**Fig. 4.6.** Perspective problem: The augmented point $A_i$ on the surface varies with perspective $P_i$

used for augmenting the object or scene of interest. In this case, the registration inaccuracies becomes more visible, as the augmentation is directly visualized the object's surface.

We mentioned earlier that three major factors affect the registration error in our SAR device: *Geometric modling*, *Calibration* (intrinsic and extrinsic) and *tracking accuracy (pose estimation)* [12]. In most cases, the object geometry is assumed to be correct and the registration errors are minimized by applying corrections to the estimated object pose or to the calibration parameters. However, in practice, especially in industrial AR scenarios, the assumption regarding the correctness of the geometry does not always hold. Multiple manufacturing processes introduce surface deformations at local level. Therefore the projected content on the "as is" model often differ from the planned pattern on the CAD model. In the case of small locally deformations the registration error only affects the deviated region. To solve this problem correction of the pose or the calibration parameters provide a further deteriorate the registration accuracy. Because these corrections globally influence the entire augmentation and not just the localized distorted regions. This is prohibitive for industrial SAR applications because erroneous augmentations are not suitable when accuracy and precision requirement in the range of sub millimetres. This motivates the development of a specific algorithm for solving the geometric assumption related registration problem in SAR. This aspect has not yet been addressed by the community, and is the primary focus of this research work.

We use three example to further elaborate the problem with figures 4.6,4.7 and Fig. 4.8. We have picked a typical use cases in the industrial scenario. To maintain simplicity we present the problems with single point projection method.

The differences between the manufactured *as is* object and the CAD model can be referred as local rigid deformations. These deviations on the surface tend to obstruct the path of the projected light ray either before or after its expected target position. As shown in Figure 4.6, the point may appear at different locations on the surface based on the angle of projection. The error magnitude $E$ is directly proportional to the magnitude of surface deviation $\Delta D$, the projection angle $\theta$ and the distance between projector and the part $D$. Another example

**Fig. 4.7.** Scale problem : An extension of perspective problem in context of two or more points.



**Fig. 4.8.** Application problem : Defined target geometry is occluded, need of adapting the projection.

in Fig.4.7 shows the scale error arising due to geometric modeling. In such cases, the augmentation visually appears correct but the relative distance of the augmented features changes significantly. Industrial applications for manufacturing assistance or quality check [139, 188, 220, 221] require projections at precisely measured distances to avoid part failure or assembly errors. This problem is specific to SAR solution because errors like this do not appear when manual jigs are used e.g. temple jig explained in sec. 3.2.1. The error becomes more significant when multiple parts (with tolerances) are assembled on top of each other. Example in Fig. 4.8 shows error in projection may lead to production failure. The geometric accuracy of the projection is a priority for assembly tasks.

An ideal SAR solution should focus on correcting the projections by accommodating the deformations to maintain the desired shape or scale. As a corrective measure new target position should be selected on the deviated surface such that projections converge to a unique

single point from any perspective. The existing application operated in an open loop fashion where pose is computed and projections are made. There is no provision for correction of the augmented points in real-time. We known that it is possible to visualize and compute the registration error using camera images. In AR, the concept of using the augmented information from the camera image as feed-back is also known as close loop AR [12]. Following this idea we propose a closed loop SAR approach to for registration correction in real-time. The core idea is to verify the registration accuracy of the augmented point by triangulating it's position and comparing it to the expected positions on the target object. Note that this idea works if we assume that pose and calibration parameters are accurate. Our approach assumes the deviations of the object geometry is the cause of registration error. We chose iterative process to adapt the projection such that registration error is minimized. To achieve this new target points are dynamically computed while keeping in mind that final projections must be free scale error. Additionally, our method also provides information about surface deviation which may be useful for the end user.

## 4.2.2  Related Work

The topic of geometric modeling is not well represented in the IAR literature. Therefore, we studied existing literature related to dynamic registration correction in AR. It is important to note that many methods are developed to improve registration accuracy of the augmented content. We will mainly discuss ideas in the context of registration correction and close loop AR.

The work of Holloway [90] illustrates the nature of registration error associated with generic AR applications. He explains that many factors play role in registration error such as hardware specific limitations, static errors (e.g. calibration, marker registration) and time delays. Bajura and Neumann[12] were the first to introduce the concept of dynamic registration correction using close loop approach for AR systems. They pointed out the three key sources of registration error, namely camera origin-to-model origin (pose), camera origin-to-image mapping (calibration) and origin-to-object (object geometry). They explained that it is difficult to determine the exact cause of error among the three error sources and thus proposed the idea of correcting for one factor at a time while assuming minimum error in the other two factors. The authors outlined two simplified strategies for registration correction HMD based AR application. Bajura and Neumann argued that close loop correction is a practical approach for solving the problem rather than trying to improve individual factors i.e. pose accuracy or calibration accuracy. So far most of the close loop AR methods have largely focused on correcting the pose or calibration parameters while assuming no error in the geometry of the model. One approach presented by Bajura *et al.* [12] focused on minimizing the registration error by measuring misalignment in 2D image space and applied correction to pose parameters. They used predefined features in the scene to identify registration error in the image space. This approach was specific to HMD based approaches and mostly applicable for 2D augmentations.

Zheng *et al.* [219] presented a generic close loop method to correct registration error for all forms of AR methods. The method encourages detection of the augmented features in the image and models them as part of tracking error minimization problem, then provides

correction for both image space rendering and pose computation. Zheng *et al.*presented another close loop approach specific to video AR [218], where the error minimization is carried out by giving weights to different image regions. In the former method, the focus is on correcting the tracking pose, whereas in the latter one Zheng aims at correcting both the tracking pose and the pixel-wise image rendering.

For SAR applications, Audet *et al.* [10] presented a projector based approach using the close loop concept to improve the registration (*alignment*) for 2D surfaces. The key focus was on modeling the light properties of the emitted projection and their reflection in the camera in order to improve alignment. Resch *et al.*[177] presented *Sticky Projections*, an approach to achieve interactive shader lamp tracking (in 3D space). The method optimizes the pose interactively by triangulating the projected features obtained as feedback in the image, followed by Iterative Closest Point (ICP) based pose optimization by comparing the transformation between computed geometry and known geometric model. The reader should note that a common assumption regarding the accuracy of the object geometry is made in all the methods mentioned above. DiVerdi and Höllerer [55] presented a hardware based methods for registration correction without modifying the pose or the calibration parameters. They pointed out inaccurate geometric modelling as one of the major sources of registration error in AR applications. The approach detects strong features of the objects (edges) in the image and corrects for the registration errors in the image space. The approach is independent of the tracking algorithm, however the scope of the work is limited to display based AR applications.

Bajura and Neumann [12] were the first to discuss the registration problem because of assumptions regarding accuracy of geometry. They proposed an experiment to simulate the problem by placing multiple objects in a reference space and one of them is not correctly positioned. They made correction in the misaligned position while keeping camera pose and calibration parameters constant. Bajura and Neumann argue that 3D measurements techniques are required to correct geometric misalignment and error in 3D alignment cannot be determined solely based on 2D image measurements. Our research is an attempt to extend this idea and provide the registration correction for errors due to inaccurate 3D model information.

On a parallel note, Menk [139] presented a method to evaluate the accuracy of augmentation for SAR applications with video projectors. Menk underlined the importance of computing accuracy of overlay (i.e. registration) in SAR applications rather than only focusing on tracking accuracy of the system. This argument holds strongly for industrial purposes (e.g. drilling, welding e.t.c) as the augmentation is eventually used to make manufacturing or assembly decisions with high accuracy (sub-millimeter ). After studying all methods we concluded that a close loop approach is necessary to verifying overly errors and correct them in dynamic manner. The problem of geometric modeling has become more relevant to increase acceptance of SAR systems, especially because the methods of achieving accurate tracking and calibration have evolved in the recent past.

## 4.2.3 Application and Assumptions

The presented problem is directly relevant to any application in industry involving marking process. We already presented example of drilling application where relative distances are important. We present another interesting application which involves unknown change of geometry. The objective is to fix two parts on top of each other as shown in the Figure.4.8. The target points are defined with respect to the part in the bottom, and the projection is expected to be on the surface of the part on top. This scenario requires changing the projection dynamically such that the correct positions are augmented. We show that such problems can also be modelled as geometric deviation or deformation problem as stated above.

We have made certain assumptions to single out the problem and focus purely on the relevant problem. The system is assumed to be calibrated in best possible manner and we assume that the pose is computed with high accuracy. We will focus on correcting registration error for a discrete set of point features (cross hairs) and not on a continuous geometric shape, i.e. line or circle. In our opinion the solution should be extended to complicated contours after understanding different aspects of the problem. Our choice is aligned with industrial processes and point features are used in many applications for spot welding or drilling [187]. Our method also assumes that information about the surface normals is available along with the 3D CAD geometry of the target object.



**Fig. 4.9.** Work flow of the algorithm, Verification and Correction module is added to identify registration misalignment and dynamically alter augmentations to improve registration.

## 4.2.4 Algorithm Pipeline

Now we will describe the software work flow designed to achieve dynamic registration correction, the components involved and our contribution to the existing AR work flow. The basic structure is inspired from the close loop concept proposed by Bajura and Neumann [12]. The pipeline mainly consists of three modules namely Tracking, Projection and Verification & Correction. The Tracking module involves two sub-modules, *Image Processing* and *Pose Computation*. We use a marker based tracking approach to have best possible accuracy. The marker positions, $O_m$, in the object space are predefined in the CAD model. Their positions in

the image space,$i_m$, are detected by image processing sub-module. The pose is computed ($P$) using $O_m \leftrightarrow i_m$ correspondences and it is supplied to the projection module. The projection module projects the designated virtual features $V_f$ on predefined target positions $O_f$ on the object. The target points ($O_f$) are selected on a digital model (e.g. CAD) and corresponding virtual features ($V_f$) are assigned for augmentation.

The above mentioned work flow is a generic SAR work flow, we have added the Verification & Correction module to achieve the closed loop functionality. The Image processing sub-module detects the projected features in the camera image. The following subsections give detailed account of the new module and forms the basis for the strategy applied to achieve better registration.

### Verification

In this sub-module, the position of projected features are reconstructed and verified against their planned positions. The input provided to this sub-module is the planned $O_f \leftrightarrow V_f$ relationship and the augmented features detected in the image, $i_f$. The image features $i_f$ are triangulated to compute points actual projected points, $T_p$. The correspondence $T_p \leftrightarrow V_f$ is known. It should be noted that triangulation is also possible using single projector-camera set-up, by considering projector as a camera [177]. The triangulated points are transformed from device space ($T_p$) to object space ($O_p$) using pose,

$$O_p = P^{-1}T_p. \tag{4.2}$$

Ideally, the triangulated features (in model space) $O_{p_j}$ should match coordinates of planned feature $O_{f_j}$, if the registration is accurate. Hence, the registration error $\delta_{R_j}$ at each point can be computed as,

$$\delta_{R_j} = \begin{cases} 0, & if\ Dist(O_{p_j} - O_{f_j}) < \delta_{th} \\ 1, & Otherwise. \end{cases} \tag{4.3}$$

where, $\delta_{th}$ is the threshold to identify registration alignment status $\delta R_j$ for each point $j$. The status is assigned to 1 (meaning registration error) for each points that do not match the threshold criteria. Geometric modeling error may not affect all augmentations, therefore each feature is compared individually. Registration status of each point is checked and misaligned points are passes on to the Correction sub-module for further processing. We also combine individual registration errors to compute a quantitative measure for describing overall registration accuracy. For the end user qualitative measure for registration is more meaningful and intuitive for decision making. This is an improvement to the existing SAR systems which can only provide tracking errors.

### Correction

In this sub-module, a correction strategy applied to the non-verified or misaligned points to regain accurate registration. The registration error $E$ at each point may varies with the projector distance $D$, the viewing angle $\theta$ and the deformation ($\Delta D$) (see Fig. 4.6 & 4.10). Therefore, each feature is dealt with individually to apply the registration correction. The goal is to find a new point on the object surface which can be assigned to the corresponding virtual feature. The new point should not suffer from perspective problem and scale problem.

**Fig. 4.10.** Pictorial depiction of the iterative correction strategy used for computing new target points for augmentation.

The problem becomes very complex at this stage, as to which direction should be selected for applying the correction, since this mostly depends on the specific application.

We followed the existing practices in the industry and learned that most times correction is applied along the normal vector. A different variation is to provide correction along a certain vector component of the normal that is parallel to one of the principle axis. The readers should note that RPS locations are defined along the principle axis for same reasons. Typically, workers use optical or tactile measurement systems and have to compute the new point on their own. For example, if X-axis is chosen as correction direction than the worker try to find a point at (*,y,z) coordinate on the surface. The choice of direction corrections is very specific and it is a decision made by design engineers. We have replicated this process conceptually and applied correction along the chosen direction. We have taken the default strategy and corrected projections along the surface normals.

The CAD model is used for the SAR applications and therefore it is reasonable to assume that 3D positions of the target points ($O_f$) is available with the corresponding surface normals $N_f$. Our applications proposes a new point on the part surface that is along the normal direction of the originally planned point. This correction strategy along the normal is supported by the argument for maintaining geometric correctness for several applications (cf. 4.7). Our process of determining new target point is iterative because the deformation of the surface is unknown.

Figure 4.10 illustrates a graphical example of the iterative solution implemented to find the expected target point. Let a point $P$ be the planned target point, $P'$ be the triangulated feature point and $\hat{P}$ be the expected target point (unknown). Let $\vec{P_N}$ be the surface normal at the planned target point $P$. The direction of the correction is along the normal and the magnitude of correction $C_i$, the distance between $P$ and $\hat{P}$, is to be determined. The corrected target

point is computed by adding correction factor to the planned target point ($P$) along its normal. The estimated correction factor is projection of vector $\vec{PP'}$ along the normal $\vec{P_N}$.

$$P_i = \vec{P_N}C_i + P \tag{4.4}$$

$$C_i = \vec{P_N} \cdot \vec{P_iP_i'} \tag{4.5}$$

The corrected target point $P_i$ is supplied to the projection module, where $i$ represents the iteration number. The new projections are triangulated as $P_i'$, verified against $P_i$ and a new correction factor $C_{i+1}$ is computed. This procedure is repeated until the stop condition ($\gamma = 0$) is reached or the minimum possible $\gamma$ is achieved for a given feature. $\gamma$ is the angle between the normal $\vec{P_N}$ and the vector $\vec{PP_i'}$. Ideally, the corrected target $P_i$ and the triangulated point $P_i'$ (i.e. projected point) should result in the same value. However, this could be true for any false positive point on the surface if the pose is kept constant, therefore we aim to minimize $\gamma$ and keep $\gamma = 0$ as stop condition for the iterations.
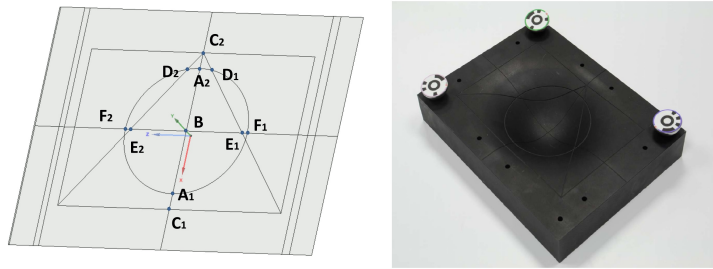
$$min(\gamma) = \angle(\vec{PP_i'}, \vec{P_N}) \tag{4.6}$$

The deviation between converged corrected point and the expected $\hat{P}$ is the final registration error. However, this cannot be computed algorithmically because $\hat{P}$ is not known. This correction strategy is applied to all the identified misaligned points and a set of new correspondences is formed for augmentation $O_f' \leftrightarrow V_f$ and supplied to the projection module.

The ideal value of $\gamma$ should be 0 for perfect convergence. In a practical implementation getting exact projection is difficult and therefore the solution may not stabilize a specific local minimum. For such conditional, acceptance criterion can be considered where new value of $\gamma_i$ must be smaller than previous one i.e. $\gamma_{i+1} < \gamma_i$. This way the correction factor is accepted only if $\gamma$ is minimized in the next iteration. This condition prevents the occasional fluctuation in projections. These fluctuations can be caused by triangulation errors or image processing errors in cross hair detection. Our method is able to perform with all type of deviations (i.e. concave, convex or linear) with the same principle. The direction of the normal has to be flipped if the direction of the vector $\vec{P_N}$ and the vector $\vec{PP'}$ are opposite, to maintain minimization problem. This would suggest inward deformation of the surface (i.e. away from the direction normal).

## 4.2.5 Experiments and Evaluation

We performed multiple experiments to test the effectiveness of our registration correction algorithm. The first experiment evaluates the performance of the correction method on an object using a known 3D geometry as ground truth. This experiment focuses on the perspective problem highlighted in the introduction. The second experiment focuses on the scale problem, also referred to in the introduction. Multiple features are projected with erroneous 3D target points to simulate scale error condition. The accuracy related numbers are discussed and important factors are highlighted. We also show an application example related to the concept of visualizing hidden geometry. All the error measurements are shown in millimeters (mm) and distances are shown in meters (m). The readers are reminded that dynamic registration problem w.r.t errors in geometry is not address before and therefore lacks comparative studies with other methods. We use a stereo camera system with a laser projector mounted together.
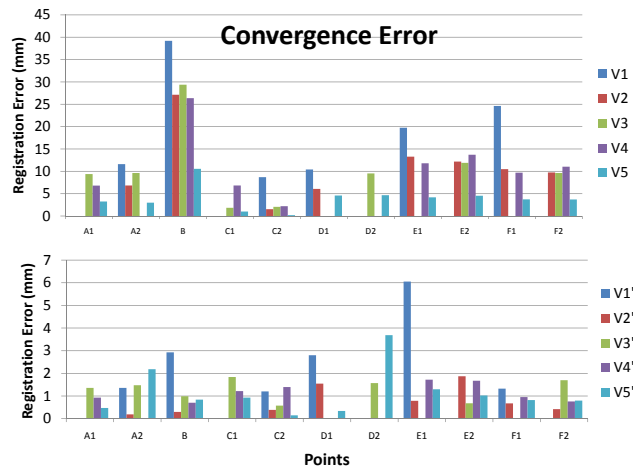
**Fig. 4.11.** A physical model with known deviations is used to compare results after correction against ground truth. (left) a planar model of the planner geometry (right) actual model with deviated geometry. *[Check with E3D]*. ©IEEE 2015, reprinted from [151].

The camera resolution of the system is 3840x2748 and tracking speed is 6 FPS due to hardware limitations of industrial grade cameras. The algorithms are developed using OpenCV and MATLAB.

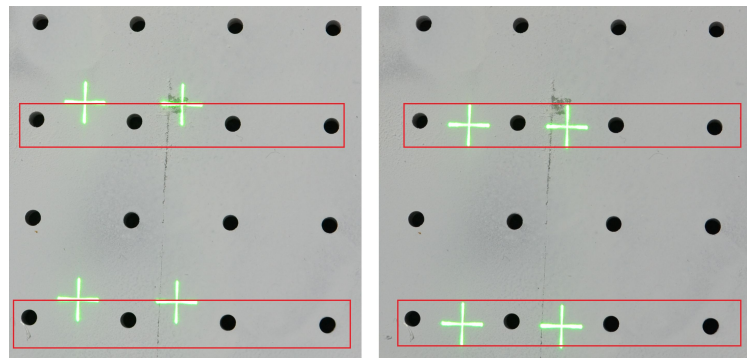### Registration correction with respect to perspective inaccuracy

A 3D object with known geometry is used for the evaluation of the proposed registration correction method. The features A' – F' on the body of the object are projected along the normal on a virtual plane as A – F (cf. 4.11 ) in digital CAD format. The points A' –F' are used as ground truth deviations of corresponding points A –F. Features A – F on the plane are considered as target points for projection planning. We move the set up around the object to project a simple cross hair feature on target points from different perspectives (V1 – V5 ). The projector distances ($D_{1-5}$) are $2, 1.4, 1.5, 1.8, 1.2$ meters and view angles ($\theta_{1-5}$) are $50, 40, 40, 35, 18$ degrees for the corresponding view. The same experiment is conducted without correction (V1 – V5) and with correction (V1' – V5') for each view and the projected positions are recorded. The ground truth positions are known, thus registration error before and after registration correction is computed to highlight the obtained gain in accuracy. Deformation values for points A –F are $11.5, 39.9, 3.61, 11.76, 18.02, 15.50$ in mm ($A1 = A2$). The goal is to validate the accuracy of the algorithm to converge to the expected point from different perspectives without any prior information of the surface deviations. Moreover, the deviation from planned target point is computed by the algorithm is also validated against the ground truth.

The graph in Figure 4.12 show that the registration error varies randomly depending the perspective (V1 – V5) and the deviation from the planned point. This behaviour highlights the perspective problem, which we want to solve. The deviation in the registration error is limited from mm to sub-mm range with dynamic registration correction. At lower viewing angles (e.g. V'5 = 18°) the registration accuracy after the correction is sub-mm ( < mm) even with larger deformations like point B ($\Delta D$ = 39 mm). The performance of the correction algorithm is less stable at higher distance (2 m) and at higher viewing angle (50°) for example V1. The points with low deviations (e.g. $C$ = 3 mm ) show highly stable convergence after the correction. The measurement value is kept zero for points that are not visible in the respective view. Moreover, the graph implicitly shows the accuracy of obtaining the deviation information. Accurate convergence means better registration and better estimation of deviation of the

**Fig. 4.12.** The graph shows the ability of the correction algorithm to converge to a unique point on the surface with unknown deformation independent of the perspective ($V$). The deviation between the actual projected point and the expected point before (Top) and after the correction (Bottom). ©IEEE 2015, reprinted from [151].

surface from planned point. The results show that the correction algorithm is able to minimize the registration error up to sub-mm between the expected point and the projected points.



**Fig. 4.13.** The images show misalignment of the features before applying the registration corrections (left) and correct alignment after (right) applying the registration correction.

### Correction for scale errors

In this experiment, we project multiple cross hair features on a metal plate with predefined relative distances. We aim to study the effect of correction at the level of individual points and the corresponding distances among the points at different deviations. The plate is planar, 3D positions of the target points is planned and the Z information is manipulated to simulate different height deviations for the target plane. Figure 4.13 highlights the misalignment between the holes and projections due error included in the target geometry. The point positions are chosen such that viewing angle ($\theta$) variation $\approx 45\pm 3°$ is minimized. The error $\Delta D$ in the height (Z dimension) of given target geometry is sequentially increased from 1 mm to 15 mm and projected positions are recorded. The pose of the object is kept constant throughout the whole set of experiment. In real manufacturing scenario manufacturing tolerances are expected between ($\mu m$ to $\pm 2 - 3\ mm$), therefore we are specifically interested

**Fig. 4.14.** Correction applied at each point with respect to subjected deformation.



**Fig. 4.15.** Comparison between corrected (top) and non-corrected (bottom) distances at different deformations.

in correcting for smaller errors. The projections are expected to be distorted both in terms of relative distances and the individual X-Y positions. We apply the correction method and records the corrected X-Y positions and the relative scale dimensions. A subset of the measured distances are AB (100 mm), BC (111.8 mm), EF (50 mm) and EH (150 mm).

The graph in Figure 4.14 shows the correction applied to each point in response to the deformation. The registration correction is almost equal to the deviation because the relation $Error = \tan\theta \times \Delta D$ holds true if the distorted surface ($\Delta D$) is planar. The results show that each point requires an individual correction even when subject to the same deformation values due to the slight variation in angle projected ray and projector-to-target point distance. This reinforces our argument that each feature should be treated individually for high accuracy registration corrections. The error in feature position (X-Y) after the correction is in sub-mm ($< 1mm$) range as shown in Figure 4.15. It is important to note that accurate position

**Fig. 4.16.** The targets defined for projection are on the lower metal plate. The image on the left shows three misaligned cross hair features, the one on the right shows aligned feature after applying registration correction.

correction implies preservation of scale but not vice-versa. This phenomenon is indicated in Figure.4.15, where the disturbance in the relative point distances is measured. The results show that the scale error remains in sub-mm range (with wrong position) for all small deviations (1-8 mm) even before the registration correction. This can be explained as the deviation for all points is same and in the same direction, therefore the positions deteriorate faster than the scale. The impact of registration correction is observed at higher deformations (e.g. 15 mm) (see Fig.4.7) where scale error is higher without the correction. The scale may deteriorate further when points are subjected to different deformation or point distances are higher. In conclusion, the experiment shows that the correction algorithm restores both the scale and the position of the projections.

### Application : Projecting hidden geometry

The aim of this experiment is to show that hidden features can be highlighted accurately using our concept. It can be occluded drilling points or planned welding spots below the assembly part (cf. 4.8). Conceptually, the experiment is similar to above mentioned grid experiment. We physically occlude the project on metal plate with another plate, the goal is to achieve the projection aligned with the holes but on the added plate automatically. The system can dynamically detect the registration misalignment and iteratively adapt the projection on the new surface. The time taken for computing the correction is 90 milliseconds as compared to pose computation time of 200 milliseconds. The computation is time is slow due to hardware limitations related to the cameras. Additionally, the high resolution images require more time for undistortion and image processing. The applications proposed are rather static in nature and fast moving target object can not be support with current implementation.

## 4.3 Conclusion and Future Directions

In the presented research we focused on two major applications of IAR i.e. alignment and marking. We proposed SAR devices as alternative to existing manual practices. We learned that industrial users accept SAR solutions only if their requirement regarding performance and workflow integration are satisfied. In this direction, we worked with the industrial partners to gain deeper understanding of their requirements and existing practices.

We learned working principle behind the existing methods of alignment and marking. Furthermore, we designed new approaches to use working principles of the existing methods with projection based SAR devices. More specifically, we proposed two tracking strategies one with markers and other without markers to achieve reliable alignment. Both solutions are based on the existing protocol for manual alignment i.e. RPS. We show that our solution is accuracy and reliable. We also worked on the problem of accurate point marking while dealing with issues such as geometric modeling. This problem is not addressed before in SAR literature and we designed a solution that is able to mimic the existing practices.

Our methods still have some limitations which lay the ground work for future work in this direction. We need a better approach to tackle RPS features when they are not defined along the principle axis. We argue that marker placement can be avoided in future completely with better markerless methods. Introduction of machine learning techniques can boost the performance of the system. The closed loop registration correction concept is currently limited to point wise feature correction which is sufficient for multiple industrial applications. However, the future application must involve registration correction for complicated geometries such as contours. Both software and hardware changes are required to implement such features. In terms of software better image processing techniques are required to increase accuracy of triangulations. Hardware changes include updated synchronization between the camera and the projector system. This would allow pictures to be taken in sync with projections, which will save time of scanning.

The goal of the project was to design AR methods to replace existing industrial processes. Our methods were well received by the industrial partners after the conclusion of the ARVIDA project [184]. These solutions are now integrated in the new commercial SAR devices of Extend3D GmbH. We believe that our findings are applicable to other application domains of SAR where used want to solve problems of geometric modeling e.g. in medicine to compensate movement of patients.

# Part III

XR Applications For Animal Behavior
Studies

# Animal Behavior Studies with Artificial Stimulation

<div style="text-align: right">5</div>

In this chapter, we will focus on applications of sensory stimulation techniques in the field of animal behavior, especially using virtual environments (VE). Biologists have been using closed-loop visual stimulation techniques for almost two decades. More recently, Virtual Reality setups (based on CAVE concept) have been designed to study fish, mammals and insects [195, 205]. However, the existing technology can only be used with a limited range of species in a smaller area. These experiments have now reached a stage where sophisticated technology is required to overcome existing limitations.

Researches at Max Planck Institute of Animal Behavior (MPI-AB) and Center of Advanced Collective Behavior Studies (CACBS) at the University of Konstanz have taken an initiative to design a unique facility to conduct collective behavior experiments using novel techniques for sensory stimulation. Researchers from biology and computer science are working together on this initiative. The author has joined this initiative as an XR and computer vision researcher at the Max Planck Institute of Animal Behavior (MPI-AB) [1]. The research is supported by MPI-AB and DFG Centre of Excellence 2117 Center for the Advanced Study of Collective Behaviour and University of Konstanz (ID: 422037984). The author was involved setting up a tracking strategy to capture the 3D movement of animal groups for both open-loop and closed-loop experiments. Specific focus was given towards designing a strategy for tracking 3D posture of birds for conducting perspective dependent visual stimulation experiments using the concepts from Augmented and Mixed Reality.

This chapter is part of our scientific contribution towards building the setup mentioned above. Closed-loop behavior experiments with virtual environments are a special use case of XR technology. These applications are designed by humans but the stimuli are exposed to the animals in a closed loop manner. Novel solutions can not be designed without understanding of the technical requirements of these experiment. At the same time, the technology must be customized to support the animal's sensory capabilities. Such discussions on behavior experiments are not reported in the technical literature of the XR community. we prepared a literature review on sensory stimulation based animal behavior studies to gain support of technology experts for building our new facility (mentioned above). We have highlighted major limitations of existing methods and the provided several ideas for designing new solutions.

*The review is published in IEEE Transactions on Visualization and Computer Graphics [149] with the title of "Animals in Virtual Environments" and it is also published in the proceedings of the IEEE Virtual Reality conference held in 2020 at Atlanta. The readers should note that the text*

---

[1]Previously known as Max Planck Institute of Ornithology

*presented in this chapter is largely based on the publication "Animals in Virtual Environments"*
*[149].* [2].

## 5.1 Animal Behavior and stimulus based experiments

A wide range of scientific disciplines use animals as a primary subjects of study e.g. medicine, neurobiology, physiology. Ethology, the field of animal behavior, is largely concerned with understanding why animals do what they do, and how. Animals exhibit behavioral strategies that have evolved to enhance its survival in the natural environment (land, air or underwater). Each animal's behavioral interactions with its own environment, and other organisms, reveals important information about ecology and evolution. Humans have studied animal behavior for hundreds of years including during domestication. In 1963, Niko Tinbergen suggested the first framework for studying behavior in form of four fundamental questions [206]; What is the survival value of the behavior? How does the behavior develop during the lifetime of the animal? How did the behavior evolve across generations? And how does it work (mechanism)? His objective was to propose a framework that defines the scope of the scientific study of behavior. It is widely accepted among behavior researchers that a comprehensive understanding of behavior can be obtained from following Tinbergen's framework [16, 206].

Different aspects of animal behavior have been studied over the last 60 years. In neuroscience, neural activity of behaving animals is recorded to find the link between sensory-motor mechanisms and neural processing [194]. The genetic basis of behavior can be studied by observing behavior in genetically manipulated animals (such as mutants). In medicine, small vertebrates (fish or mice) are preferred because they exhibit some fundamental behavioral traits that are consistent with other vertebrates, including humans. Their behavior can be closely monitored during experimental drug trials to study the progression of the disease and to measure the resulting effect on the animal's behavior [53]. Revealing the behavioral strategies of the animals is useful for solving problems in the fields of engineering and technology. For example, behavior of animals has been studied for various applications in robotics [7, 97, 202]. Biologists and engineers have benefited greatly by working together on novel interdisciplinary projects where robots are used to investigate the principles of decision making in animals [112, 120, 213](details in Sec 5.7).

Animal behavior is studied using various experimental methods. Behavior is investigated in both indoor (lab, cage) and outdoor (wild, open area) environments depending on the research questions. Outdoor environments are more suitable for the observation of realistic behavioral patterns. However, experiments in natural environments can be time-consuming and expensive. Outdoor experiments are also prone to unplanned disturbances from external factors which may alter the behavior during the experiment e.g. weather conditions, human disturbances. Indoor environments provide more control over the experimental conditions and minimize the influence of external factors. It is thus easier to develop standardized and repeatable methods for such behavioral experiments. Indoor environments are suitable for

---

[2]The publication is available as open access under CC-BY 4.0 license. The lead author of this dissertation is the copyright owner of the published text along with co-authors.

carrying out detailed studies, but the range of behaviors displayed in such environments may be limited. Many wild animals do not exhibit natural behaviors in indoor environments. Some species, termed as *model species*, are preferred as they exhibit an ability to perform naturalistic behaviors in indoor environments e.g. zebra fish, fruit flies etc.

Artificial sensory stimuli are often used in experiments to invoke behavioral responses from animals. In natural conditions, animals constantly receive sensory stimulus (visual, auditory, haptic etc.) from their environment and must react to it appropriately. Stimuli are often designed artificially to mimic natural conditions. For example, temperature and light manipulation is sufficient to artificially simulate day and night cycle for insects and birds. Niko Tinbergen used cardboard models of adult gulls to invoke begging behavior in gull chicks [207]. Artificial stimulation is a powerful technique for achieving repeatable behavioral observations [37]. The experimenter can plan the timing of stimulus delivery and change properties of the stimulus between different trials to observe changes in behavior. Such experiments provide a deeper understanding of the decision-making of animals. After Tinbergen's initial findings, more advanced techniques were developed to stimulate sensory systems of different animals for behavioral experiments. Technological innovations such as cameras, speakers and projectors have made a major impact in behavioral studies. They are used in novel ways to manipulate the information received by the animal about its surroundings environment (e.g. audio or visual stimulation).

In the late 90s, researchers studying human behavior, psychology, and perception started exploring the Virtual Environments (VE) as a tool for manipulating the human perception of reality by artificial stimulation of human sensing [201]. The concept of CAVE VR [45] was introduced with the idea of creating an immersive experience for the viewer by means of visual stimulation. The head position of viewers is tracked in 3D and the stimuli are always rendered on the walls from the perspective of the viewer (for details see sec.2.3.2). This approach made it possible to create and maintain the illusion in real-time without wearing cumbersome HMDs. Around the same time biologists had shown that the method of displaying virtual stimuli on a screen was useful for studying behavior but it was limited due to lack of interactivity. Biologists started adopting the CAVE VR design and introduced the concept of interactive virtual environments for animals, almost two decades ago [83, 185]. Their goal was to design a novel experimental approach where the animal behaves as if freely navigating in its natural environment. Since then many techniques have been developed for studying behavior of freely moving animals (fish, mammals and insects) in the virtual environment e.g. FreemoVR [195], FlyVR [194]. Animal VR systems can thus be considered as a cleverly modified version of human VR systems.

The technology used for designing VE for animals is similar to that used for designing XR applications for human users (see sec. 2.4). However, the sensory perception of animals is different to that of humans which means that they may sense the environment in a different manner e.g. UV vision in birds, ultrasonic hearing in bats. The methods developed for stimulating humans may therefore only be suitable for some animals. Virtual environments for animals are limited by our ability to produce the technology that matches the animal's sensory input. This being said biologists have shown that circumventing some limitations is possible by adopting new methods, for example developments of real-time tracking and realistic graphic rendering was crucial for developing VEs. We argue that stronger research collaborations

between biologists and the technology developers from the XR community must be promoted to push the research further. In the following text, we trace the journey of stimulus-based behavior experiments from simple non-interactive models to fully interactive VR systems.
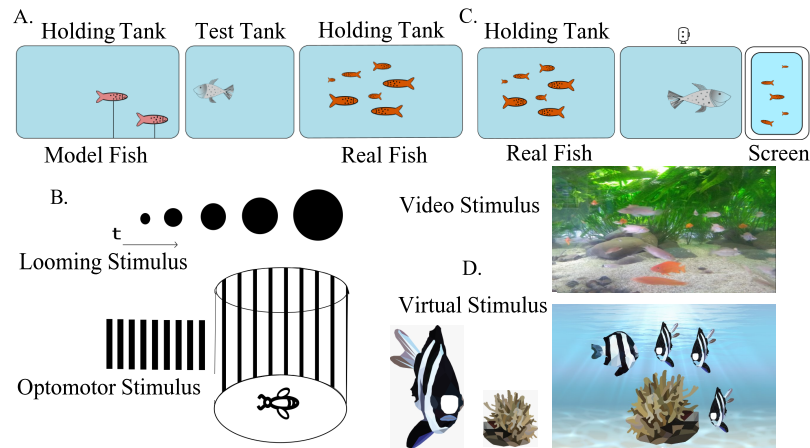
## 5.2  Review Method

We have collected the scientific literature from different research domains associated with the study of animal behavior e.g. ethology, neurology, psychology. We started from other review papers that focused on virtual environments and virtual stimuli for animals [31, 56, 58, 123, 194, 205, 216]. These publications are informative and extensive but they are prepared for readers with backgrounds in biology and behavioral experiments. Often the experiments are discussed with specific focus on an application (e.g. neurobiology [31, 58]) or a type of animal (e.g. rodents [205]). Therefore, such reviews are customized for biologists rather than technology developers.

Our review provides a more general overview of the methods that is suitable for technical experts. We also studied papers with a strong emphasis on the limitations of using virtual stimuli with animals [41, 194, 205] and collected literature on non-interactive methods for artificial stimulation [46, 216]. These methods are commonly used for studying behavior and their success is one of the strongest arguments in favor of developing virtual environments for studying behavior. Overall our review is designed to serve as a guide for readers of the computer science community (especially XR) who wish to later explore more detailed literature in the field of behavioral studies. We have mainly covered experiments that use visual stimuli since these represent the largest and most diverse body of work to date in stimulation based behavior studies.

**Tab. 5.1.**  Overview of artificial visual stimuli used in open loop experiments

| Type | Stimulus | Animal-Behavior | Key Attributes |
|---|---|---|---|
| Static | Model, Image, Color Filter, Paint conspecifics | Birds - Feeding[207], Vigilance [66], Mate choice [18], Social hierarchy [179]. Fish - Mate choice [141] | Configurable properties: shape, size etc., reusable method, non-interactive |
| Abstract | Patterns with points, lines, circles | Fruit fly - Perception and navigation [185], Movement and physiology of eye [32, 119], Motion control [76], Trajectory correction [204]. Locust - Motion parallax [192], Insect locomotion [202]. Moth - Navigation [83]. Mice - OMR [1] | Setup can be mechanical design with pattern cylinder or with screens, popular for studying vision induced motion e.g. OMR, OKR [1, 113, 119] |
| Video | Video recording displayed through screen or projector | Lizard - Courtship [98], Communication [165] Jumping spider - Recognition [42], Birds - Alarm calling [66], Fish - Laterality and cooperation [28]. Review - [46, 156]. | Stimulus can be edited and customized, reusable setup for multiple behavior experiments, can display abstract stimuli, non-interactive. |
| Virtual | Computer generated content through projector or screen | Fish - Mate preference [14, 116], Predator response [80, 94], Communication [89]. Review - [34, 216] | Stimulus programmable, reusable setup, semi-interactive or rule based interaction, can display abstract stimulus. |

**Fig. 5.1.** Common type of visual stimuli used in open-loop behavior experiments. A. Static stimuli: An experimental setup designed to study preference of the fish in the test tank when given a choice between model fish and real fish, B. Abstract stimuli: Common type of patterns used to study vision induced motion and pictorial depiction of mechanically controlled Pattern Cylinder (PC) with striped patterns, C. Video stimuli: Experimental setup for preference study with video playback experiment, similar to figure A., the screen is used to display videos fish as shown in example image, D. Virtual stimuli: Example of an image displayed as virtual stimulus. It consists of some fish around a coral, all of which are separately designed graphical models (left) combined with an underwater background. This stimulus is displayed using a screen similar to the setup shown in figure C.

# 5.3  Artificial Stimuli and Open-Loop Experiments

In this section, we aim to introduce the readers with a brief history of stimulus design and open-loop experiments. This section is included for readers from non-biological backgrounds to understand how biologists came to use virtual environments for studying behavior of animals. First, we explain the fundamental ideas behind using artificial sensory stimuli and then relate these ideas to the framework of studying behavior [206]. We cover the four different categories of visual stimuli that are commonly used for behavioral experiments: static stimuli, abstract stimuli, video stimuli and virtual stimuli. Each category represents one or more types of visual stimuli which share some common properties in terms of design and/or the method used to present them. We will describe the intuition behind designing each type and with suitable examples we will show how these stimulus were implemented. In the end, the knowledge gained from these experiments is summarised. Review of these methods has helped authors understand nuances of behavior studies for designing the 3D tracking facility presented in chapter 6

It is difficult to design interactive stimuli without understanding which stimuli to use and how to present them. Every novel idea of visual stimulation was first tested with an open-loop (non-interactive) approach. The stimulus does not change or react to the animal in open-loop methods. However, the method is sufficient to check the feasibility of using a stimuli and to learn the right technique for displaying it. Experiments with the open-loop approach are still being used in behavioral studies and they have provided much of the fundamental understanding required to build the modern closed-loop behavior experiments.

## 5.3.1 Static stimuli

Static stimuli were some of the first employed in behavioral experiments. Here animals are presented with a static object such as a model or an image, usually of an animal (figure 5.1A). We refer to this type as static because the properties of stimuli do not alter or change during the experiment. It was hypothesized that animals may perceive a static model as a real animal and react to it. It was shown that in some cases such stimuli were sufficient to invoke a response such as fear or attraction. Tinbergen and Perdeck [207] used a model of a bird to invoke begging behavior in chicks of gulls. The chicks responded naturally to the models i.e. as if begging for food from a parent. Following such studies, other researchers tested artificial objects extensively using different variants, which differed both in terms of visual properties of the model (e.g. different colors, shapes), and timing of stimulus delivery (e.g. time of the day or frequency) [46]. Evans and Marler [66] studied alarm behavior in chickens using a model of a predator. They created a setup where a model of a predatory bird would move above the cage on a rope. This simulated a typical behavior of a predatory bird gliding in the sky looking for food. It was observed that chickens made alarm calls when the model was moving over the cage. Images, photographs and slides were also used as static stimuli [46]. Other examples of static stimuli are environmental modifications e.g. light filters, which allow or reject of specific wavelengths [18, 141], and visual modification e.g. painting conspecifics [179]. Robotic animals are increasingly used as visual stimuli [108, 112, 120, 213], but to maintain focus on virtual visual stimuli we do not cover them in this thesis.

Static stimuli have proved to be reliable and repeatable means for conducting behavioral experiments. The main advantage is that the same stimulus are used for different animals and its visual properties are modified between trials. Moreover, the timing of stimulus delivery and frequency of displaying the stimulus are also controlled [46]. A major limitation to this approach is that there is no scope for feedback between the stimuli and the organisms. Consequently a problem that arises is that individuals can habituate to, and stop responding to, stimuli over time [46].

## 5.3.2 Abstract stimuli

The primary intuition for the design of abstract stimuli was to design a stimulus that is minimalistic yet sufficient to drive behavioral decision-making process in animals (such as movement). These are the most widely used stimuli for behavior experiments, especially for studying mechanisms related to visually-induced locomotion. For example, a common abstract stimulus consists of simple patterns designed from primitive geometric shapes such as points and lines (see Fig. 5.1 B) e.g. stripes or circles. The idea is to measure the movement of the animal in response to the patterns displayed to it. The mapping between the features observed by the animal and the animal's movements reveal the underlying process of behavioral decision making in the context of such stimuli. This experimental concept is designed to investigate fundamental questions regarding the behavioral and neural basis of visually-induced locomotion. Small invertebrates, such as fruit flies or honey bees, and relatively simple vertebrates such as fish, are the preferred model species as they possess less complex nervous systems and relatively fewer behavior patterns.

One example of visually-induced locomotion behavior is the optomotor response (OMR), which is the property of moving the body and/or head in concert with the features in the environment for image stabilization [113, 119]. Similarly, the property of moving the eyes in concert moving features is called optokinetic response (OKR). One of the first experimental setups for studying OMR and OKR was a mechanically controlled stimulus delivery system. It consisted of a stage for placing the animal and a cylindrical drum was designed to be rotated around surrounding the stage using a motor. Its inner walls were painted with abstract patterns (stimulus) cf. figure 5.1B. The movements of the animals are restricted to a small area and sometimes tethering is used to keep the animals fixed in one spot which simplified measurement of head or eye movements. The movement of the animal is recorded using video cameras or using a simple array of photodiodes [83] or torque motor [79]. It is shown that the animals typically display a tendency to move in the direction of the rotation. The width of the striped pattern, the rotation speed of the cylinder, and the direction of the rotation are also influential in decision making [1].

Another example of visually induced motion is avoidance of or flight response to, a rapidly-expanding (or 'looming') shape on the retina, typically a black expanding circle with a white background as shown in figure 5.1B. In early designs, looming stimuli were simulated by mechanically moving a dark circular cardboard cutout closer to the animal. Electronically controlled display methods (e.g. LED grid, LCD) have replaced the mechanical methods for displaying the stimulus and computer vision techniques are now deployed for movement measurements [185, 197]. It is also possible to conduct abstract stimuli-based experiments in a fully automated closed-loop manner [75, 76, 194] (covered later in virtual stimuli).

Experiments with abstract stimuli are relatively easy to design and their results tend to be reproducible and the method has opened doors for reverse-engineering the process of visually induced motion. Abstract stimuli have been successfully used for more than 70 years for research studies by biologists and engineers alike. The experiments have been conducted exhaustively with small insects (e.g. fruit flies, bees) and vertebrates (e.g. zebrafish, mice). Some examples include studying the vision properties such as depth perception [185], motion parallax [192], movement of the eye [119], physiology of the eye [32], locomotion mechanisms such as flight behavior [202] and motion control [76], and behavioral patterns like navigation [83, 185], and trajectory correction [204] (see table. 5.1). Mechanical designs have suffered from mechanical limitations such as motor speed or requirements of sticking patterns on cylinder. These modifications require manual interventions, which is not suitable for conducting high-throughput experiments. Another limitations is that temporal resolutions of measurements is low.

### 5.3.3 Video stimuli

This category includes experiments in which a recorded video footage is used as visual stimuli. These experiments are known as video playback experiments in the behavior literature. For experiments investigating social behavior, video cameras are typically used to record the activity of an animal and this is then used as a stimulus during the experiment [46, 156]. The focal animal is usually placed in an enclosure or an arena and is shown a video sequence through a screen or projector placed at a reasonable distance (figure 5.1 C). The responses of

this animal is recorded using a video camera, which are later used to map behavioral decisions of the animals to the visuals presented in the stimuli. The stimuli may involve another animal of the same species (conspecific) or a different species (heterospecific) behaving in a certain way. It was hypothesized that animals do not comprehend the concept of video screens and thus will react in a natural and instinctive manner to the presented stimuli. The intuition behind this method is to simulate the natural environment in the photo realistic way [46].

Jenssen [98] designed a mate-choice experiment with female lizards using video playback method. He displayed video sequences of male lizards performing courtship display and reported that female lizards did react appropriately to such stimuli. The video playback technique is commonly used method for studying a wide range of behavior responses such as: aggression, attraction, fear etc, in species including arachnids [42], birds[66], reptiles [165], and fishes [28]. It is a reliable technique for quantitative analysis of behavior. The movements are measured in 2D or 3D space using multiple cameras [170]. The video playback methods have provided many insights into the questions related to the survival value of a specific behavior.

Video stimuli have some clear advantages because a customized sequence of behaviors can be shown. The same setup can be used to display wide range of behaviors; for example the setup in figure 5.1C can be used for studying either mate preference or aggression. Camera and display technologies required for the experiment are typically commercially available which makes this method accessible to researchers. Various software tools are designed to quantify behavioral response from the video sequences [51]. Detailed behavioral studies are possible because the entire stimulus sequence is known and it could be mapped to the response of the animal after the experiment.

Video playbacks methods also have some disadvantages. It is often assumed that the animal is reacting to the stimulus and it perceives the stimulus as being real. This assumption may not always hold. Stimuli are customized yet they are mostly pre-programmed sequences [46]. The animals in the videos do not interact with the real animal and the lack of interactivity may lead to habituation or unnatural reaction from the test subject. The use of display and video technology add other limitations to the playback experiments. The sampling rate of the camera used to create the stimulus must match the physiological visual properties of the animal, otherwise the motion of the animal in the video stimuli may appear blurry, discolored or distorted to the animal. Technical specifications of the display screen or the projector must be considered to avoid similar problems e.g. resolution of the display, refresh rate of screens, etc. There are some common limitations shared by all screen-based methods of visual stimulation, which are covered in detail later in section 5.5. Video playback methods are limited to those behaviors which are possible to record. Further details on feasibility of using video stimuli can be found in the review paper of D'eath [46]. In summary, video playback methods made a strong case in favor of using technology provided the researcher considers carefully its limitations and makes reasonable assumptions.

## 5.3.4 Virtual stimuli

Virtual stimuli are the most advanced category of artificial visual stimuli. The setup is more or less similar to video playback methods but the content of the stimulus is created virtually i.e. using computer graphics and animation technology. The initial motivation for using virtual methods was to increase immersion and remove limitations of earlier methods. In a computer-generated stimuli, the user can configure fine details, which is not typically possible with raw video stimuli (see figure 5.1D). Virtual stimuli are programmable and therefore various properties can modified independently. For example the shape, the size, the color,the background, and the even movements of the animals can be individually changed for each experiment. Graphic design and rendering are inspired from techniques developed by the video game and animation industry.

Virtual stimuli are considered a major improvement over other stimulus types. These stimuli are faster to design and modify, which means multiple experiments can be conducted with different variants. Initially, video playback and abstract stimuli-based experiments were also performed with virtual stimuli for cross-validation [216]. Virtual stimuli are proposed as an alternative to other screen-based methods and this method has been widely adopted for studying different behavior patterns (fish) e.g. mate preference [14, 116], predator response [80], and visual communication [89]. It is not always necessary to design a complex looking stimuli. For example, Ioannou *et al.* [94] projected small dots onto a surface of a fish tank to simulate the movement of very small prey. The fishes attacked the projections as if they were real prey, which helped the authors to understand the hunting strategy of the fish, as well as to conduct artificial evolution of the prey. The methods are constantly updated to create realistic looking animals and scenes. anyFish [211] and FishSim [145] are software packages to simulate 2D projection of a 3D animated fish. Joysticks are employed to define the motion of the virtual animal in semi-interactive manner, or to introduce perturbations in the stimuli [123, 145].

Abstract stimuli-based experiments have benefited significantly from digitization techniques. Software packages have been designed to automate the workflow i.e. stimulus delivery, behavior measurement (locomotion) and data analysis. Fry *et al.* [76] designed fully automated setup for conducting open-loop experiments with abstract patterns. They used an optical tracking approach for computing 3D trajectory of a freely flying fruit fly in real-time. Most importantly, the combination of virtual stimuli and real-time tracking methods provides an opportunity to design closed-loop experiments. We cover closed-loop methods with virtual stimuli separately in the next section.

Virtual stimuli-based methods have many advantages when compared with previously discussed categories. Modular software design allows different modules to be changed as and when the new versions are developed. This provides better control over display, rendering or measurement (tracking). Software based methods are easier to distribute and share with other scientists in the community, which promoted standardization. Most importantly, the stimulus can be programmatically controlled which was not the case with any other method. The method suffers from display technologies based limitations i.e. the displays are made for the human visual system and may lack sufficient spatial resolution, spectral resolution. Technological problems are inherent to the virtual stimuli-based methods are discussed in

**Tab. 5.2.**  Overview of artificial visual stimuli used in closed loop experiments

| Type | Design | Feedback method | Animal-Behavior | Key Attributes |
|---|---|---|---|---|
| Mechanical | Arena, Pattern cylinder | Torque meter, Treadmill | Fruit Fly - Pattern recognition [54], Motion perception [33] | Motion based rotation of pattern cylinder, features not configurable. |
| Hybrid | Arena, LED Screen, Projection | Optical Sensor, Photodiode, Optical Tracking, Treadmill | Fruit fly - Depth Perception[185], Moth - Neurophysiology [83], Rodent - Navigation [130], Neural activities [57, 59, 60, 86], Review - Neuroscience [31, 58], Primate Cognition [56], VR for animals [194], Rodent - [205] | Animals are restricted or tethered, fully interactive, built with open source software frameworks, setup configurable for multiple species. |
| Digital | Arena, Projection | Optical Tracking, Active treadmill | Fruit fly - Real-time 3D tracking [76], vision induced motion [194], Flight pattern [195] Spider - Navigation [165, 194], Ant - Foraging [47], Fish - Social behavior [195], Rodent- [50, 195], Review - Neuroscience [31, 58], VR for animals [194], Rodent - [205] | Free moving animals, real time perspective correction, underwater projection, arbitrary surfaced arena, support for multiple species, configurable software. |

detail with the limitations of the closed-loop experiments. Although programmable, the open-loop experiments with virtual stimuli lack interactivity. Thus, repeated trials with the same animals are not advised as they may get used to the stimuli [94]. Butkowski *et al.* [34] and Woo and Rieucau [216] have reviewed the use of animation for open-loop behavior experiments in more detail.

# 5.4 Closed loop Experiments and Virtual Environments

In this section, we will focus on behavioral experiments that use artificial visual stimuli in a closed-loop. When navigating in a three-dimensional environment, the features visible to the eye change according to the perspective and movement of the viewer. The fundamental idea of a closed-loop experiment is to constantly update the visual stimuli according to the movement of the animal. This design has two major components, *tracking* and *stimulus delivery*. AR and VR technologies used with the humans also come in the same category. It is necessary to synchronize these two components in real-time for a realistic appearance. Real-time tracking and perspective correction for a freely moving animal is a difficult problem. Over the past two decades, different techniques have been developed to circumvent this problem mostly by restricting the movement of the animal. Behavioral experiments with virtual stimuli are often referred to as Virtual Environments (VE) or Virtual Reality (VR) interchangeably in the behavior literature. For the sake of clarity, we will use the term VE generally for closed-loop experiments with virtual stimuli. We reserve the term VR (as a subset of VE) for closed-loop experiments where stimuli are rendered in a perspective correct manner i.e. visuals change with the motion of the head. Based on the designs of the closed-loop experimental setups we have divided them into three different categories: mechanical design, hybrid design and digital design.

Most setups are largely inspired from the CAVE VR design [45] and their designs have remained more or less the same from last two decades. An enclosure is designed for the
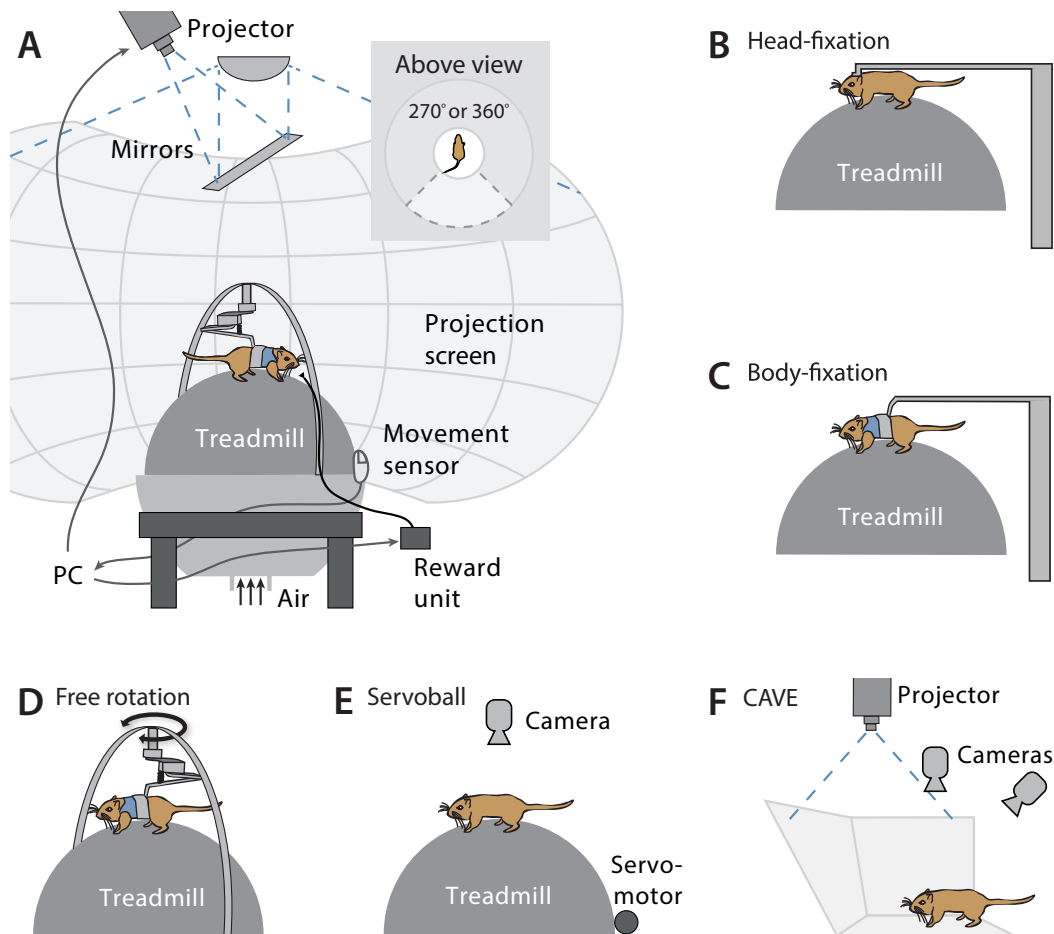
animal and it is placed on a stage or a platform. The platform is surrounded by a toroidal or cylindrical screen that is preferably matching the field of view of the animal, for example see Fig. 5.2A[205]. The experiments mentioned in this section are selected based on the novelty of the approach. Our aim is to highlight critical improvements in the behavior experiments using virtual environments. We show that many of these improvements are largely dependent on the methods developed for computer vision and XR applications.

## 5.4.1  Mechanical design: Restricted animals in fixed (non-virtual) environment

One of the earliest designs of a closed-loop experiment was a mechanically designed flight simulator for insects. The setup is similar to the rotating pattern cylinder design explained earlier in section 5.3.2. In the flight simulator the rotation of the pattern cylinder is coupled to the motion of the insect via a torque motor. This way the motion of the insect triggers the rotation of the cylinder in its visual field, which emulates a real-life flight conditions. Dill *et al.*[54] used a this setup to study visual pattern recognition in fruit flies and showed that flies could remember patterns based on how they appear on the retina from a specific perspective. The fly was tethered and its head rotation was immobilized to restrict the movement of head independent of thorax. The turning response was recorded by measuring the movement of thorax. Often head fixation is used to force the insect to turn its body instead of the head and this modification simplifies the tracking problem. Tethering is used to control the sensory experience and allows experimental recordings to be made during the experiment.

A treadmill with a styrofoam ball is another type of mechanically designed setup for closed-loop experiments. In this case, a tethered animal is placed on the ball and walking motion of the animal rotates the ball which is used as a feedback to turn the pattern cylinder. The rotation of ball is converted to electronic signals which serve as input to the servo motor responsible for rotating the cylinder. The balls rotation is measures with infrared tracking. The ball is painted with a pattern of infrared reflective dots. An infrared LED is placed near the ball and its reflection is picked up by a photodiode which further decomposes the motion into rotation and translation using sequence detector logic. The mapping between the animal's movements and the visual pattern is stored in the computer for further analysis. Bülthhoff [33] used this setup to study the genetic link between vision and motion perception in fruit flies. He used genetically modified flies (mutants) and wild fruit flies in the flight simulator to perform a navigation task. He showed that mutants showed defect in visual orientation and therefore concluded that optomotor response may be encoded in genetic experession of the animal.

Mechanically designed closed-loop experiments were mainly used to study visually induced motion with abstract stimuli. These experiments did convincing demonstrations about the preference an animal for some patterns and showed that they actively move the ball in order to get towards their preferred patterns. However, the patterns remained unchanged or fixed during the experiments which was considered a major limitation of this approach. This limitation is alleviated in the modern closed-loop designs which use projectors and screens instead of pattern cylinders. In the modern VE setups, the concepts of using tethering and treadmills are also adopted from the flight simulator experiments.

**Fig. 5.2.** Examples of different type of VR systems (credit : cf. Thurley and Ayaz [205]). Figure A-F show different techniques used for fixation of animals, recording of movement and display of stimulus. Details covered in text.

## 5.4.2 Hybrid design: Restricted animals in VE

Closed-loop experiments with the hybrid design were motivated by the success of virtual stimuli in the open-loop experiments. In hybrid designs, virtual stimuli are displayed using screens or projectors (instead of rotating pattern cylinder), and treadmills and/or tethering techniques are used to restrict the animal's movements and to simplify the problem of tracking (cf. figure 5.2). In the late 90s, the behavior researchers started the development of closed-loop experiments with virtual stimuli. It was easier to configure virtual stimuli to show desired patterns and the appropriate display technologies started becoming commercially available at the time. Treadmills based techniques were readily available and useful for precise perspective correction while rendering the stimuli on the screen. These experiments are often referred to as the first experiments which put animals in Virtual Environments. VE allowed researchers to try different types of stimuli which extended the scope of research to other behaviors in the three-dimensional world i.e. navigation, foraging, etc.

Schuster *et al.* [185] designed one of the first experiments with fruit flies in VE. The fly was placed on a stage surrounded by a 360° panoramic LED screen which displayed the stimulus

pattern. The 2D movement of a walking fly was measured in the X-Y plane using a simple blob detection technique from computer vision. The fly was tethered and its wings were clipped to restrict its movement to a plane. The authors claimed that they were able to study depth perception in fruit flies with the system which was not possible in previously designed open-loop methods.

Another notable approach is from Gray *et al.* [83], where they designed a VE to measure neurophysiological activity in moths for studying foraging behavior. Moths navigate in a complex 3D environment to find the source of odor. The authors simulated similar conditions in VE by designing a multisensory stimulation (visual, olfactory and mechanosensory) mechanism. A wind tunnel was placed in front of the moth for olfactory and mechanosensory stimulus. A 3D scene consisting of a textured floor and vertical pillars was rendered using techniques from a computer game Descent III. The moth was tethered and multichannel neural recording was obtained by probing the ventral nerve cord of the moth. It was assumed that flight is at a constant altitude (fixed Z) and navigation was restricted to a horizontal plane. The movements of the moth's abdomen were measured using optical sensors. An Infrared (IR) light source with photodiode array was used as the optical sensors. It was shown that the moth navigated in the virtual space by turning towards the odor emanating from the wind tunnel. The authors demonstrated that the turning sequences in VE were consistent with findings of optomotor response observed with freely flying moths. This way a previously trusted technique was used to validate effectiveness of the virtual flight simulator. Generally, fixation of the animal is considered a limitation of this approach. However, the researchers studying the neural link between stimuli and behavior have preferred fixation of the animals to be able to measure the neural activity in a reliable manner [194, 205].

Numerous variations of treadmill-based designs have been used to study the movement of rodents in virtual environments (cf. figure 5.2). Each technique has imposed different degrees of constraint on the movement of the animal e.g. body fixation, head fixation, etc. Hölscher *et al.* [91] designed one of the first VE setup for mammals (rodents), similar to figure 5.2A. In this experiment, the movement of the treadmill was restricted to the horizontal axis, the body of the rodent was fixed but the head position was not. Rotation of the treadmill was computed using optical sensors, similar to tracking the ball in a computer mouse. A DMD projector was used to display stimuli on a screen. Two mirrors were added to reflect the projections for covering a wide field of view (360° azimuth and -20°to +60° elevation). The VE contained cylinders hanging downwards from the ceiling and no features on the floor. This design was meant to avoid giving any tactile feedback that the rat may expect from visual cues. An open source rendering engine (OpenGLPerformer) was used to render graphics in real-time with the support of NVIDIA graphic cards. The stimulus was presented at a fixed distance and stereoscopic depth cues were not considered. The experiment demonstrated that rodents could be trained to navigate in a 3D virtual environment using a 2D stimuli. Such findings had been shown in primates and humans. The rats were trained in a real maze for navigation task to compare their learning ability in the real and the virtual world. They learned to operate the treadmill to navigate "closer" to the objective in VE to earn a food reward. They got better with the number of attempts and consistently minimized the distance to, and thus the time taken to reach, the reward. This method was a significant improvement in comparison to classical lab-based experiments. It removed the limitation of constructing actual mazes in order to study navigation. However, restriction of movement and lack of other
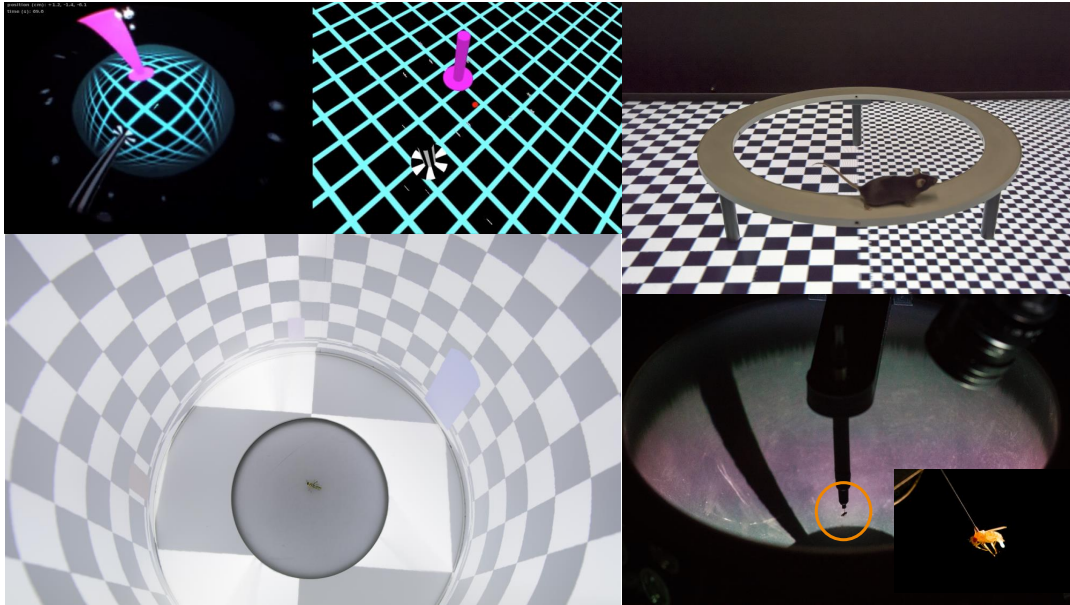
stimuli such as vestibular, tactile or olfactory cues are considered a major limitation of this approach. The ability of rats to learn and adapt to a new environment while suppressing lack of information from other sensory inputs, was nevertheless considered positively for the use of virtual environments. As mentioned previously, head fixation and body-fixation techniques (cf. figure 5.2B,C) were used for head stabilization during measurement of neural activity. For example, recording membrane potentials [86], two photon microscopy [57], two photon calcium imaging [59], patch-clamp recording [60].

VE designs are modified to match the visual properties or motion properties of the animals. These modifications are necessary to investigate specific questions and improve the realistic appearance of the stimulus. We have picked some examples to clarify this point. Free motion treadmills (cf. figure 5.2D), for example, were designed to introduce vestibular information about rotation, which was missing in the earlier designs [130]. Takalo *et al.*[200] increased the field of view and the temporal resolution to customize stimuli for fast-moving American cockroaches (tethered). Dahman *et al.*[47] used hollowed styrofoam design for accurate registration walking speed of desert ants (tethered). They showed that ants changed their pace significantly between different approach and search phases, they slowed down significantly while approaching nest position. Stowers *et al.*[194] designed visual stimuli with configurable chromatic properties to increase the naturalistic appearance of the scene for jumping spiders. Stowers *et al.*claim that such systems are well suited to study visual features important for decision-making behavior such as target selection or predator avoidance.

## 5.4.3  Digital design: VR for freely moving animals

In this subsection, we cover experiments with *true* Virtual Reality designs, which also allow free movement of the animal. The stimulus is rendered in such a way that it creates an illusion of a three-dimensional space from the animal's perspective, even though the projections are on a 2D surface at a fixed distance. This compensation of view is known as perspective correction in human XR literature. Perspective correction is achieved using sophisticated computer vision techniques of employing real-time tracking (with multiple cameras) of the animal's head in 3D. The tracking data is fed to the rendering engine with a minimum delay to provide real-time projection, correctly rendered from the perspective of the animal as it moves through the virtual space. The stimulus is displayed on flat 2D surfaces or arbitrary shaped curved surfaces (see figure 5.3) using high framerate projectors or screens. The VR systems use advanced concepts such as multi camera-projector synchronization and calibration, real-time 3D tracking and rendering [76, 194]. These are from computer vision and XR research (see sec.2.4).

One of the method for designing VR systems with freely moving animals is to use an active treadmill. Active treadmills are used to compensate motion of the animal and keep them in stationary position (cf. figure 5.2E,5.3). Real-time video tracking is used to keep the animal in the same physical location The tracking feedback is fed to a servo motor, which generates counter-motion of the treadmill. This type of design facilitates behavior studies without creating very large arenas e.g. navigation behavior in jumping spider [168, 194] or foraging behavior in desert ants [47]. Treadmill-based solutions are not suitable for all animals and therefore the development of novel 3D tracking methods was crucial for the

**Fig. 5.3.** (clock-wise) Example of stimuli from the fishVR system [195] rendered from perspective of fish (left) and human (right), MouseVR system with free moving rat on circular platform [195] (credit: https://strawlab.org/freemovr), FlyVR setup with tethered fly (credit: Simon Gingins), top view of VR arena made for terrestrial insects (credit : Centre for the Advanced Study of Collective Behaviour, Konstanz )

development of VR solutions for animals. Fry *et al.* [75, 76] created TrackFly framework to conduct high throughput closed-loop experiments with unrestrained flying fruit flies. They tracked free moving flies at 50 Hz using a multi-camera setup. Building upon this tracking approach Stowers *et al.* [194] built a FlyVR system. Markerless tracking was done with infrared filters to facilitate fast image processing and block the visible light from stimulus screens. Stowers *et al.* [194] showed that it was possible to combine real-time 3D tracking and stimulus delivery to induce flight movements in the desired 3D trajectory. They introduced the concept of a modular and reconfigurable framework designed for animal VR systems. The FlyVR framework supported the configuration of multiple camera-projector systems along with accurate geometric and photometric calibration for arbitrary surfaces. This is an advantage over previous methods as the same framework can be used with different configurations (tethered, free-flying and treadmill) for different animals. Open-source frameworks such as ROS and OpenSceneGraph were used, which are well known in the robotics and graphics community. Additionally, they also showed a new approach where multiple flies could be tracked while the stimulus was delivered by focusing on one of the flies.

Stowers *et al.* [195] extended their research and designed VR setups for rodents (MouseVR) and fish (FishVR) with FreemoVR platform. FreemoVR platform is designed to provide flexibility to the user. The platform supports a wide range of stimuli, naturalistic and abstract, and adapts to different display configurations based on the species. The FishVR system is the first underwater VR application, where visual stimuli are projected on a fish bowl from below, and perspective correction (see figure 5.3) is supported by tracking position of the head using infrared 3D tracking. The study confirmed that fish responded to the virtual stimulus as if they were real. It was shown that fish avoid virtual obstacles placed in the fish tank and swimming around it as though obstacles were really present. In addition, when a virtual

**Fig. 5.4.** The image displays three different animals placed in Virtual Environments. (left) A tethered fruit fly placed in a VR experiment to study decision making (©MPI-Animal Behavior). (center) A freely moving locust placed on an active treadmill to study decision making ©Center of Advanced Study of Collective Behavior, Konstanz. (right) Praying mantis with glasses used in an open-loop experiment by Nitynanda *et al.*[161] to study depth perception (©Newcastle University, UK).

conspecific (same-species fish) was introduced, they swam with them as though in the real world. VR for freely behaving animals offers new avenues for research in the field of social and collective behavior. The MouseVR setup is designed to allow mice to move freely on a raised circular platform (see figure 5.3). In this experiment, stimuli were displayed on the floor, a similar setup is also used by Del Grosso *et al.*[50]. Experiments with checkered patterns show that freely moving mice estimate height using motion parallax, a finding that was not possible to test in earlier mention treadmill based systems [195]. Experiments with freely flying flies in FlyCave indicated that flight control of flies is fundamentally altered by tethering, even without head fixation [195]. The authors used this study to stress the importance of developing new methods for free-moving animals.

## 5.5  Limitations of Virtual Environments

In this section, we will discuss the limitations faced by researchers while designing virtual environments for studying animal behavior. Virtual environments are designed to create believable experiences for animals by artificial stimulation of their sensory apparatus. The main challenge is to create realistic simulations which change continuously based on the behavioral response of animals. Currently, this is done primarily by displaying visual stimuli to the animal by using screens or projectors, and tracking their response by using a camera and treadmills. The stimulus can also be controlled externally to introduce perturbations. This approach is limited to some animals because existing technology is not capable of solving tracking and simulation related challenges generally, for all species. Most of these limitations can be attributed to the physiological properties of the animals. We examine the limitations of existing tracking and stimulus delivery methods and link them to the physiological properties of the animals. Our discussions are inspired from other reviews which focus specifically on limitations of using screen/display based artificial stimuli i.e. video playback [46], animation [216] and VE [41, 194].

### 5.5.1  Limitations of stimulus design and delivery

Animal vision has evolved for enhancing survival rate, therefore, different animals have different visual properties such as color vision, the field of view, etc. The stimulus employed must

therefore be compatible with the requirements of each animal's visual system. The stimulus design and delivery approach are also crucial for the success of behavioral experiments. In these respects, all commercially available technology has limitations. Certain technical considerations are made when such technologies are used for behavior experiments with animals. We have already discussed vision properties in the introduction (see sec.2.1.1) with focus on human vision properties. In the following text, we outline some important visual properties and relevant technological considerations.

**Multispectral vision** is the ability to visualize different spectra of light. It is also known as spectral resolution in the literature. Humans and some primates are trichromatic, which means that a combination of three colors (Red, Blue, and Green) is sufficient to cover the entire color spectrum seen by humans. In the case of animals, some are dichromatic (most mammals) or tetrachromatic (birds, reptiles), and some animals see completely different hues e.g. UV, UV with red, UV with green [46]. Invertebrates commonly use polarized light for navigation and therefore show a preference towards it. Failure to reproduce such properties may affect the experiment.

Technical considerations: The LCDs and projectors designed for human vision are useful for some animals with trichromatic vision and dichromatic vision. In some cases, lighting conditions are changed [194] or color filters are used to match the colors on the display with the color perception of animals [160](see figure 5.4). Creating realistic colors for animals with multi-spectral vision (e.g. UV) is difficult and should be considered when designing an experiment. Invertebrates have attractions towards some wavelengths of light, which should be considered in order to avoid unintended manipulation of the behavior.

**Flicker fusion threshold** is the threshold beyond which a flickering pattern appears continuous to the observer. Visual information from the environment is integrated a certain time before it is experienced. This integration time is varies in different species. Fast-moving animals typically have a higher flicker fusion threshold. A movie displayed at a frequency of 25 Hz is sufficient for humans to perceive continuous motion but the same movie would appear flickering for the bees. The illusion of motion may be broken by the slow or glitchy movements of the stimuli and should be avoided.

Technical considerations: Flicker fusion is an important criterion when selecting the display and lighting for illumination of experimental arenas. The light source may appear to flicker if animal's flicker fusion threshold is higher than operational frequency of the light source. Inger *et al.*[93] have discussed potential biological affects caused flickering of artificial lights. The same criterion is applicable for choosing operation frequency of the display or projector, and also while choosing rendering rate of the stimuli. Most existing methods use displays at 120 Hz and they use GPUs to accelerate rendering process. It is also possible to manipulate different factors in the setup to decrease the flicker fusion threshold of the some animals. For example, manipulating size of stimulus, luminance of display, brightness of surrounding and region of retina involved in image formation [46].

**Visual acuity** of an animal is it's ability to resolve spatial detail. It can also be defined as the spatial resolution of the eyes. It is measured in degrees; some animals have very high acuity (e.g. eagles, falcons) and some have very low acuity (e.g. fish or ants) [46]. Animals with

very high acuity may see the displays as pixelated surfaces or surfaces with holes, which may affect the experiment.

Technical considerations: It is considered for selection of the display screen or projectors. The distance between the screen and the eyes is also a an important factor. The distance from the screen must be greater to avoid problems with acuity. If this requirement is unmet, the illusion of continuous color might be broken, which may be an important consideration for the experiment [46].

**Field of view (FOV)** refers to the area/volume observed by the eyes at any given moment. It is usually measured in degrees and can vary widely between different animals. FOV depends on the position of the eyes and the construction of the eye. Animals with front-facing eyes and overlapping vision (e.g. primates, cats) have considerably smaller FOV than animals with eyes on the sides of the head (e.g. birds or insects). FOV may change slightly for animals that can rotate their eyes in the socket. It is also notable that most animals do not have sharp vision in all parts of the FOV i.e. high resolution at fovea and less at the periphery.

Technical considerations: FOV is considered while deciding the display area and shape of the screen for the stimulus. Most of the time curved or cylindrical display screens are selected for small insects and rodents [205]. Projectors are preferred over LCD screens because curved LCDs are difficult and expensive to produce. A larger FOV is ctypically overed by using multiple projectors, which adds the additional complexity of synchronization and calibration of projectors.

**Depth Perception** Most animals have some mechanism to perceive depth in the environment. Different animals use different cues such as stereopsis, motion parallax or focusing, overlap, shadow, vertical distance to the horizon, retina to image size ratio, perspective and texture [46]. Biologically stereo vision is not always necessary for all species, many species use non-stereoscopic depth cues because they have limited overlap between field of view. Failure to accommodate some depth cues may reveal the 2D nature of the stimulus [43, 216].

Technical considerations: Depth cues are considered while designing appearance of the virtual stimuli. Most artificial stimuli based methods display stimulus on flat or curved surfaces. It is likely that animals can perceive the flat or curved surface of the screen if the stimulus is rendered without correct perspective correction, which may affect their behavior. Recent VR methods use 3D head tracking or body tracking to maintain perspective of the animal but do not offer stereoscopic depth cues. Researchers must be careful while designing experiments which require the animal to may be use depth cues. Stowers *et al.* [195] suggest that tracking eye movements is important for introducing depth cues. Recently, Nityananda *et al.* [161] glued color filters to study stereoscopic depth perception in insects (see figure 5.4). They show that it is possible to use such modification when the research question is chosen appropriately.

## 5.5.2 Limitations of tracking

Tracking the movements and the perspective of the animal is essential for designing a VR experiment. Tracking freely-moving animals is challenging and most of the experiments still require tethering or other restrictions. Adding markers on animals may affect their natural behavior, but recent advances in computer vision have shown promising results for markerless tracking [169]. Existing tracking limitations often stem from physiological properties which are explained below.

**Locomotion properties** Animals possess diverse abilities to move in their environment, such as flying, swimming or jumping. Often the speed of locomotion may vary and some movements (e.g. jumping) have to be restricted in order to keep the animal in the designated space. Accurate movement tracking is necessary for mapping each decision (in terms of movements) of the animal to the visual features rendered in the virtual world. Mismatch in mapping can potentially invalidate behavioral findings.

Technical considerations: Locomotion properties of the animal play a big role in selection of the tracking approach. Ideally, the animal should be freely moving but restriction may be necessary depending on the need of experiment (e.g. neurophysiology). Tethering or treadmill based approaches (figure 5.3) may be selected if feedback from other sensory systems can be compromised or disregarded for the purpose of the study. In both cases, optical tracking is used for tracking movement of the animal. In the case of treadmills, the motion of the treadmill is measured to infer the movement of the animal. The cameras selected for the purpose must operate at higher frame rates than the rate of rotation of the treadmill. Additionally, the mechanism of the treadmill must sensitive towards variations in the movement of the animal. For example, the styrofoam ball must accelerate and decelerate in sync with the animal's motion otherwise it may create an unwanted perturbation for the animal [47]. Similar considerations must be made when selecting cameras for tracking motion of the animal in tethered cases. For example, Stowers *et al.*[195] used sampling frequency of 100 Hz to compute motion of the fruit flies.

Computer vision algorithms are used to track motion of freely moving animals. The performance of such algorithms is dependent on the visibility of the animals in the images. Camera properties such as frame rate, resolution, opening angle, rolling/global shutter, must be considered to capture the movements of the animal. 3D tracking requires multiple cameras which adds technical complexity regarding calibration and synchronization of cameras. Active treadmills are designed to restrict animals to one particular spot to reduce tracking complexity. Often lighting configurations are selected to enable real-time tracking. For example, infrared (IR) is preferred because it is easier to add much more light to the scene without disturbing behavior of the animal [194].

**Appearance** of many animals differs in terms shape and structure which presents novel challenges for purely image-based tracking of animals. Some animals do not have any conspicuous features, and some have repetitive patterns which makes different parts of the body appear identical. Such confusions lead to fluctuations in the tracking results, and consequently in the presented stimuli.

Technical considerations: The appearance of the animal is an important consideration for the selection of tracking software, camera and light conditions. The software must process the image in real-time and detect the animal. Paint or reflective markers can be used to add features to have seamless tracking results. Many recent improvements in marker-less tracking methods can allow real-time tracking of seeming featureless objects. Light conditions are often changed to increase the detection rate of markers or animals, while high-resolution cameras are useful for capturing fine details of the animals. However, high-resolution images require a longer time for processing and storage, and therefore the selection of camera often involves a tradeoff.

### 5.5.3 Latency

Latency of a closed loop system is the overall delay between movement of the animal and the change of stimulus on the display. To allow for an interactive experiment in real-time the latency must be very low. Multiple computational steps are involved between these two events such as image processing, data storage, graphics rendering, etc. Each of these steps introduces a time delay in the system. Overall latency of the system is caused by both software and hardware components. Designing real-time tracking solutions is very difficult. Hardware requirement such as powerful computers, higher bandwidth data transmission, and responsive displays are difficult to meet both technically and economically. The overall latency of the system must be considered while selecting an animal for a study.

### 5.5.4 Lack of technical expertise

VEs for animals are designed by biologists using the technology developed by engineers. Modern VEs are rely heavily on software and algorithmic methods from computer vision and XR communities, because their work available through open-source distribution. These methods are complex and certain expertise are required to tweak these methods to be able to use them with animals. Biologists are forced to develop engineering and programming skills to develop new concepts for behavioral experiments in virtual environments. Only a few biologists have successfully bridged the gap between these fields. We consider that lack of technical expertise, from CV & XR community, is one of the biggest limitations for the development of VR for animal behavior experiments.

This point is also highlighted in the next chapter, in which we have presented a survey of researchers studying behavior. A relevant finding was that most researchers want to collaborate with technology experts from computer science (92% survey participants). About 85% of participants think that closed-loop experiments are important for studying behavior. Details discussion on results are in the next chapter and overall results of the survey can be found in the appendix A.

**Fig. 5.5.** (top) Mixed-Reality with multiple animals : A school of fish in a large fish tank where bottom of the tank is illuminated with multiple projectors to provide desired visual stimulus to a group, (bottom) Tracking multiple animals : Locusts with markers and 2D posture tracking of locusts using DeepPoseKit [82]

## 5.6 Future Directions

In this section, we will discuss the future of interactive behavior experiments based on the current research trends in robotics and computer vision. While it is true that technology experts in XR are not closely involved in behavior studies yet. The research in computer vision and robotics is also relevant to closed-loop behavior studies. The aim is to show that these two communities have started showing interest in behavior studies. We want to capitalize on this fact and motivate experts from XR to join in the effort. In this direction, we will suggest a few ideas for the development of novel XR applications for behavior studies that may be realized with support from experts in the XR community.

We believe that research in closed-loop behavior studies will thrive from the contributions of an interdisciplinary community of engineers, computer scientists, and biologists. We are motivated to attract attention in this direction because of the infrastructure that we are creating for conducting research in collective behavior (presented in chapter 6). Our research is merely a starting point that may encourage further discussion on this topic.

### 5.6.1 Involvement of robotics community

Animal behavior studies are moving towards an increasingly interdisciplinary approach. The behaviors and mechanisms studied in the animal VR are gaining attention in the field of robotics. Vision induced locomotion and navigation studies in small insects are appealing for designers of nature-inspired robots [13, 112, 223] and self-navigating drones [97]. Studies focusing on the understanding of sensory mechanisms of small animals are gaining interest in the field of smart sensor design. In 2018, DARPA launched a robotics challenge to design small,

lightweight and power-efficient micro-robots for use in disaster relief scenarios of the future. The recent developments in free moving VR systems have opened doors for conducting new types of studies in collective behavior and social behavior, with developers of self-organizing robots already using the theories developed in collective behavior studies [203, 215]. The robotics community is actively involved in development of new methods for studying behavior using interactive robots [108, 112, 120, 213].

We have designed a new setup for studying collective behavior using a motion capture system (see chapter 6). In the next chapter, we show that our setup is capable of conducting closed-loop experiments with animals. We have remained focused on XR technologies with visual stimulation, however, the use of robotic agents should also be possible in our setup. In the near future, we would like to collaborate with technology experts in the robotics community for designing novel experiments in a larger space with a wider range of animals.

## 5.6.2  Involvement of computer vision community

Recent publications in computer vision literature show that the community is taking interest in challenging problems involving animal tracking. Many researchers have proposed easy to use methods for 2D tracking of keypoints [82, 135] (see Fig. 5.5). 2D posture computation in animals is considered vital for video-based activity recognition of complex behavior patterns e.g. courtship display or aggression display. Marker-based approaches are also used for tracking the position of multiple individuals at the same time [82] as shown in Fig. 5.5. For VREs, real-time methods for tracking of the eye positions and head orientations of animals is missing. In the case of small animals the head position is inferred from the animal's position and orientation, however, it is difficult to do the same animals with articulated bodies. Extracting 3D posture of animals from images and videos is an emerging topic in the computer vision community [84, 136, 162, 225]. In our opinion, more research in this direction is required but equal involvement of biologists is necessary. We will present our research on 3D posture problems in the next chapter.

## 5.6.3  Introducing new concepts of XR

**Spatial Augmented Reality (SAR)** applications are not fully explored in animal behavior experiments. Projectors are readily used in behavioral experiments but often their use is limited to displaying stimuli in open-loop e.g. predator-pray interaction [94]. Ioannou *et al.*[94] used open-loop approach due to the lack of methods to track fish in real-time. Now, it is possible to perform similar experiments in closed loop with the help of real-time tracking methods. One possible application is the use of dynamic projection mapping with robotic animals. It was shown that social behavior can be studied using robots instead of real animals, and animals can and do interact with robots [112, 120]. Landgraf *et al.*[120] tested robotic models with different appearances and claimed that appearance was crucial for the acceptance of a robotic agent by real fish to study the social behavior of fish. We argue that dynamic projection mapping techniques [143, 152] can be deployed to alter features of robotic stimuli. The projector can project different patterns on a robotic agent while maintaining the perspective of the real animal using real-time 3D tracking. Experiments with SAR could open

new possibilities such as training animals for navigation or memory experiments using virtual agents projected on a wall or a robot.

We will revisit this topic again in sec 6.3, where we present a concrete example of using Spatial AR with displays to study collective vigilance.

**Collective behavior studies** related to the decision making of a group and the effect of the individual decision on the group may be studied using VR. It is shown that real animals do interact with virtual conspecifics in the VR e.g. fish[195]. Multiple VR systems can be plugged together to create conditions for social behavior in a virtual manner i.e. collaborative VR space for animals [121]. In such a virtual social scenario, each animal may interact with a group of virtual conspecifics which are projections of real animals from other VR systems. The visual information available to each individual can be controlled in such environments and manipulated based on the needs of the experiment. If the animals start to swarm in virtual environments the principal's governing their decisions can be studied in much detail.

Another idea to study collective behavior is to use multiple animals in the same space. This requires creating a larger CAVE arena or developing an entirely new concept involving Mixed Reality or Augmented Reality solutions. As of now, there are not many ideas presented in this direction. The discussion in the next chapter are entirely dedicated to this particular problem.

## 5.7 Conclusion

In this chapter, we discussed a broad range of behavioral experiments with a focus on closed-loop experiments. We learned that virtual environment based behavior experiments are an exciting application domain of XR technology. These behavior studies support fundamental research from vision and cognition to robotics and artificial intelligence. We also learn that there are many limitations in this field of research. The existing VEs are not yet suitable for studying all types of behaviors. The animal's sensory feedback is generally restricted to prefer one sensory input using methods previously described such as wing clipping, body fixation. Experiments with multiple animals are not possible and most setups cater to a small range of species. We also show that researchers in other fields have started taking an interest in behavior studies. However, most projects are individual efforts and more collaborations are required to improve the state-of-the-art in VEs for animals. We argue that some of these problems can be alleviated in the future by starting new collaborations between experts in computer science and biology. Overall, we hope that the discussions presented in this chapter should provide a solid foundation in favor of the research presented in the next chapter.

# XR Techniques For Collective Behavior Studies

In this chapter, we will introduce the readers to a novel experimental setup, "The Barn" that is specifically designed to conduct collective behavior experiments with interactive sensory stimulation. Although, open-loop applications are also supported we will mainly focus on exploring capabilities of the setup for conducting closed-loop experiments. The original idea of the setup is conceptualized by leading researchers in the field of collective behavior at MPI-AB and CASCB, University of Konstanz. The larger objective is to design a multi-sensory environment that allows a wide range of experiments with multiple individuals of the same or different species. The real-time tracking approach is specifically chosen to develop a new type of closed-loop experiments for interactive sensory stimulation. To best our knowledge this is a first attempt to build a 3D tracking facility for tracking multiple species of wild animals in a relatively large area (15 meters x 7 meters x 4 meters).

In the previous chapter, we outlined some of the major limitations of behavioral research with VEs (sec. 5.5). In this project, we have addressed two of these key limitations i.e. tracking and lack of technical expertise. The proposed setup is designed with a contribution from a multidisciplinary team of researchers from biology, physics, and computer science backgrounds. The long term vision of this project is to promote collaborative research activities between researchers in collective behavior, computer vision, machine learning, and XR. The complete project includes intellectual contribution many researchers and therefore providing details of all possible experimental ideas is beyond the scope of this thesis. The scope of this chapter will be limited to the author's contributions to the project as a computer vision and XR researcher.

The primary task was to choose a tracking strategy that is suitable for conducting closed-loop experiments. We collected the requirements from biologists and selected a suitable tracking strategy. A commercially available motion capture system was installed and preliminary tests were conducted with several species of insects, birds, and mammals. We will introduce the system to the readers and briefly mention observations of the preliminary tests. Perspective tracking is crucial for visual stimulation based closed-loop experiments. 3D posture tracking of animals is an open problem in computer vision and it is extremely relevant for behavior experiments in general. Therefore, we focused on using the motion capture system for developing approaches for posture tracking. We used birds as model species and used 6-DOF object tracking functionality of the motion capture system to show that real-time posture tracking is achievable with our setup. Furthermore, we show that the marker-based tracking approach can be used for preparing a dataset to solve the problem of marker-less posture (2D or 3D) estimation. We further comment on the strategy we are currently using to develop the dataset and the methods for predicting 3D posture using a single camera. Finally, we propose

use cases of our marker-based posture tracking methods for developing XR applications with both humans and animals.

This chapter is written for both computer scientists and biologists. We believe that both groups will appreciate our effort to show that complementary understanding is vital to find the right solution to study behavior. Working with wild animals in a large area is a challenging task because of unpredictable factors that affect the performance of the tracking approaches e.g. animal species, size of the animal, locomotion properties, group size, etc. In behavior studies, a perfect tracking solution is always desired but difficult to realize. The task of a technology expert is to design the best possible solution with minimum loss of data and the task of behavior researcher is to design the experiments with considerations to the limitations of the system. With this intention, we have provided general guidelines for biologists who want to study behavior using the setup and present open problems to invite collaborators from computer scientists and the XR community.
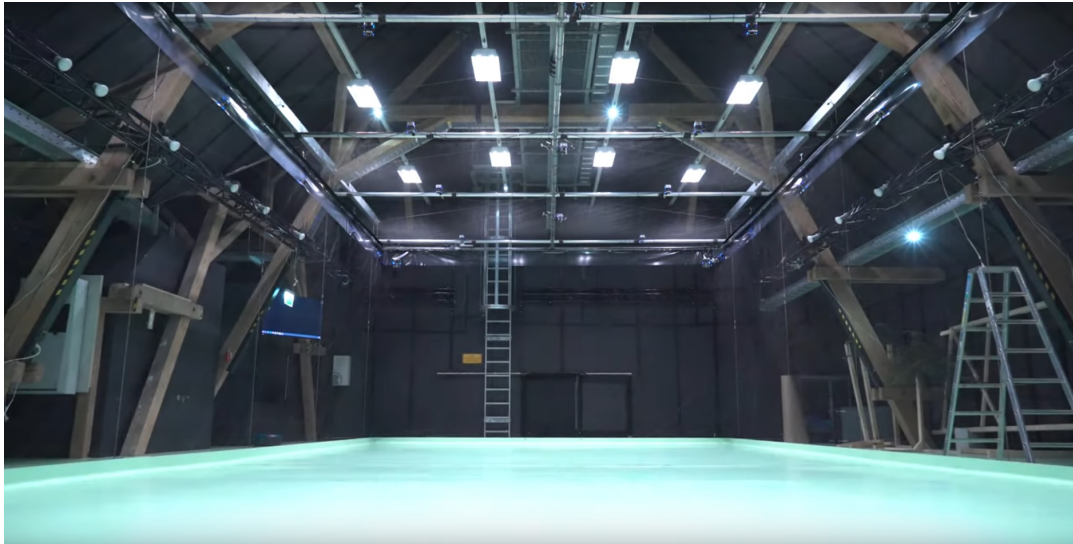
## 6.1  Introduction : The Barn

### 6.1.1  Background

A typical collective behavior experiment involves measurement of behavior patterns (e.g. movements, vocal calls, postures) of multiple individuals [147, 196, 198]. Scientists often customize the setup design based on the available technical resources, study species, and their technical expertise. Ultimately these choices dictate the range of experiments that can be performed. The measured behavioral patterns, study species, and the duration of recording depending on the experiments. Therefore, the technologies chosen to measure the behavior patterns vary from one experiment to another i.e. microphone for acoustic tracking or optical cameras for visual tracking. Often commercially available solutions are tweaked to meet the purpose of the experiment. Developing experiment specific setups are demanding in terms of space, monetary investment, and technical skills. If the ultimate goal is to accommodate a wide range of experimental ideas, it is important to set up a multi-sensory environment in a large space that may allow the measurement of multiple quantities for a group of same or different species. To the best of our knowledge, such environments do exist for studying human behavior patterns but not for animal behavior.

Collective behavior experiments with closed-loop stimulation are a special case. Interactive robotic agents have been used as visual stimuli in collective behavior experiments [120]. However, these setups are limited in terms of the experimental area and the range of species that can be tracked. Experiments with virtual environments are mostly performed with a single individual in relatively smaller arenas [149, 195]. Tracking animals in larger areas require suitable tracking infrastructure for getting real-time feedback. Moreover, it is important to solve the problem of 3D posture (mainly head orientation) tracking to introduce new methods of interactive experiments to other species such as birds.

To overcome limitations of the existing experimental approaches we designed a 3D tracking set up in a cage of dimensions 15 meters x 7 meters x 4 meters. The cage is situated in a

**Fig. 6.1.** The image shows a side view of the facility designed to study collective behavior at lage scale. The dimensions of cage are 15 meter x 7 meter x 4 meter. The height of the roof can be changed from 4 meter to 2 meters. The cage is covered with a special net to prevent animals from escaping the area. The stucture is designed to specially to offer multiple mounting points for tracking technology.

barn at the Max Planck Institute of Animal Behavior, Radolfzell, Germany. One of the first objectives was to choose an outside-in (see sec. 2.4.3) tracking strategy suitable for open-loop and closed-loop collective behavior experiments. We collaborated with behavior researchers and defined some key tracking requirements.

## 6.1.2  Tracking requirements

The facility is designed with the input from several collaborators and experts at MPI-AB and CACBS. We will provide details of some important requirements in this section. Additionally, we also conducted a short survey where we asked several behavior researchers (n = 13) to score different features that they would like to have in a setup that is designed for collective behavior studies. The generic scale of scoring is 1 (low) to 5 (high). The results are mentioned along with explanation of the requirements, complete results can be found in the appendix A.

### Target species

Generally, the setup is designed to conduct experiments with multiple species of animals i.e. insects, birds, and mammals. The selection of technology for such a wide range requires a lot of flexibility and trade-offs. We learned that collective behavior experiments with insects in such large areas may be beneficial but not absolutely necessary. In the case of mammals, bigger mammals (except humans) are difficult to maintain or control in large spaces. Such large spaces for 3D tracking of birds are not designed yet. It was concluded that the setup should be customized for tracking birds with the assumption that any technology capable of tracking fast-moving birds would also be able to track some mammals or insects.

**Survey notes:** we asked the researchers about importance of developing a generic setup for studying multiple species. In response, 12 participants scored 3 or above and 11 participants scored 3 or more for having the feature in their setup.

### 3D Position - Trajectories - Posture

Tracking problems in collective behavior experiments are reduced to two-dimensional spaces [52, 82, 94, 147] because suitable methods are not available. However, most animals move in three dimensions and therefore 3D tracking is identified as a key requirement to extend the scope of collective behavior studies. Movement trajectories play an important role in the analysis of behavior as it allows researchers to answers questions about emergent properties in highly dynamic environments e.g. bird flocks [44, 71]. Therefore, the sampling rate or tracking speed is also considered an important criterion (discussed below).

One of the most challenging requirements is 3D posture tracking for a group of individuals. Perspective tracking as a subset of posture tracking problem is a major limitation in existing closed-loop setups. Most VE built for animals work with smaller animals and do not track *true perspective* of the animal. Perspective is rather inferred from 3D tracking of the head position. We have given special attention to this requirement while building the setup.

**Survey notes:** We asked different questions about requirement of 3D tracking. 9 participants (score > 4) thought that measurements of 3D movements was important for collective behavior studies. However, 11 participants thought that setting up experiments with 3D measurement was difficult and rated difficulty of 4 or above (0% participants scored < 3). When questioned about reasons of not conducting research in 3D, 9 participants agreed that methods were not available, 11 agreed that technical complexity of implementing methods was very high and 6% felt lack of collaboration with technology experts was a reason. Interestingly, 10 participants scored 3 or above in favor of having 3D tracking in their setup. The same number of participants scored 4 or more in favor getting automated tracking and trajectory results from their setup.

In terms of real-time tracking and closed-loop experiments. 11 participants think that it is important to conduct experiments with closed-loop techniques (score 3 or above). Same number of participants scored difficulty level of designing such experiments as 4 or higher. More than half (7) of the participants (score 3 or more) wanted to have 3D posture tracking in their setup. Same number of participants scored 4 or more in favor of having real-time tracking in their setup.

### Individual detection

Individual recognition is shown to be crucial for many collective behavior experiments involving decision making or leadership [52, 196, 198]. This is one of the most important and challenging requirements. The designated space is capable of fitting many individuals and therefore identity must be maintained in an automated manner for high-throughput experiments. Manual annotations will require hundreds of man-hours which is not feasible. It is also essential to solve the problem for closed-loop experiments, especially for cases when a particular focal individual is targeted for displaying of the stimuli. It is a challenging computer

vision problem due to the complex appearance of animals (discussed in the previous chapter). This requirement increases the data capturing and processing requirements dramatically.

**Survey notes:** We asked participants to rate individual tracking as a desired feature in their setup. 10 participants scored 4 or more in favor of having individual tracking within the session. Whereas, 11 participants scores 4 or more for having individual tracking between different trials.

### Tracking speed

As we mentioned earlier, Tracking speed is an important technical requirement and therefore discussed separately. Tracking speed for our applications includes both data capture speed and processing speed. Capture speed defined the sampling rather and affects both (open-loop and closed-loop) experiments. The processing speed is rather crucial for closed-loop experiments.

The locomotion properties of the animal play an important role in the selection of the capture rate for tracking technology. Birds can achieve remarkable speed while flying and therefore capturing movements of flying birds in real-time is a challenging task. It should be noted that all birds do not behave in a similar manner and their movement speed is much lower while they are walking on the floor. For closed-loop systems, the sampling rate and processing rate must complement each other to reduce system lag (see chap. 5). We concluded that we should select sensors that allow us flexibility for the selection of operating speed based on study species and experiment type.

**Survey notes:** We have already discussed survey results regarding need of real-time tracking requirements (7 participants scored 3 or more). It is not only important requirement in terms of closed-loop experiments, it is also important for getting faster results.

### Data storage and transmission speed

Data storage and transmission speed are influential in designing the duration of the experiment. Behavior experiments may require several minutes to days depending on the behavior being studied. Furthermore, most behavior experiments require replicates where different parameters are varied to observe the corresponding variation in the behavior. Therefore limited data storage and transmission speeds can heavily affect the planning and execution of the experiments, especially when a large number of sensors are required. Transmission speed is always an important consideration for real-time performance. We concluded that this requirement must be considered during the selection of sensors for tracking technology. The implemented solution must allow different strategies for data storage depending on the needs of the experiment.

### Accuracy

Accuracy requirements do vary with the type of experiments and the quantity being measured. This requirement may also vary with the number of individuals being tracked in the space and their size. If the tracking error larger than body size could lead to errors in trajectory assignment which is important while studying social animals who have a tendency to interact within close proximity of each other. Ideally, the accuracy of the system should be at least an

order of magnitude lower than the size of the animal. Closed-loop experiments require very high accuracy to display visual stimulation based on perspective. If the stimulation is based on the locations of the animal, accuracy requirements are relatively lower.

**Survey notes:** All participants in our survey scored 4 or above for having a highly accurate tracking accuracy in their setup.
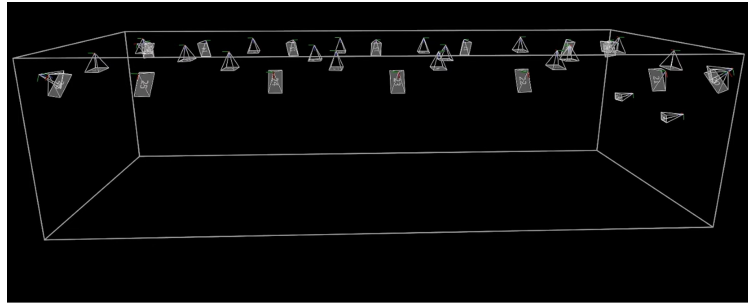
## 6.1.3 Sensor consideration

We considered commercially available solutions because they are easy to deploy and come with sufficient technical support. First, we considered motion capture systems with IR sensors e.g. VICON, OptiTrack. Commercially available solutions are extremely fast and used widely for applications where real-time tracking is vital (covered in sec. 2.3). These systems are capable of tracking a large number of objects with sub-mm accuracy. They also provide complimentary software solutions for customizing the tracking solution for posture tracking applications. One major limitation is that these systems operate with IR reflective markers.

We also considered motion capture systems with video cameras (RGB). This major advantage is that this approach does not require markers and video data captures a lot more information which could be potentially valuable for behavior studies e.g. behavior classification, markerless identification. RGB motion capture technology is still at a nascent stage and limited to relatively smaller areas. Many systems require a green screen and data processing can not be done in real-time. There are very few commercially available RGB motion capture systems suitable for our purpose. The research in markerless human tracking and behavior analysis is driven by marker-based (IR) motion capture systems [95]. We learned that multi-camera video setups designed for computer vision research are specially designed to solve problems in 3D human posture and require large investments for development and data management. Finally, the methods developed for image-based animals detection or posture tracking are not suited for high-throughput behavior experiments in large spaces [103, 104, 225, 226]. As of now, it is only possible to capture the posture of animals in 2D [82, 135], and 3D posture tracking for behavior are limited to smaller environments [72, 84, 105].

After consideration of available technologies, it was clear that investment of time and resources is required to develop a purely video-based customized solution. Therefore, we decided in favor of developing a hybrid tracking approach and installed IR-based motion capture technology coupled with synchronized video cameras. The initial strategy is to use marker tracking as a primary solution to conduct experiments and use the video data as supporting modality. The long term strategy is to leverage the marker tracking solution and develop methods for video-based tracking and behavior analysis. The data collected from marker-based experiments can serve as ground truth for problems in 2D-3D posture tracking, identity detection, behavior classification, etc.

There are several key reasons for going with this strategy. Motion capture systems have extremely fast response, reliable, and accurate tracking performance. The technology is readily deployable which is crucial for starting behavior experiments right away without doing any technology developments. Our setup is an indoor environment with artificial lights which

**Fig. 6.2.** The image shows the optical tracking setup using VICON motion capture system. The outside-in tracking configuration involves 36 cameras (30 IR and 6 RGB).

makes a good argument in favor of using IR based tracking. Synchronized video cameras are useful for recording additional information from the experiment and to visualize tracking data recorded by IR sensors.

It is worth mentioning that we considered using depth sensors as well. They are very successful with the human posture tracking approaches in real-time [30, 122, 189]. We used Microsoft Kinect and Intel Real-sense (D-435) in a test setup (1 sq. meters) with 2-4 birds. We concluded that the accuracy and range of the depth sensors make them unsuited for our setup. Moreover, calibration and synchronization of multiple depth sensors can be a complex process and require a lot of customization.

## 6.1.4  Setup specification

### Cameras

The setup consists of 36 cameras in total - 26 Vero 2.2 (IR), 4 Vantage5V (IR), and 6 Vue 2(RGB). 26 Vero cameras operate at a maximum of 330 Hz (max 330) with full resolution of 2048 x 1088 (2.2 MP). Vantage5V cameras operate at a maximum of 420 Hz (max 2000) with a full resolution of 2342 x 2048 (5 MP). Vue cameras operate at 30 - 60 Hz (max 120) with full resolution of 1920 x 1080 (2.1 MP). The cameras can operate at higher frequencies at cost of resolution and accuracy. The system is fully synchronized in terms of data capturing, therefore operating frequencies of the cameras can be different but only as multiple of each other i.e. IR at 100 Hz and RGB at 50 Hz. The readers are encouraged to check data sheets provided by VICON for specific details of cameras.

### Sensor arrangement

Figure 6.2 shows the arrangement of the sensors. 34 cameras are rigidly attached to the ceiling. Ideally, their locations can be moved according to the needs of the experiment. IR cameras are the main tracking modality and therefore placed in a way that whole cage is covered. Special attention is given to support the fact the ceiling of the cage is variable and can be raised from 2 meters to 4 meters. The Vero cameras are mounted on the ceiling, 12 looking top-down in a grid fashion and 14 on the edges with tilted inwards, such that they cover the entire volume of the cage (see Fig. 6.2). They are expected to track all movements in the lower half volume of the cage (at full height) with very high accuracy. This limitation is due to

the triangulation based tracking strategy. The frustum (filed of view) of the camera becomes smaller as the tracking object gets closer to the sensor and therefore it is difficult to triangulate objects closer to the ceiling. To overcome this problem Vantage cameras are installed in the four corners. Vantage cameras have a larger field of view and tracking range because they have higher resolution and powerful IR strobes compared to Vero cameras. Vantage cameras have an overlapping field of view with each other and the Vero cameras on the ceiling (looking top-down). Therefore, this arrangement supports the triangulation of markers closer to the ceiling.

4 RGB video cameras with a wide-angled lens are placed on top of the cage along the central axis. They do not have an overlapping field of view but arranged such that all views can be stitched together to cover the entire length of the cage. This placement is strategic and has two major advantages. 3D tracking data from the entire cage can be augmented in the video footages and behavior traits can be recorded from a top-down view. Machine learning-based object detection techniques can be deployed in the future to support of VICON tracking, especially for connecting lost trajectories. The other 2-RGB cameras are not mounted at a specific location. They are to be used *on demand* basis to collect dataset from desired perspectives (discussed later with posture problem) or to record behavior at a specific region of interest during the experiment e.g. collective feeding or mating perch.

## Software Interface

The setup is operated with two commercial software: **Tracker** and **Nexus**. Both have different functionalities for collecting and processing data. Tracker is optimized for identification and 6-DOF pose computation of unique marker patterns (explained later). It is a light-weight software designed to support real-time communication of tracking results via ethernet.

Nexus software is rather customized for collecting and annotating 3D motion capture data with humans. The software is also used to post-process tracking data e.g. labeling 3D trajectories, skeleton fitting, and filtering the 3D data. NEXUS is designed to execute user-defined scripts for data filtering. Further, the Nexus is useful for augmented visualization of tracking data on the RGB data.

These softwares are customized for the primary customers of the motion capture systems and therefore options for visualization 3D data are limited. The users have to develop a customized software interface to augment 3D data on RGB images in the desired fashion. This is essential for computer vision applications using the tracking data, videos, and calibration files [95]. To support such customizations a Software Development Kit (SDK) is provided to get access to raw data generated by the system. Users can enhance the performance of the system by developing customized algorithms for processing data from the IR cameras. A detailed explanation of the software and data formats is beyond the scope of the thesis and readers are encouraged to follow the specification guide provided by Vicon.

## Data Storgae and transmission

The tracking data is stored as a recording session directly on the master computer (in a proprietary data format). Each session is stored with relevant files (calibration, raw 2D detections, marker patterns, etc.) for reusing the data for visualization or post-processing.

Video data is dumped directly on a dedicated SSD drive (512 Gb/camera) and a separate post-processing step is performed after the recording session is complete (using Nexus software). Post-processing involves conversion from raw to *.mp4* or *.avi* format (lens distortion is not compensated). The storage capacity for video data is limited because the raw data files are rather large. A 10 Gigabit connection is set up between the master computer and back storage device. Technically, this arrangement is sufficient for scheduling daily experiments and archive the data.

The tracking data is processed in real-time and streamed over the network. Tracker uses the UDP protocol to stream the data to the desired IP address. The data can be stored separately or used to design a closed-loop experiment. The response time of the system is in a range of milliseconds.

### Calibration

A 3D calibration object with active LED markers is used to calibrate the entire system (both IR and RGB cameras). Once the calibration is done it is ideal to not change the sensor positions or ceiling height. The accuracy of the system (triangulation, pattern identification, etc.) depends highly on the calibration process. The calibration is sensitive to large shifts in temperature ($\pm 8 - 10° C$) and it is important to recalibrate the system to avoid errors.

The calibration routine for the setup requires roughly 15-20 mins. For behavior experiments, it is recommended to calibrate the system at regular intervals to maintain consistency. The software is able to assist the user with color-coded information about the accuracy of calibration (Green-Red). The core idea and principle of multi-camera calibration are covered in the introduction (see sec 2.4.6).

The calibration files are stored separated with each session recording. This file contains all necessary calibration parameters i.e. intrinsic and extrinsic parameters. The software allows users to define a reference coordinate system after the calibration is over. All camera positions in the calibration file are stored w.r.t the reference space. It is important to know that VICON follows a proprietary method to undistort video files. This step is extremely important for the development of any computer vision applications, we obtained a special tool from the company to perform undistortion of videos.

## 6.1.5 Tracking capability

The setup operates in two modes: on-line and off-line. The system is in on-line mode when it is recording activity in the cage. This mode is designed for closed-loop experiments and therefore the system streams 3D markers positions and 6-DOF pose of marker patterns (see Fig. 6.4) in real-time. In off-line mode, the user can view the recorded session and process the data or filtering it before exporting the tracking results. Additionally, the users can access in-built features to visualize trajectories, pose patterns, velocity, and acceleration of patterns of tracked objects.

To best of our knowledge, such a system is installed for the first time to track a wide range of wild animals. We conducted pilot experiments with birds, insects, mammals to understand

**Fig. 6.3.** Different animals used in the pilot experiments to check the performance of the system for a different types of animals. Except for the bee, all animals shown in the image are able to move freely in the barn with markers attached to them. The example of a bee with a 6 mm marker is displayed to highlight the importance of correct marker selection. The pigeon has 6 mm markers attached to body parts and 9 mm markers on the backpack. The starling is wearing a backpack of 6 mm markers. The hawkmoth has one 6 mm marker and other half-spherical 3 mm markers. The choice of the marker depends on the species. The performance of the system is affected by the choice of markers.

tracking features of the motion capture system. These trial experiments have shaped our understanding in terms of the range of experiments that can be performed using the space. We learned that performance of tracking varies based on the species, mainly size and speed. Optical tracking systems suffer from line-of-sight issues and the tracking is lost when animals leave the tracking volume or congregate at places causing occlusion. It is possible to tune software and hardware settings to optimize the tracking performance in a limited capacity. It is cumbersome to change hardware settings and a good understanding of technology is required to do such changes. There, it is recommended that the hardware setting of the sensors are chosen as global settings and software settings are varied during the experiment for enhanced performance.

### Marker selection

Makers come in different shapes (spherical or hemisphere) and sizes (3 mm to 25 mm diameter) (see Fig. 6.3). The weight of the marker depends on the size and manufacturer. Bigger markers increase the detection rate and assure stable tracking accuracy at larger distances. The type of markers also affects detection rate, in our experience markers with smooth reflective coating perform better than markers with layered reflective tape. The placement of markers on the animals is crucial, the behavior of the animal should not be altered by the placement of the markers and it should not hinder the range of motion of the animal. We have used 3 mm hemispheres, 6 mm and 9 mm spherical markers for preliminary tests with animals.

Markers are used as an individual unit or as part of a pattern for 6-DOF tracking (use cases are covered below). Careful consideration of pattern size and weight is required before conducting any experiment with animals. Often animals may require customized pattern designs based on their physiology. It is possible that animals will damage patterns or lose the markers causing an unexpected drop in the tracking performance. Additionally, wear and tear design factors

**Fig. 6.4.** The image shows a starling with a custom-designed backpack. The backpack has 4 markers of 6 mm size arranged in a unique pattern. It should be noted that the system can handle both 2D and 3D patterns, therefore the 3D pattern is achieved by raising the height of one of the markers. The base is made of special material to keep the backpack lightweight. The backpack is placed on a slightly raised platform to avoid occlusion of markers while flying. The markers attached close to the back of birds may show irregular tracking patterns because the line of sight is broken with one or more cameras. The image is part of a study conducted by Nora Carlson to understand the collective behavior of starlings ©Nora Carlson.

must be weighed against the robustness of tracking to while considering reuse of markers or old backpacks. It is up to the user to decide the impact of the loss of tracking data on the experimental results.

### 3D-Point Position tracking

3D position tracking of a single retro-reflective marker is one of the most basic features of the motion capture system (see Fig. 6.3). The triangulated positions of multiple markers are computed with sub-millimeter accuracy. These positions are streamed in real-time as they are computed and therefore a specific order is not maintained. However, off-line mode uses a nearest neighbor like algorithm to find temporal consistency and exports trajectories in *.csv* format. It should be noted that single marker tracking is not a commonly used feature of motion capture systems. Single markers are almost always used for posture tracking applications (discussed later).

**Preliminary trials:** A single marker can be tracked but identity can only be maintained if the tracking results are continuous. In collective behavior experiments, it is common to use trajectories as a proxy for identity (only within the same trial) [52]. This feature could be very useful for researchers who want to maintain identity within the sessions and not between sessions. Keeping this use case in mind we tried using a single marker-based tracking approach with several insects as shown in figure 6.3. We choose insects because they are small and can accommodate only one marker, for bigger animals, other solutions are most feasible.

**Preliminary observations:** We tracked insects (moths, locust, spider) in the cage for several sessions and observed the tracking data for preliminary understanding. We learn that the system is indeed capable of tracking small insects with several limitations. The exported data is in form of a series of short trajectories, where each trajectory is categorized as a separate point. We do not know the exact criterion used by the Vicon software to assign these trajectories. From studying the exported data, we can conclude that trajectories are lost if tracking is lost. Tracking is lost due to line of sight issues or due to the movement of the animals. Triangulation is weak near the corners and near the sensors because of lack of overlap between the field of view of the cameras. Flying insects (moth) tend to go outside the tracking volume or their wing motion causes occlusions. Similarly, terrestrial insects can go to the corners of the cage occlusions.

The limitations of triangulation and trajectory mapping can be improved by modifying the experimental approach or software settings to accommodate the needs of the experiment. Marker size makes a big difference as bigger markers reflect more IR light and detections are much smoother. The existing sensor arrangement is not best suited for insect tracking when we use small markers and sensors are at 4 meters height. It should be changed to limit tracking to a smaller volume and which will enhance the chances of triangulations with smaller markers. In the software, the threshold settings regarding the detection of markers and triangulation can be optimized for improving performance. For example, the distance consideration between two points for trajectory assignment can be reduced. Sampling rate of the cameras can also be increased for mored detections. Such changes may enhance false trajectory mapping if several points are moving very close to each other.

VICON system is not built for this type of application and therefore the software is not optimized for combining trajectories of a single point. We conclude that, a customized algorithms is required to fuse the exported data and create continuous trajectories. This feature is required to evaluate the performance of insect tracking in quantitative manner. It is better to quantify the tracking loss when identity is provided. We can get a sense of tracking success with individuals markers in table 6.2, where we have reported results of tracking multiple markers as part of a skeleton. In our opinion, there is a huge scope for improvement by post-processing the exported set of triangulated points. If the points are triangulated, the assignment of the trajectory can be improved with knowledge of movement patterns, speed of animals, and the total number of markers in the scene.

**Data Format:** The data can be exported as a *\*.csv* file with the 3D position of markers. Each observation is given a unique *Point ID* and stored with corresponding *Frame ID*. In the subsequent frames, *Point ID* is maintained if the measurement is identified as the same points by the trajectory matching algorithm otherwise a new ID is assigned. This way the system accumulates data in the form of an upper triangular matrix where new columns are added with each frame.

## 6-DOF Pose tracking

This is one of the most used features of the motion capture systems. A user-defined pattern of 4 or more markers is used for 6-DOF pose tracking. The markers are directly attached to the tracking object or they are attached to a base that can be attached to any target object. The marker pattern is stored as a target object with the desired name in the software. This process includes defining a local coordinate system for the pattern, $O_{obj}$, where one marker is selected as the origin, and two others are used to define the primary axis and the secondary axis. The position of each marker ($P_{obj}^{1..4}$) is stored in this local coordinate system as a *\*.mp* file. During the recording session, markers are triangulated ($P_{vicon}^i$) in the vicon space (reference space) ($O_{ref} = O_{vicon}$). A matching algorithm is used to identify the cluster of points and correspondence is established between triangulated points and saved patterns i.e. ($P_{vicon} \leftrightarrow P_{obj}$). The unique constellation of markers is used to assign these correspondence automatically and compute pose of the patten i.e. $P_{obj} = Rt_{vicon}^{obj} \cdot P_{vicon}$.

The marker constellation can be two dimensional or three dimensional. The geometric properties of the pattern such as distances and relative angles are used as unique descriptors to assign correspondences during the matching process. The ability of the software to match accurate correspondences relies heavily on the pattern design. 2D patterns have weak descriptors because the range of variation is limited to a plane. Additionally, two-dimensional patterns have to face the cameras, and the visibility of markers is poor at oblique orientations.

**Preliminary Trials:** For collective behavior experiments, we intend to use a pattern-based approach for tracking multiple individuals and maintaining their identity. We used backpacks made of planar patterns with two species of birds i.e. pigeons and starlings (see Fig. 6.3). 9 mm markers were used with pigeons in several trials with 1, 2, 4, and 8 individuals. Starlings were given 6 mm markers with groups of 10 or 15. Backpacks used with pigeons (35 mm x 91 mm) were larger than starling backpacks (30 mm x 50 mm). To evaluate performance of pattern tracking with birds, we rely on the trials where only one bird was tracked. For this

experiment, we used two styles of pattern attachments and recorded 10 different sessions with 4 different pigeons. One pattern is made of 9 mm marker and attached with a planar backpack. Another pattern is made of 6 mm markers and attached on the head directly. The same experiment is later used for evaluation of posture based tracking.
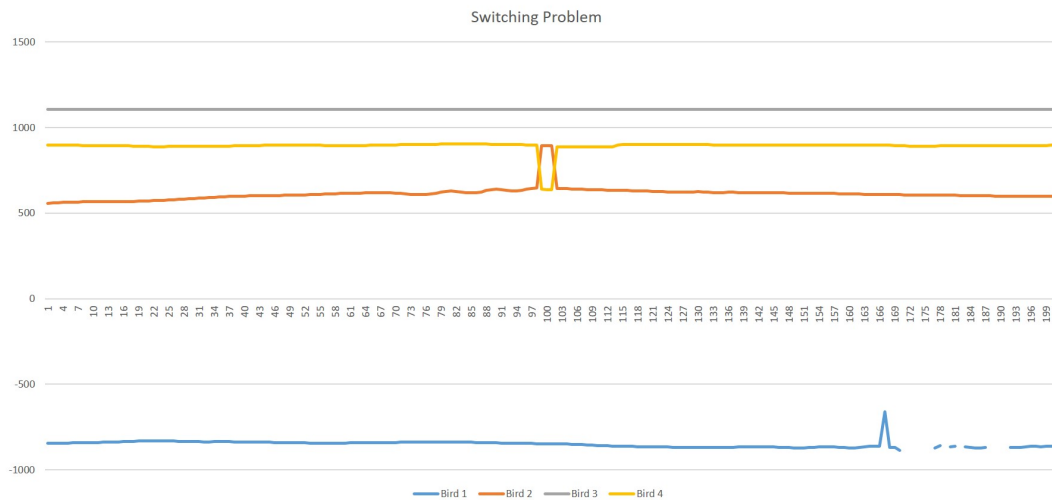
**Tab. 6.1.** Evaluation of 6-DOF pattern tracking with a single bird. The table reports number of frames where tracking result was not obtained for a particular pattern.

| Total Frames | Backpack Pattern | Head Pattern | Both |
|---|---|---|---|
| 34456 | 0 | 39 | 0 |
| 6451 | 0 | 0 | 0 |
| 38686 | 0 | 392 | 0 |
| 25584 | 0 | 1118 | 0 |
| 11369 | 0 | 68 | 0 |
| 10403 | 444 | 617 | 335 |
| 591 | 0 | 2 | 0 |
| 2893 | 354 | 538 | 257 |
| 17817 | 1 | 341 | 0 |
| 19729 | 4 | 240 | 4 |

**Preliminary observations:** In table 6.1, we have reported the tracking performance of pattern tracking with a single bird from 10 different sessions. The results were recovered from the export function provided by the Tracker software (post-processing). For the sake of simplicity, we have reported the frames where pattern tracking did not yield any results. We learn that tracking is fairly stable and results are obtained for more than 95% of frames. The head pattern is lost more often than the backpack. This is possible because birds often move or shake their head at a rapid pace and the system may not be able to capture such data or discard it due to quality of detection. It is also seen that tracking is lost for both patterns at times. The reason can be occlusion of markers while flying or the bird moving out of the tracking range (i.e. closer to sensors).

It should be noted that have a tracking result does not mean that results are always accurate. There are two types of errors observed: jittering and abrupt flipping of rotation angles flip by 90 or 180. These problems are observed when marker detection is unstable or triangulation results are erroneous (see Fig. 6.6). It can also occur when the sampling rate is not matching up to the speed of the motion. Jitter effects are removed by smoothening the trajectory using filters or increasing the sampling rate. Flipping of angles is generally the result of wrong correspondence matching for pose computation. The occurrence of these problems may depend on several factors such as bad pattern design, wrong detection, degraded calibration or bird behavior, etc. Therefore, performance may vary a lot between two different trials. We have identified probable causes of such errors in the appendix section (see appendix B).

For experiments with multiple patterns, the identity of birds remains stable over time when multiple patterns are used. One common problem is the switching of identities between different patterns, which is the result of using similar marker patterns. An easy method to spot this problem is to plot the position of each backpack as displayed in Fig. 6.5. The translation parameters of the backpacks are interchanged when identity switch happens, thus, it appears as if the bird has "*teleported*" to a new location. Discussion on this problem is also covered with other errors in the appendix section.

**Fig. 6.5.** The image shows a typical example of identity switching. We have plotted translation parameter (Tx) of four different marker patterns. It can be seen that translation parameters of *Bird 2* and *Bird 4* are swapped by exact values between frame 97 and 103 frames (6 frames).
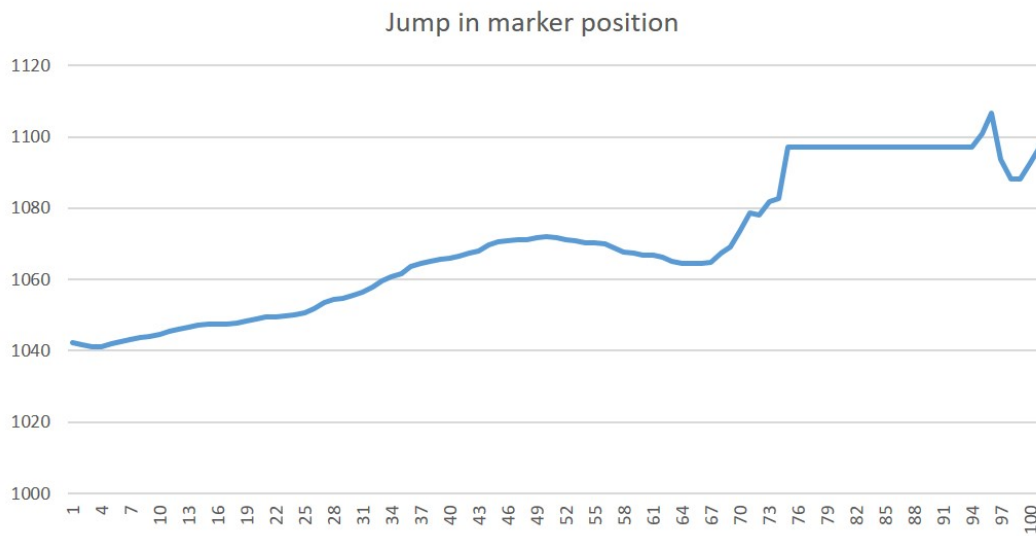
Pose jitter and switching are observed more with starlings than with the pigeons. One explanation is the behavior of birds is very different, starlings fly a lot more actively, perch outside tracking volume at higher heights, and often on the cameras which deteriorates the calibration. The probability of occlusions is higher in the case of starlings because they congregate in close vicinity of each other. This behavior may also lead to identity shifting, especially because starling backpacks offer a considerably smaller area to create unique pattern designs. Moreover, the use of smaller (6 mm) markers may also play a role in degraded tracking results.

The problems of jittering and flipping are observed more in the on-line mode. In off-line more, post-processing reduces triangulation errors and pattern matching errors dramatically. Often position and identification tracking are enough and orientation may not be required, thus flipping problems can be ignored. It should be noticed that all recordings are done at speed of 100 Hz. The errors discussed above occur for a fraction of a second (few frames) and such errors can be removed from the system without significant loss of data.
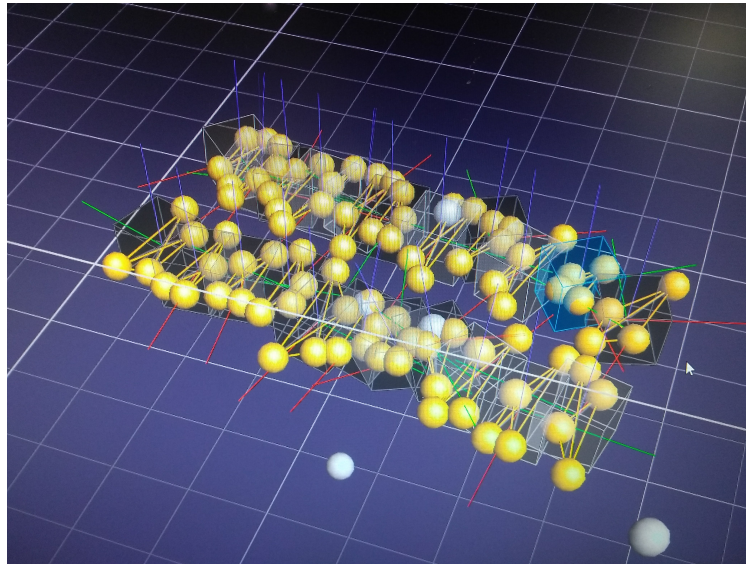
**Data Format:** The 6-DOF tracking data can be exported in multiple ways. The Tracker software exports rotation (quaternion and Euler angles) and translation parameters (in mm). The NEXUS software rather provides 3D position of each of the markers in the pattern (in vicon space $O_{vicon}$). Identity of each marker is maintained and therefore each point is labeled in the *\*.csv* file.
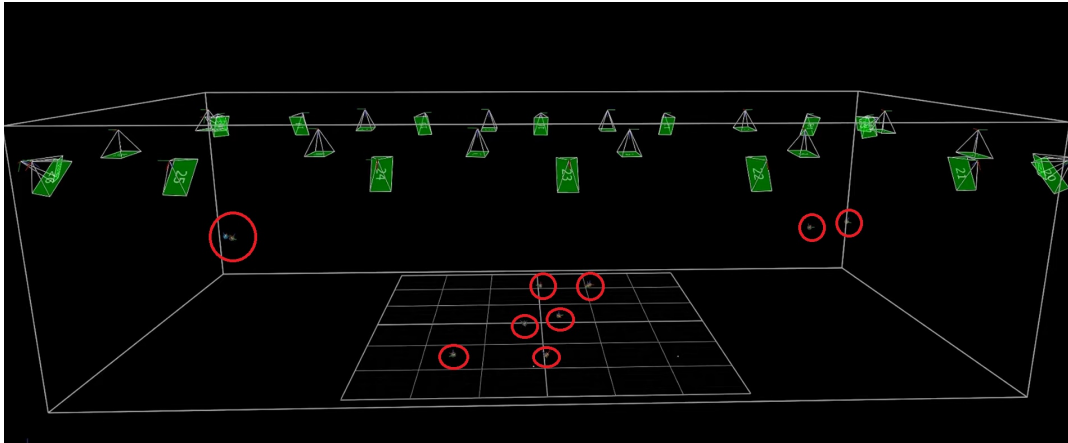
### Posture tracking

One of the major applications of motion capture is the 3D posture reconstruction of articulated bodies. The system and operating software are fully configured to support applications of human posture tracking. It is used extensively in the medical and entertainment field for capturing the full-body motion of humans (covered in chapter 1). Markers are attached to different parts of the body and assigned to an existing template of the human skeleton.

**Fig. 6.6.** The image shows a typical example of rapid change in position of a marker. We have plotted the x-coordinate of a single marker that belongs to a pattern. It can be seen that the translation value changes steadily when the birds is moving. We can also observe rapid change of position in last 6 frames which can be due to fast movement or error in triangulation.



**Fig. 6.7.** The image shows different marker patterns placed next to each other to check the performance of the system in terms of the unique identification of the patterns. The phenomenon of ghost points (in gray) is also visible in the image. Such points are created when the 2D matching algorithm of the motion capture system does not come to a consensus and creates artificial points. This is possible due to several reasons ranging from hardware errors such as bad focus or strobe intensity values to software errors such a matching threshold or calibration.

**Fig. 6.8.** The image shows a screenshot of a pilot experiment showing multiple animal tracking. 10 starlings were tracked in the cage with each starling having a unique pattern of 4 markers of 6 mm size. The image shows the 6-DOF position of the tracked backpacks.

A one-time initialization step is performed where the subject is asked to move different limbs in sequential order. The initialization step is used to add subject-specific information optimization to the skeleton i.e. position of markers w.r.t joints. After setting up the system desired sequences are recorded and exported using the NEXUS software. The final report is exported in various formats based on the requirements of the application e.g. positions for VR applications, joint angles for sports, or medical studies.

Posture tracking of animals with motion capture is attempted before [2] but not in the context of large scale collective behavior experiments. The classical approach requires a large number of markers to be attached to the body of the animal. Animals are trained and only one individual is used at a time. Adding a large number of markers is not feasible and perhaps not required for behavior studies.

**Preliminary test:** We argue that posture tracking for behavior experiment does not require full 3D reconstruction. For many experiments, it would be sufficient to have body orientations and head position. Especially, for closed-loop applications head tracking is most important to deliver the stimuli. We decided to attach markers on the body and the head for creating a simple skeleton with only 1 joint. Typically, each segment of the skeleton has to be supported with at least 3 markers. We concluded that developing a skeleton model with more segments is not feasible for birds and perhaps overkill for the task at hand. For the preliminary evaluation of posture tracking, we have used the same 10 sessions that we used for the evaluation of pattern tracking. We attached four 6 mm markers on the head and another pattern with 9 mm as a backpack (see Fig. 6.9).

**Preliminary observations:** Figure 6.9 shows the results from tracking a single pigeon with markers on the head. The image shows augmented images in the two camera views and the 3D reconstruction in the top right corner. The tracking is seamlessly maintained and trajectories of the markers can also be displayed in the offline mode. The results of preliminary evaluation are displayed in table 6.2. These results can be directly compared to the results of pattern tracking mentioned in table 6.1. The only difference is that results are obtained through a different pipeline. In this case, we linked the two patterns as part of a skeleton, a feature

**Fig. 6.9.** The image shows results of posture tracking with a single bird using the Vicon system. The bird is wearing a backpack with marker patterns and we have additionally attached markers to the head and defined them as another marker pattern. The image shows two different views of the augmented 3D data on the RGB images. We also show 3D view generated from the reconstructed marker positions.



**Fig. 6.10.** The image shows the results marker-based posture tracking with multiple birds. The pattern on the head and the body are linked to each other as part of one single skeleton and therefore accuracy of matching is high. The view on the left is the 3D view of reconstructed points and view on the right shows augmentation of the tracking results on the video images.

unique to Nexus software. Skeleton tracking results in marker positions, therefore we have reported the number of times each marker is lost. We can observe that tracking loss for the head pattern is reduced (except session 8). Also, the rate at which an individual marker is lost is very low compared to the complete pattern being undetected. There is a possibility that more two or three markers are lost in one frame, we have not looked at such combinations yet.

One interesting observation is that tracking loss in backpack pattern has increased in all cases, while the same for head pattern has decreased. These results can be due to different software implementations in Vicon. We know that the company claims Tracker to be optimized for 6-DOF pose tracking (results in table 6.1. Whereas, Nexus is optimized for posture tracking or rather skeleton matching. They use different constraints to get the results and fill gaps in

Evaluation of posture tracking using patterns with a single bird. The table reports number of frames where tracking result was not obtained for markers in a specific pattern. Two or more markers can also be lost in a same frame, however to maintain simplicity we have included only the instances where all markers were not detected.

| Total Frames | Backpack Pattern | | | | | Head Pattern | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | bp1 | bp2 | bp3 | bp4 | All | hd1 | hd2 | hd3 | hd4 | All |
| 34456 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6451 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 38686 | 230 | 79 | 79 | 79 | 79 | 118 | 79 | 79 | 79 | 79 |
| 25584 | 222 | 9 | 0 | 0 | 0 | 341 | 254 | 303 | 222 | 1 |
| 11369 | 0 | 44 | 0 | 0 | 0 | 71 | 79 | 119 | 79 | 66 |
| 10403 | 388 | 395 | 389 | 386 | 380 | 734 | 597 | 602 | 608 | 596 |
| 591 | 0 | 0 | 0 | 0 | 0 | 6 | 1 | 0 | 12 | 0 |
| 2893 | 530 | 531 | 584 | 542 | 511 | 949 | 856 | 978 | 862 | 822 |
| 17817 | 0 | 1 | 0 | 0 | 0 | 104 | 53 | 132 | 355 | 48 |
| 19729 | 5 | 19 | 5 | 43 | 5 | 363 | 73 | 129 | 70 | 65 |

trajectories. Our takeaway message was that raw data exists for such manipulation, which leaves room open for the development of customized algorithms.

Figure 6.10 shows a group of 8 pigeons being tracked with the system. Each pigeon is identified correctly and tracked while moving through space. We have combined the pattern-based pose tracking with skeleton mapping. The backpack allows us to get the 6-DOF pose of the body which can be used as a proxy for body orientation. Similarly, the pattern on the head allows us to compute the 6-DOF pose of the head. HMD based VR systems use outside-in tracking (e.g. motion capture) approaches to track the head of the user and generate perspective correct views assuming the location of the eye (2.3). Similarly, we can safely assume that the pose of the head can be used to infer the viewing direction and its field of view provided that we know the eye position of the bird. Therefore, we conclude that this arrangement is suitable for closed-loop experiments because the pose of the head pattern can be streamed in real-time through the on-line data streaming facility. We have further developed a technique to measure eye position w.r.t head pattern and track eye location which is discussed in the next section.

It is useful to note that the performance of the tracking varies depending on the configuration of the skeleton. Head of the pigeon has limited space and therefore markers can have very similar patterns when multiple birds are used. This case can cause confusion between patterns and the switching problem can be observed when birds are in close vicinity. In such cases, it is important to learn that tracking is more stable when the body and head patterns are defined as part of one single skeleton. It is possible to define both patterns separately and treat them as a different object (normal 6-DOF tracking). However, defining them as part of the same skeleton adds additional constraints for matching algorithm and this significantly improved the matching results during the post-processing pipeline. This has been observed qualitatively, we have concluded that more rigorous tests should be performed to test posture tracking with multiple birds. However, it is suitable to conduct such tests after improving the backpack designs. In the future, we plan to design an on-line experiment which will give us a better idea of performance.

**Data format:** Posture data is exported using the Nexus software as a *.csv* format. It is possible to export data in multiple ways however we export marker position for our application. The eventual file has 3D position of each marker along with frame no of the captured marker.

## 6.2  Solving computer vision problems in the Barn

The multi-sensory setup is built to conduct collective behavior experiments with animals. However, the application of the setup is not limited to behavior research. This section will focus on discussing the use of the setup for solving computer vision problems related to animal tracking. This idea overlaps with our long term strategy to extend the tracking capabilities of the existing setup by developing new methods using RGB cameras. The existing setup surely provides novel features for behavior experiments e.g. 3D position, posture, identification. Adding markerless functionality will allow us to overcome many shortcomings of the existing setup. We believe that exploiting the advantages of marker-based methods is the best strategy to develop reliable markerless solutions with ground truth. In this spirit, we have proposed a strategy to use the existing setup for solving the 3D posture tracking using a single RGB camera.

### 6.2.1  Background: Markerless tracking for behavior

In general, there are many advantages to developing markerless techniques for behavior studies. The approach is less intrusive and allows animals to behave naturally. The information obtained through videos is richer than marker-based approaches and therefore the same sequence is used for extracting multiple traits i.e. movements, postures, with temporal consistency. Initially, classical computer vision techniques were used for such tasks i.e. feature detection and image processing techniques. These methods are limited to a specific problem or a specific dataset but not useful for different species or backgrounds. Also, most techniques differ between indoors and outdoor environments.

Recently, biologists have started leveraging CNNs for quantitative animal behavior studies using computer vision in the lab [82, 136] and in the wild [82]. Research on vision problems related to detection, tracking, and posture computation is limited to 2D space. Some researchers are focusing on identification of animals using natural patterns [69, 169]. Most of these methods are developed using supervised machine learning techniques. Supervised training approaches require carefully curated datasets with a manual intervention that may not be always practical or scalable. Most importantly, generating datasets for 2D annotation are possible manually but obtaining 3D ground truth locations require completely different approach [95, 190].

### 6.2.2  3D Posture for Animals

The problem of computing the 3D posture of animals is relatively less studied. A small community of computer vision researchers has worked on this problem in the last decade. The methods and approaches depend on the input type i.e. single image [38, 163], video sequence [20], multi-view images [84, 225], depth image [172] etc. Some approaches are specific to

species e.g. flies [84], cats [103] and some approaches are applicable to multiple species e.g. quadrupeds, [226], birds [104].

We chose to focus on the problem of computing posture from a single view (images or videos) for birds. This choice is based on the current availability of the hardware and consideration of scalability in the future developments of the setup. Our existing setup has only two cameras with side views and others are mounted on the ceiling. We can use the two cameras as a stereo setup to create a multi-view posture approach. However, this will not be scalable because adding more cameras will dramatically increase the processing time for video conversion, storage, and extraction. Multi-view approaches are suitable consideration for smaller areas, however, the method entire volume of the setup is not practical in terms of setup, cost, processing requirements, etc. Our intension was to aim for a feasible solution that requires minimal processing time.

### State of the art: 3D posture prediction

Most methods using single view based posture prediction use silhouettes and keypoints to predict 3D posture [20, 226]. Zuffi *et al.* [226] introduced a model based method to compute full 3D reconstruction from a single image. They created a generic shape space model (SMAL) to represent quadrupeds using 3D scans of various toy animals. The algorithm can generate a 3D posture by deforming the SMAL model but it requires manually selected keypoints and silhouettes. The evaluation is qualitative and the method focuses on getting complete posture of the animal. This approach is inspired by an approach called SMPL where 3D scans of humans are used to create a shape space model of humans [125]. The authors used toy models because detailed 3D scanning of live animals is a difficult task. Biggs *et al.* [20] used the SMAL model and proposed a new approach for extracting posture from a video (sequence of images). The method detects keypoints and segments the shape of the animal automatically to recovered 3D shape. The final evaluation is based on the accuracy of silhouette matching with synthetic sequences. The methods presented above-used shape constraints for recovering 3D shape. Zuffi *et al.* [224] developed another approach for recovering complete posture with texture and camera position from a single image. The network is trained on a synthetic dataset created using the texture from real images. The method presents an end-to-end solution without computing a keypoint or silhouette. The evaluation is done in 2D space based on comparing the results with manual annotations.

Kanazawa *et al.* [104] proposed a category-specific approach to recover a deformable shape with texture mapping but do not model pose. The method uses annotated images to compute a mean shape and applies deformation to the mean shape during inference. The model uses keypoints, silhouettes, and texture for optimizing the loss function of the network. The model is trained on CUB-dataset [212] of birds which provides only 2D annotations and therefore evaluation is qualitative. Another notable approach is from Novotny *et al.* [162] where the 3D pose is predicted from 2D keypoints. Novotny *et al.* show that their approach works for predicting 3D pose from 2D keypoints for any category e.g. birds, humans. The approach is evaluated with Human3.6M [95] posture dataset which includes 3D ground truth of human posture. They use CUB-dataset [212] and therefore can not evaluate 3D prediction with ground truth.

We can observe that most recent work is based on qualitative evaluations where manually annotated images are considered as ground truth and error is measured in 2D space. The ground truth for such an approach is very difficult to obtain. Additionally, the existing methods focus on recovering detailed 3D posture or shape of the animal (except [162]) from the images. We argue that recovering detailed 3D posture or shape is not always necessary for our application. A simple wireframe representation of 3D body axis and head orientation may be sufficient for many behavior experiments. For supporting closed-loop applications, it is rather beneficial to select an approach that requires minimal processing time. Our problem is comparable to the problem of predicting 3D human pose (joint locations) from 2D keypoints [133, 167]. Therefore, we decided to investigate this direction for our research.

Prediction of the 3D pose from monocular image or video is a difficult problem. Recent popularity of the approach for humans is motivated by two factors *a*. Availability of large datasets with accurate ground truth [95, 190] and *b*. Robust methods to detection 2D keypoints for single and multiple humans [6, 36]. Research in 2D human posture tracking has already inspired the development of similar approaches for tracking 2D posture of animals in the lab and the wild [82, 135, 136]. However, the problem of 3D posture prediction from monocular images in not attempted yet because suitable datasets are not available. Therefore, we decided to first create a dataset suitable for solving this problem.

## 6.2.3  Generating datasets for posture prediction

For preparing the dataset we took inspiration from Human3.6M dataset [95]. It is one of the largest datasets created with the intension of supporting a set of problems related to computation of human posture e.g. posture prediction (2D, 3D, single view or multi-view), video analysis, activity recognition, etc. The dataset is captured with a marker-based motion capture system, similar to ours. 11 professional actors (6 males and 5 females) with different body sizes were used to record various activities such as walking, sitting, etc. The markers attached to their joints are used to track full-body motion in 3D. 3D joint locations and their 2D projections (from 4 different views) are recorded automatically for each sequence.

We decided to follow a similar approach to create ground truth for the prediction of bird posture. However, we were forced to chose a slightly different strategy because we model the problem differently. The human posture model is defined using 17 joint positions. The dataset containing millions of images is used to learn the relative movement between joints and predict the joint location in 3D w.r.t a root joint i.e. generally pelvic joint. Our motion capture setup is not suitable for defining such a complex skeleton model for birds. Mainly because bird physiology presents a completely different type of articulation due to wings. We decided to focus on scenarios where wings are closed i.e. walking or standing birds. For these activities, we can compute a simple posture model that is defined by body orientation and head orientation (explained earlier). Therefore, we decided to prepare a dataset that is suitable for predicting the 3D posture in terms of head and body orientation.

One key challenge was to decide the features that can be used for the prediction of the posture. For articulated models generally, joint positions are selected and the 2D-3D dataset is directly achieved using motion capture. In our case, walking birds do not have many articulated

Concept                                    Stereo Annotation

**Fig. 6.11.** Annotation tool:(left) The image displays the basic concept behind the annotation tool. The head and the body are represented as two rigid objects with distinct coordinate systems defined by vicon patterns. It is assumed that features on the head and body (marked in blue) can be tracked reliably by computing the 6-DOF motion of patterns on the head and body. The image shows that the 3D position of marker patterns are already augmented in the image. (right) The image displays a snapshot of the initialization process using stereo images. The desired features are annotated manually using the video cameras. Their positions are triangulated in camera space ($O_{cam}$) and transferred to the respective coordinate system defined by patterns on the head ($O_{head}$) and the body ($O_{body}$).

movements. Moreover, it is ideal to select features that are easily detectable in the images. For birds, such features are eyes, beak, tail, wings, etc. and we already know that these features are used in other datasets (with 2D keypoints) [212] to predict 3D shape [104] and pose [162] of birds. We decided to use the distinct features on the head and body for the prediction of 3D posture. The major problem was to get 3D locations of these features directly with motion capture systems because we markers can not be attached at these locations due to practical and ethical reasons. To solve this problem we designed an alternative strategy for generating automatic 2D-3D annotation for posture problems.

### Annotation method

We have designed an automated approach to get 2D and 3D locations of desired features. For the first iteration chose 3 features on the head (beak, left eye, and right eye) and 3 features on the body (tail, left-wing, and right-wing) as key points. However, our method is not limited to the annotation of these features. It can be extended to create annotations for more features and for other computer vision problems. Our method relies on high accuracy 6-DOF tracking of motion capture system for producing large scale datasets with minimal intervention. We have explained the main concept and pipeline in the following text.

*Concept:* Our idea is inspired by human VR approaches where the eye position of the subject is inferred from 6-DOF tracking of the headset. This concept is also applicable to birds as features like eyes and beak are attached rigidly to the head. All features on a rigid body experience the same translation and rotation. If the head is represented with a coordinate system $O_{head}$ (see Fig. 6.11) then the position of all features on the head will remain consistent in that coordinate system of the head. Similar logic can be applied to features on the body of the bird, assuming that the bird has its wings folded and it is walking or standing. We already know
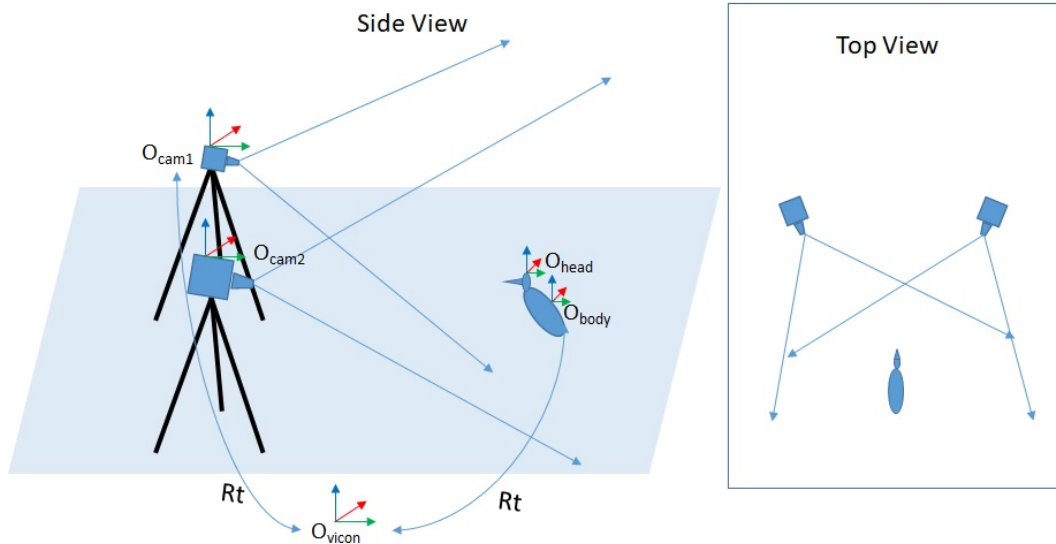
that we can attach marker patterns to track the 6-DOF motion of head and body (sec.6.1.5). Therefore, it is reasonable to conclude that features on the head and body can be represented in the coordinate system of the corresponding marker patterns. This means that the 6-DOF pose of the patterns can be used as a proxy for the motion of the head and body provided that pattern remains in the same position. However, this is only possible when the features are registered in the coordinate system of the pattern.

For the registration process, we decided to use the same video footage that is used to collect the posture dataset. This way we found a way to avoid duplication of effort and arranged the setup to support the registration process. The cameras are arranged in a way that field of view of the cameras is overlapping (see Fig. 6.12). We developed a manual annotation tool to mark the positions of features on the head and body (see Fig. 6.11). These features are triangulated and transferred to the coordinate system of the patterns for one-time registration. This is possible because the cameras are part of the vicon system and calibrated along with the IR sensors. We use the calibration information to triangulate the points and transfer them from camera space ($O_{cam}$) to vicon space ($O_{vicon}$). Furthermore, the vicon system provides the pose of the marker patterns ($Rt^{obj}_{vicon}$) for that particular frame and that is used to transfer the points from vicon space to ($O_{vicon}$) to object space ($O_{object}$). This registration remains consistent as long as the physical location of the marker patterns does not change.

*Pipeline:* The annotation pipeline is used with a recorded tracking session. We assume that backpacks and head patterns are already defined and tracking pipeline from vicon is used to generate pose of both the patterns. It is ideal to have filtered patterns to avoid tracking errors (jitter or switching) affecting the annotation process. The annotation tool reads the tracking data and projects the information of markers on the video images. The user then selects the visible key points in the image. After sufficient annotations (ideally 1), the keypoint locations are triangulated and transferred to the coordinate system of the tracking object. The pose of patterns in that particular frame must be correct. This is easy to verify while annotating the data because marker positions are augmented on the image. If the augmented positions are incorrect the user can select another frame and annotate that particular frame. The registered points are stored in a *.txt file and used with the tracking data to create automated annotations from the vicon tracking data. We have designed a script to support end-to-end annotation. Ideally, the user has to manually annotate each feature only once. The accuracy of registration depends on the accuracy of calibration and 6-DOF tracking.

*Features of existing tool:* We provide a short list of existing features of the annotation tool. The use of these features can be done beyond the scope of the posture problem (details in sec. 6.2.4)

- Compatible with vicon data formats and data structures (i.e. calibration, *.csv files etc.).

- Augmentation of tracking data on video images (marker positions, 6-DOF pose, bounding boxes etc.).

- Manual annotation (in video) of 2D features and 3D stereo triangulation.

- Registration of custom 3D features to marker patterns.

**Fig. 6.12.** The pictorial illustration of the setup used for collecting dataset for the posture problem. Two video cameras are placed approximately at a distance of 2.5-meter from each other, at the height of 2 meters from the ground. The cameras point towards the floor and have sufficient overlap in the field of view. Marker patterns are attached to the head and the body and the bird are placed in front of the cameras.

- Transformation of features between spaces (i.e. object space, camera space, vicon space).

- Automated processing of vicon tracking data for one-click annotation.

- Exporting annotations (2D keypoint, bounding box, 2D-3D correspondences).

### Setup and dataset collection

Figure 6.12 illustrates details of the setup used for collecting the dataset. We place two video cameras at an approximate distance of 2.5-meters from each other. The cameras cover bird movement from two different perspectives. The cameras are placed to have a side view of the bird, similar to human posture tracking. We chose this perspective to maximize the possibility of the overlapping field of view between the cameras. Cameras placed closer to the floor have very small overlap. The overlapping field of view is one requirement for the annotation tool. However, the benefit of having such an orientation is that the same dataset applies to the problem of multi-view pose prediction.

The bird is placed in front of the cameras and encouraged to move naturally in the area. The birds move in the range of 2 meters to 5 meters in front of the camera. Unlike humans, the motion or activity of the birds can not be controlled or dictated, therefore we distributed food on the floor to motivate the bird to move naturally. In total 12 sessions were recorded with 4 different pigeons at 100 Hz (IR cameras) and 50 Hz (video camera). We have obtained tracking data for a total of 201846 frames. Figure 6.13 shows results of the annotation tool.

### Challenges

We faced some challenges while creating this dataset which may affect the performance of the predictions. One of the major challenges is that the behavior of birds is unpredictable. The

**Fig. 6.13.** The image shows results from the annotation tool. For better visualization we have represented the eyes as one single point by taking mean of the right and left eye positions.



Tracking error        Self occlusion        Self occlusion        Rapid motion

**Fig. 6.14.** The image shows example cases where automatic annotation fails. The failure can be due to error in tracking caused by rapid motion of birds or self occlusion.

range of motion captured in our sessions is limited to walking, standing, and feeding scenarios. During these activities, birds often make certain movements which reflect as annotation errors in the dataset (see Fig. 6.14). Birds often tend to clean feathers and move ahead rapidly. Tracking of the head pattern is lost in these scenarios because of self-occlusion or poor detection (see Fig. 6.14). Currently, we are working on designing strategies to filter out the frames with tracking errors. It would be ideal to segregate the dataset based on the position of the birds in the image and the distance of the birds from cameras. This will allow us to have more control over the data given to the network for prediction of the posture.

Working with wild birds proved to be a challenging experience. We often recorded sessions where birds did not move at all and these sessions are not used for the dataset. This is where knowledge of experimental practices in biology is valuable for computer scientists. In behavior experiments, it is common to allow the animals to be familiarized with the environment. This process requires time and repeated exposure to the setup. After initial failure, we were able to capture sessions where birds performed naturally. It was difficult to restrict the motion of the bird in the camera field of view which resulted in blank images. It was equally challenging to motivate the birds to perform all possible range of behaviors. Therefore, it is likely that the dataset is biased towards certain postures. The next step is to eliminate similar postures and create a balanced dataset.

Another major challenge is to understand the data requirement of the prediction network i.e. how large the dataset should be? Our strategy is to deploy existing machine learning models used for posture prediction in humans and then expand the dataset gradually based on the results.

## 6.2.4 Ongoing work on posture prediction

We are currently using the dataset for the prediction of 3D features using a single view. This problem is being attempted for the first time and therefore it is difficult to estimate the requirement of the dataset. Therefore, we have decided to first use the dataset on the existing methods that are used for the prediction of human posture from a single view. After computing initial results we will be in a better position to understand the development of datasets for the same of other species.

We have identified two possible approaches that may be most suitable for our purpose. The first approach is "*lifting*" the 2D keypoints to predict 3D locations [133]. The approach of Martinez *et al.* [133] results in avg. error of 47.7 mm with Human3.6M dataset when 2D ground truth is used and the error increases to 67.5 mm when 2D keypoints are detected using stacked hourglass approach [157]. The second approach is to take advantage of the temporal consistency in videos [167]. Pavallo *et al.* [167] showed that errors can be reduced to 37.2 mm with ground truth and 58.6 mm 2D stacked hourglass keypoint detector. Arguably, the temporal consistency approach makes is likely to have less errors but may not be capable for real time applications. Therefore, we want to explore both option and prepare a dataset that will be useful for motivating computer vision community to work on this problem. It is worth nothing that errors in the existing models are too large when compared to the bird size ( pigeons are 300-460 mm). However, it is worth nothing that we have to design a new error model since the existing model for humans predicts avg. error using 17 joint locations.
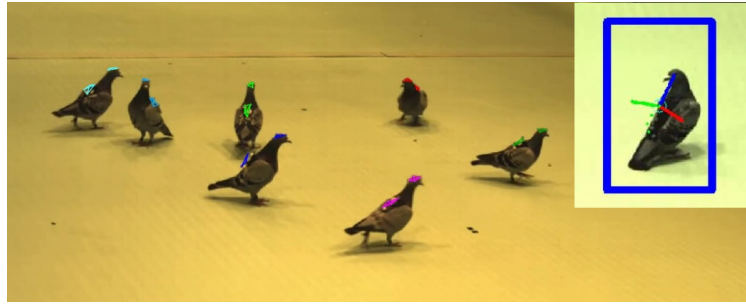
## 6.2.5 Extended use of the annotation tool

We have largely focused our discussions on capturing the data for the problem of posture prediction. However, the scope of the annotation tool extends beyond generating datasets for the posture problem. Our annotation tool is useful for applications in XR and for creating automated datasets for other computer vision i.e. multi-view posture, object detection, identification, etc.

Manual annotation is directly useful for conducting closed-loop experiments with marker-based perspective tracking. Our method can be used as a tool to compute the eye position of birds before the experiment. The time required to record a session and manually register the eye locations is less than 30 minutes (including setup, video conversion, and tracking processing). Once registered the eye position can be computed in real-time using the streaming facility provided by the software. In the next section, we have used this approach as a basis for proposing novel XR experiments with birds in large spaces.

The datasets prepared with our tool are also suitable for solving other computer vision problems. Our dataset automatically generates ground truth for 2D keypoint locations which is useful for training and validating results of 2D keypoint detectors for birds. The dataset already contains multiple views of the bird which is useful for solving the problem of 3D posture prediction from multi-view either using triangulation based methods [84] or deep learning models [96]. Furthermore, it is possible to create segmentation masks using techniques applied by Biggs *et al.* [20] for their dataset (not implemented yet). This would

**Fig. 6.15.**   The image shows multiple birds being tracked by the vicon system. The pattern on each bird is unique and identity is maintained (depicted with color code). The inset image shows a bounding box drawn around the bird. The box is generated automatically with our annotation tool. Using the same method we aim to extend our tool to annotate scenarios with multiple birds.

allow us to provide segmentation masks and key points for posture estimation. The head and body tracking allow us to create a bounding box around the pigeons (see Fig. 6.15). Therefore, the tracking data is also useful for training object detection model for bird detection. Such models will be extremely useful for our setup to eliminate pose switching errors in the case of multiple birds.

We are currently working with our collaborators to extend the capabilities of our annotation tool and create a framework for generating problem-specific datasets. Most importantly, we want to extend our method for creating new datasets with multiple animals Fig. 6.15). Besides the above-mentioned problems of detection and posture, datasets with multiple birds will be extremely useful for natural pattern-based identification. This framework development is part of our long term strategy of developing computer vision tools to reduce dependency on marker-based approaches.

## 6.3  XR applications in the Barn

In the final section of this chapter, we will present a set of novel XR applications that can be deployed using the marker-based tracking solution. Our intention behind suggesting these applications is to encourage collaboration between XR researchers and behavioral scientists. We have proposed three applications with different configurations. The first application is a behavior experiment with animals, the other two involve XR applications for humans for visualization of behavioral experiments. This way we show that the setup fits the idea of "*XR for all*", where XR applications can be developed for human and non-human animals.

### 6.3.1  Mixed Reality Application for Collective Behavior Studies

In the previous section, we explained a method to capture the 3D position of different features of the animal using a marker-based tracking approach. We already know that our setup is capable of streaming head tracking data in real-time. Therefore, it is reasonable to assume

**Fig. 6.16.** The illustration shows an experimental setup for conducting closed-loop visual stimulation based experiment with multiple birds. The setup includes a 2D display screen that is placed at a fixed location and its position is computed using the motion capture system. The idea is to track the head orientation of multiple birds and provide a stimulus to a focal individual (represented in green) when the individual is near the screen. The field of view is inferred from the eye location computed in an off-line process using the annotation tool.

that eye location and field of view can be computer in real-time. Based on this assumption we propose a simple closed-loop experiment.

We propose a small experimental setup for studying collective vigilance in birds (figure 6.16). The idea is to display simple visual stimuli (e.g. looming stimuli) to one or more individuals in a group and to record the reaction of the group. The stimulus can be triggered in closed-loop using the real-time perspective tracking solution. Generally, it is difficult to configure such an experiment in closed-loop with multiple animals because real-time perspective tracking of multiple animals is difficult. Therefore, such experiments are designed in open-loop configuration and the stimulus is triggered manually or at specific time intervals. We can also design complex visual stimuli if we assume that the response time of the system does not degrade with the additional step of eye position calculation. This approach is simple to implement and it can be extended to other species relatively easily. Especially, with animals that can carry a marker pattern on the head.

## 6.3.2 AR applications for trajectory visualization

Augmented Reality applications using outside-in tracking is not a new concept (see sec.2.2). Technical requirements for setting up such AR applications are already explored in the XR literature [49, 154, 184]. AR and VR applications are now advertised as core applications of commercial motion capture solution providers e.g. VICON [1], ART [2]. Our setup is an ideal space for designing Augmented Reality applications for displaying experimental data. Figure 6.17 shows two example use cases for the proposed application. We propose using handheld displays or optical see-through HMDs for visualization of augmented content (see Fig. 6.17). This proposed application will partially solve the limitation of the existing setup

---

[1]https://www.vicon.com/applications/location-based-virtual-reality/
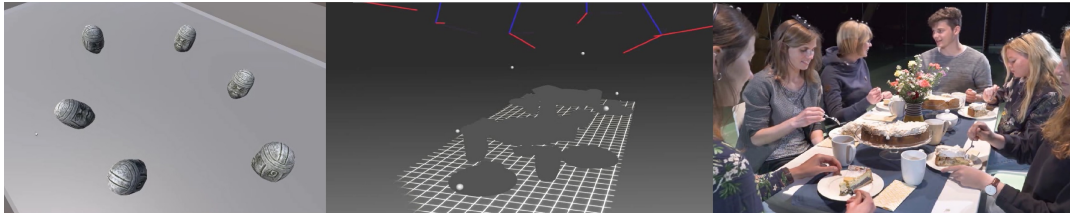[2]https://ar-tracking.com/applications/augmented-reality/

**Fig. 6.17.** The illustration shows two concepts of using Augmented Reality (AR) for visualization of experimental data. (left) The display medium is a handheld display (e.g. mobile, tablet) that is tracked in 6-DOF using the motion capture system. (right) The display medium is optical-see through HMD that is tracked in 6-DOF by the motion capture system. The augmented view in the display may contain information on bird identity, location, and motion trajectory.

i.e. augmented data is visible only in off-line mode. Another limitation is that augmented content is only viewable in proprietary software with limited options that are not customized to visualize behavior data.

There are many advantages of using such visualization techniques in the context of behavior studies. Our setup is one of the first setup designed to track the 3D movement of a group of animals. Therefore, new visualization techniques are necessary for a better understanding of data which can be provided with AR. The user can walk into the cage and point the display at the animals to see experimental data in the augmented view. The augmentations may include movement trajectories, id numbers, experiment duration, the health of the bird, etc. The same information can be directly augmented in the field of view of the user with optical-see through headsets e.g. hololens or Aryzon cardboard. In a more complex version, the user can also see the visual field of view of the animal in the augmented view if the eye position is computed in real-time (covered in the previous section). The tracking data can also be visualized in an empty cage after the experiments are over. There are many possibilities of using the augmented visuals for experiment and beyond. For example, organizing science outreach activities, especially for young scientists and school children where a group of users can visualize data together and interact with it. Such implementations can be cost-effective when developed with smartphones or smartphone mountable cardboard headsets (e.g. ZapBox, Holokit, Aryzon).

## 6.3.3  Remote demonstration of experiments using XR

In the previous applications, we described XR applications using the real-time tracking of the facility. It is also possible to use the experimental data for organizing a remote demonstration

**Fig. 6.18.** The image displays the concept of remote VR visualization of recordings in the barn. We show three different visualizations of the same scene. (right) The image displays a scene recorded in the barn where multiple humans are dining on the table. (center) The image displays the head orientation of each person in the augmented view. (left) Image displays the scene recreated from a different angle using the Unity engine and 3D tracking data of the head movements. The images are part of a pilot experiment conducted by the Psychology dept., the University of Konstanz in collaboration with Max Planck Institute of Animal Behavior.

of the experiment using VR or AR (see Fig. 6.18). The most common way of demonstrating experimental setup and the findings is the use of 2D video data. We argue that this experience can be enhanced by using VR or AR techniques and development with our setup does not require significant investments e.g. flying with storks application [110]. The experiment can be viewed by a group of collaborators in remote space or a smaller version of the setup can be created in the living room using inbuilt AR tool-kits of the smartphones. The techniques required to do this are already available. The virtual view displayed in Fig. 6.18) is created by directly importing the motion capture data in the Unity framework. Modern motion capture systems already provide plug-ins for rendering engines such as Unity. Furthermore, Unity provides options to directly generate VR and AR applications for smart-phones or VR headsets.

Let us build up on the previously discussed example of collective vigilance experiment with a group of birds. Assuming that 6-DOF head and body movement of multiple birds is already available. We know the dimensions of the cage and have complete calibration data of the sensors. Therefore, we can use digital models of the animals and recreate the complete experiment virtually. The virtual view can be fully interactive and include information about field of view or movement trajectory of the animal. It should also be possible to generate virtual views from the perspective of the animal. We believe that such rich visualizations will not only benefit researchers but also contribute towards research activities. For example, it will be possible to generate synthetic datasets with real animals and arbitrary backgrounds for solving computer vision problems in the wild [95, 224].

## 6.4  Conclusion

In this section, we provide a short summary of the work presented in this chapter and conclude with an outlook towards future research activities. We introduced a new setup for conducting collective behavior experiments. We provided details of the setup and preliminary observation. The setup is not perfectly customized for our needs, but we have obtained a reasonable understanding of using it for behavior experiments. We have also identified the developments required to further support behavior experiments i.e. new algorithms for filtering data, removing switching errors, and designing customized marker patterns.

The setup is also built with the intention of pursuing sensory stimulation based experiments in large areas. Closed-loop experiments with visual stimuli require perspective tracking and this problem is not well studied with a focus on animals. We treat perspective tracking as a subset of the problem of posture tracking. We adopted a two-step approach to work on this problem with a primary focus on birds. As a first step, we designed a marker-based approach for tracking the posture of birds. The marker-based solution works accurately and it can be already used for designing closed-loop experiments. The second step of our strategy is to develop a markerless method for tracking posture. We learned that there is no ground truth available to train machine learning models for this problem. Therefore, we designed a method to compute large scale datasets by leveraging the marker-based posture tracking. We have a reliable method to compute ground truth with minimal manual intervention. Using our posture problem as an example we also show that our setup can be used for collecting ground truth for solving other computer vision problems specific to animals. Currently, we are using the dataset created for posture to predict the posture of birds using images from a single camera. In the near future, we aim to publish the dataset and the results of our posture prediction techniques.

On a broader scale, we can claim that our setup is one of the first steps towards designing a space for conducting a wide range of experiments with a wide range of species. As of now, the performance of the motion capture system fulfills the requirements of collective behavior experiments. However, there is a strong need for collaborations with computer vision and the XR community to overcome existing limitations and extend the range of experiments. We show that the existing setup is capable of providing a reliable framework for solving challenging problems in computer vision and developing exciting applications for XR.

# Part IV

Conclusion

# Conclusion

<span style="color:blue; font-size:2em; text-align:right;">7</span>

We worked on two different projects to make tools for closed-loop visual stimulation for human and non-human animals. In the first project, we developed AR solutions for improving industrial manufacturing practices. And the second project involved the development of tracking strategies for designing new XR solutions for animals. We have already discussed problem-specific conclusions earlier in Part 2 and Part 3. Therefore, we would like to present our research from the perspective of an XR tool designer and guide the discussion towards workflow specific development of XR applications.

At first glance, applications in industrial manufacturing and animal behavior do not seem to have much in common. One common link between them is that both use tools based on XR technologies. The researchers in these fields have remained positive about using XR technologies for various applications. In Part 2, we already discussed that industrial researchers have made significant monetary investments to deploy AR solutions. Similarly, we also discussed, in Part 3, that biologists have made significant investments in customizing VR solutions to conduct fundamental research with animals. Interestingly, both the fields have experienced very different trajectories regarding the implementations of XR solutions. Industrial applications have been one of the key application domains for XR technology developers since the beginning of XR research. Yet, the technology is used sparingly in the industrial domain [70] (more details in chapter 3). In contrast, XR (mainly VR) technologies are used extensively with animals [149] but contributions from the XR technology developers are rare which has resulted in the slow growth of novel experimental techniques or setups (see chapter 5). In our research, we highlighted these facts and proposed that both the fields will potentially benefit much more from working in close collaboration with XR developers. This recommendation is also valid for technology developers in XR that are trying to promote the use of their technology.

Our research can be considered as an attempt to put this advice into practice. The main focus was given towards developing solutions customized to the needs of the end-users. This involved working with the end-user and learning details of the existing working practices in each of these application domains. In both projects, we found two questions to be very important: what are important considerations of the user while selecting a technology? And why the user wanted a solution to be designed in a specific manner? For example, industrial users wanted to implement new AR solutions but did not want to change existing workflows or compromise inaccuracy. At the same time, they wanted solutions to be scalable, cost-effective, and user friendly. We studied established practices in manufacturing (e.g. RPS for alignment) and worked with industrial partners to replicate the core working principles of existing practices to design a stable alternative solution using AR. We have not conducted user-specific tests to prove the efficacy of the solution, however, our solutions have been part of a commercially sold product for the last few years and they are well accepted by industrial customers globally [67].

On a similar note, a major technological limitation in behavior research is to modify the technological components that are developed for the visual stimulation of humans (i.e. tracking, visualization techniques) and use them with animals. Therefore, it is ideal (for behavior researchers) to work with technology experts to make the best possible customizations. At the same time, it is vital (for technology experts) to understand the requirements of the experiment along with functional knowledge of sensory organs of the animals to propose the best possible XR solutions. We realized that there was a substantial gap in the knowledge exchange between the two communities. With our project, we have tried to highlight the importance of bridging this knowledge gap. First, we conducted a technology review where we learned that existing experimental techniques with closed-loop stimulation are not extended to larger areas and more animals due to lack of suitable tracking infrastructure (see chapter 5). After learning core requirements (see appendix A), we proposed a new setup for studying collective behavior. We proposed modifications in motion capture technology and customized existing tracking approaches to suit the needs for animal tracking. Moreover, we proposed several new directions for developing XR applications (for humans and non-human animals) using the tracking setup. As of now, we have started using this facility for conducting open-loop behavior experiments. We hope that collaboration with technology experts will allow us to conduct closed-loop experiments in the near future.

In closing remarks, we can say that existing solutions and practices must be studied to replacing them with novel solutions and techniques. Every step taken in the direction of understanding the user's needs will bring the XR developers closer to developing a novel and improved solution.

# Part V

Appendix

# Survey on requirements of collective behavior studies

We took a survey to understand requirements of the researchers studying collective behavior. These questions were intended to learn their preferred choice of methods while doing an experiment. We also included questions regarding their technical capabilities and their willingness to collaborate with technology developers. All the questions are listed below with responses.

## Question 1

Rate the importance of features you wish to have in a setup that is developed for studying collective behavior? (1 less important and 5 very important)

**Tab. A.1.** The table shows votes given by each researchers for their preferred features in the setup.

| Features | Score | | | | | |
|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | N/A |
| Larger tracking volume ($> 3$ $m^3$) | 0 | 1 | 0 | 5 | 6 | 1 |
| 3D Position measurement ($>$ 50 hz) | 0 | 2 | 3 | 3 | 4 | 1 |
| 3D Posture measurement | 0 | 5 | 2 | 3 | 2 | 1 |
| Real-time tracking | 3 | 2 | 1 | 6 | 1 | 0 |
| Automated tracking results | 0 | 0 | 1 | 3 | 7 | 2 |
| Automated movement trajectories | 0 | 0 | 1 | 4 | 6 | 2 |
| Ability to track different species | 0 | 0 | 5 | 3 | 2 | 2 |
| Individual tracking within trial | 0 | 0 | 2 | 1 | 9 | 1 |
| Individual tracking between trials | 1 | 1 | 0 | 5 | 6 | 0 |
| Tracking accuracy | 0 | 0 | 0 | 5 | 8 | 0 |

## Question 2

How important is it to maintain identity of individuals across different trials? (1 = less important, 5 = very important)

- 1 - 0%

- 2 - 7.7%

- 3 - 30.8%

- 4 - 38.5%

- 5 - 23.1%

## Question 3

How important is it to maintain identity of individuals across different trials? (1 = less important, 5 = very important)

- 1 - 0%

- 2 - 7.7%

- 3 - 15.4%

- 4 - 38.5%

- 5 - 23.1%

## Question 4

Do you think that measurement of 3D movements is important for conducting collective behavior experiments? (1 = less important, 5 = very important)

- 1 - 0%

- 2 - 15.4%

- 3 - 15.4%

- 4 - 38.5%

- 5 - 30.8%

## Question 5

To the best of your knowledge, rate the technical difficulty of setting up an experiment with 3D measurement techniques? (1 = less difficult, 5 = very difficult)

- 1 - 0%

- 2 - 0%

- 3 - 15.4%

- 4 - 30.8%

- 5 - 53.8%

# Question 6

Do you agree that existing research in collective behavior studies is not oriented towards 3D measurement of movement or posture? Why?

- 1- Lack of methods for 3D tracking - 75%

- 2 - Technical complexity of implementing methods - 91.7%

- 3 - Lack of collaboration with technology experts - 50%

- 4 - Lac of training datasets - 16.7%

- 5 - No, I do not agree - 0%

# Question 7

How important is it to conduct experiments with closed-loop stimulation techniques e.g. VR? (1 = less important, 5 = very important)

- 1 - 0%

- 2 - 15.4%

- 3 - 46.2%

- 4 - 23.1%

- 5 - 15.4%

# Question 8

To best of your knowledge, rate the difficulty of designíng a setup that is capable of conducting closed-loop experiments with animals using visual stimulation (VR or AR). (1 = less difficult, 5 = very difficult)

- 1 - 0%

- 2 - 7.7%

- 3 - 23.1%

- 4 - 23.1%

- 5 - 61.5%

# Question 10

Do you use video recordings for studying collective behavior? If yes, Do you prefer to use machine learning techniques for automated detection of animals or tracking the animals?

- Yes, I prefer automated methods - 100%

- No, I prefer manual tracking - 0%

- No, I do not use video recordings - 0%

# Question 11

Have you created annotations for machine learning algorithms? If yes, did you manually annotate images?

- Yes, I did manual annotations - 58.3%

- No, I used unsupervised learning - 0%

- No, I do not use machine learning - 41.7%

# Question 12

Would you prefer to collaborate with computer scientists for developing various methods to study collective behavior?

- Yes, I prefer automated methods - 92.3%

- No, I prefer to develop my own methods - 7.7%

# Tracking errors and open problems <span style="color:blue">B</span>

In Chapter 6, we introduced a new setup for studying collective behavior. We installed a motion capture system for tracking the animal movements in 3D. We encountered several errors when using marker patterns with animals, especially 6-DOF pose computation errors such as jitter or pose flipping and pattern identification errors. We know that some of these errors can be solved by developing customized methods to process the data. These customizations are necessary for obtaining meaningful output for behavior research. In the next sections, we will discuss the problems observed and their possible causes. After that we will discuss a set of open problems that require attention in near future to boost the performance of the system.

## B.1 Tracking Errors

The accuracy of pattern identification and pose computation depends on two factors: marker detection and correspondence matching. These two factors are affected by a range of settings in both hardware and software. Hardware settings include operating focus of the cameras, opening angle, aperture etc. as well as physical location of the cameras and calibration parameters. For motion capture systems hardware settings are optimized while setting up the system by moving markers of various sizes in the tracking volume. We did not change hardware settings during our trials, although we will predict probable hardware solutions while discussing tracking problems. In software, the detection and matching results are controlled by a wide set of threshold settings e.g. circle detection settings, triangulation settings, strobe intensity, frame rate etc. There is no specific setting for all tracking problems and a trade-off must be made based on requirements of the experiment. It is convenient to capture the data and change settings during post-processing to get best possible results. Detail explanation of each settings is provided in the software manual.

The following points include discussion on each problem and show their relation with detection and matching error.

- **Incorrect 2D detection** of marker positions contribute towards 3D triangulation error which in turns results in pose with higher residual errors. Detection error is due to unstable centroid detection in cases where markers appear blurry or small due to fast motion or distance from the sensor. At oblique angles markers in the 2D pattern often occlude each other or merge as on single point based on a certain view. Such errors are cause very deviations in pose and which leads to problem of jittering. The problem can be corrected to an extent by filtering algorithm in post-processing (provided in software). For real-time applications filtering can be done using information from previous frames with kalman or other similar filters. Altering detection settings or strobe intensity in the

software may alleviate the problem. Larger markers and 3D patterns may bring huge improvement in this type of errors.

- **Ghost points** are false 3D points that actually do not exist but created by the system due to incorrect triangulation. Figure 6.7 displays this problem clearly (gray points). This occurs when multiple markers are placed in close vicinity of each other, the rays projected from two different 2D projections intersect in the 3D space and considered as a valid marker. Often 2D projection of multiple markers are very close to each other and they are detected as single marker. These point act as noisy detections resulting in ghost points that are very close to real points. *Ghost points* jitter effect if they are considered for pose computation. They also create confusion for pattern matching algorithm which results into abrupt switches in identity. This affect is seen often in starling dataset, because the markers are very close to each other. Creating 20 unique patterns in a 30 mm x 50 mm area has proven to be a challenging problem.

  Ghost points perspective dependent and therefore not consistent but appear as noisy points between frames. The ghost points also be reduced by adding constraint on triangulation parameters which may result into loss of data. One of the major reasons is error in calibration. While operating at low margin of error when we use very small 2D patterns and therefore calibration accuracy becomes paramount. We observed that temperature changes increased error in calibration. This problem is well known and vicon provides internal settings for compensation of temperature. We have not evaluated exact accuracy of these compensations and have relied on their recommendations in the manual. Also, birds perching on the cameras or flying into the cameras is not good for the maintaining accuracy. These considerations are absolutely essential when considering study species.

- **Incorrect 3D-3D correspondences** within the same pattern leads problem of flipping. Let us assume that marker $P_{obj}^1$ defined the origin of object and the correspondence order used by the system is $P_{vicon}^{1,2,3,4} \leftrightarrow P_{obj}^{1,2,4,3}$ instead of $P_{vicon}^{1,2,3,4} \leftrightarrow P_{obj}^{1,2,3,4}$. This can occur due to wrong triangulation and use of symmetric patterns. The probability of this error depends completely on pattern design (2D-3D, size), number of patterns in use and most importantly movement of the animal. It is also possible that one of the points is occluded, i.e. $P_{vicon}^4$ is not detected and $P_{vicon}^{1,2,3} \leftrightarrow P_{obj}^{1,2,4}$ are used for pose computation. In both cases the behavior of error is predictable. The rotation error of angles will be large in consecutive frames and translation errors with be small. This is because pose computation algorithm always give results that minimize the distance between points after transformation i.e. $dist(P_{vicon}, (Rt_{obj}^{vicon} \cdot P_{obj}))$. In special cases when the origin point is matched correctly ($P_{vicon}^1 \leftrightarrow P_{obj}^1$) and other are matched incorrectly the rotation error flips with 90° or larger angles. In such cases translation error can not be larger than distance between the two most distant points considered for pose computation. This is because of the distance constraint applied during pose optimization. We can numerically prove this however, the error caused by ghost points or presence of many 3D points from different patterns in same vicinity. In such cases, pose error may behave like flipping problem or jitter problem depending on the constellation. Corrective measure to avoid this error is to design distinct patterns and increase the

thresholds that affecting 3D-3D matching. Patterns with more markers can also solve this problem but that would lead to increased size and weight of the backpack.

- **Incorrect pattern matching** is basically identity detection problem where one pattern is mistaken as another. This error usually leads to abrupt *pose switching* between two or more patterns present in the tracking volume. This error can be caused by lack of diversity in pattern design and/or because of ghost points. VICON softwares are able to identify symmetric looking patterns and indicate it to the user. However, during execution ghost points or temporary occlusions may also cause pattern mismatch. Best measure to avoid this issue is to design unique patterns.

# B.2  Open problems

Motion capture systems are considered a standard tracking solution for many applications that require high accuracy and real-time performance e.g. medicine, industrial manufacturing. We learned that the system is driven to its limitations when it is used for animal tracking. Following strategies can be considered to obtain error free tracking results. However, correction of some tracking errors in on-line mode remains a challenging task.

**Improved marker design:** As a first step, a clear strategy is required to design new marker patterns and verify their performance. The setup is rigid and camera positions are fixed, therefore it should be possible to design a mock simulation to estimate performance of pattern at different distances and angles. Evaluation of marker patterns is well studies in AR and VR research and existing concepts can be applied in collaboration with AR/VR researchers. We proposed a simple idea to create 3D patterns by modifying height of one of the markers (see figure 6.4). This method has yielded stable results with pigeons in posture tracking experiments (mentioned later). One of the collaborators has recently conducted a set of experiments with starlings using these patterns.

**Develop data filtering algorithms:** The flipping and switching problems show systematic patterns which can be detected and removed by designing data filtering algorithms. Additional conditions can be applied to constrain the data based on properties of animal motion. For example, the back of the bird will never be facing the floor and therefore normal of the plane can be monitored The switching often happens between two known patterns which means that displacement must be mutual and detectable.

**Data fusion and markerless tracking:** The existing approaches rely heavily on detection of marker and identification of pattern. The position of animal is lost when markers are not detected or animals go out of tracking range. We argue that it is possible to use RGB videos to recover the lost data. Currently, Moreover, RGB based animal detection and identification technique may also help in removing identity switching problems. Strong contribution from computer vision community is required to deploy these solutions.

# List of Authored and Co-authored Publications

<span style="float:right">C</span>

**2020**

[149]  **Hemal Naik**, Renaud Bastian, Nassir Navab, and Iain D Couzin. "Animals in Virtual Environments". *IEEE Transactions on Visualization and Computer Graphics 26 (5), 2073 - 2083*.

**2019**

[82]  Jacob M Graving, Daniel Chae, **Hemal Naik**, Liang Li, Ben Koger, Blair R Costelloe, and Iain D Couzin. "DeepPoseKit, a software toolkit for fast and robust animal pose estimation using deep learning". *eLife, 2019*.

**2018**

[111]  Jens C Koblitz, Oren Frokosh, Mate Nagy, Nora Carlson, **Hemal Naik**, and Iain D Couzin. "Turning birds into bats—Multi-modal tracking to study collective behaviour". *The Journal of the Acoustical Society of America 144 (3), 1886-1886*.

**2016**

[148]  **Hemal Naik**, Mahmoud Bahaa, Federico Tombari, Peter Keitler, and Nassir Navab. "Frustration Free Pose Computation For Spatial AR Devices in Industrial Scenario". *IEEE International Symposium on Mixed and Augmented Reality, 2016, Merida, Mexico*.

**2015**

[178]  Christoph Resch, **Hemal Naik**, Peter Keitler, Steven Benkhardt, and Gudrun Klinker. "On-site Semi-Automatic Calibration and Registration of a Projector-Camera System Using Arbitrary Objects With Known Geometry". *IEEE Transactions on Visualization and Computer Graphics 21 (11), 1211-1220*.

[151]  **Hemal Naik**, Federico Tombari, Peter keitler, and Nassir Navab. "A Step Closer To Reality: Closed Loop Dynamic Registration Correction in SAR". *IEEE International Symposium on Mixed and Augmented Reality, 2015, Fukuoka, JPN*.

# Abstracts of Publications not Discussed in this Thesis

<div style="text-align: right;">D</div>

## Exploiting Photogrammetric Targets for Industrial AR

Hemal Naik, Yuji Oyamada, Peter Keitler, Nassir Navab.

In this work, we encourage the idea of using Photogrammetric targets for object tracking in Industrial Augmented Reality (IAR). Photogrammetric targets, especially uncoded circular targets, are widely used in the industry to perform 3D surface measurements. Therefore, an AR solution based on the uncoded circular targets can improve the work flow integration by reusing existing targets and saving time. These circular targets do not have coded patterns to establish unique 2D-3D correspondences between the targets on the model and their image projections. We solve this particular problem of 2D-3D correspondence of non-coplanar circular targets from a single image. We introduce a Conic pair descriptor, which computes the Eucledian invariants from circular targets in the model space and in the image space. A three stage method is used to compare the descriptors and compute the correspondences with up to 100% precision and 89% recall rates. We are able to achieve tracking performance of 3 FPS (2560x1920 pix) to 8 FPS (640×480 pix) depending on the camera resolution and the targets present in the scene.

# Turning birds into bats—Multi-modal tracking to study collective behaviour

Jens C Koblitz, Oren Frokosh, Mate Nagy, Nora Carlson, Hemal Naik, and Iain D Couzin.

Acoustic localization has been used to track numerous vocalizing animals, including whales, bats birds by measuring the time of arrival differences of a sound recorded by multiple receivers. This method, however, requires the species of interest to vocalize regularly to achieve decent temporal resolution. In order to track the movements of birds that do not vocalize regularly, a miniature, radio controlled ultrasonic speaker is attached to the birds and constantly emits ultrasound chirps. The chirps from various tagged animals can be discriminated based on information encoded in the signal. A 30 microphone array covering the ceiling of a large aviary ($14.8 \times 6.6 \times 3.9$ m) is used to record the chirps and provides the basis for accurate localization of the sources. As the chirps do not overlap with the frequency of the animals vocalizations, these can be localized and assigned individually in a flock of moving birds. In addition, a VICON motion capture system provides a simple and accurate method to ground truth the acoustic localizations as it records very precise movement information of the individuals while line of sight is between the marker on the animal and a number of cameras is established. The acoustic and visual tracking systems working together provides a unique and novel answer to consistent and precise localizations of non-vocalizing individuals (and groups of individuals) as they move through space.

# On-site semi-automatic calibration and registration of a projector-camera system using arbitrary objects with known geometry

Christoph Resch, **Hemal Naik**, Peter Keitler, Steven Benkhardt, and Gudrun Klinker

In the Shader Lamps concept, a projector-camera system augments physical objects with projected virtual textures, provided that a precise intrinsic and extrinsic calibration of the system is available. Calibrating such systems has been an elaborate and lengthy task in the past and required a special calibration apparatus. Self-calibration methods in turn are able to estimate calibration parameters automatically with no effort. However they inherently lack global scale and are fairly sensitive to input data. We propose a new semi-automatic calibration approach for projector-camera systems that - unlike existing auto-calibration approaches - additionally recovers the necessary global scale by projecting on an arbitrary object of known geometry. To this end our method combines surface registration with bundle adjustment optimization on points reconstructed from structured light projections to refine a solution that is computed from the decomposition of the fundamental matrix. In simulations on virtual data and experiments with real data we demonstrate that our approach estimates the global scale robustly and is furthermore able to improve incorrectly guessed intrinsic and extrinsic calibration parameters thus outperforming comparable metric rectification algorithms.

# DeepPoseKit, a software toolkit for fast and robust animal pose estimation using deep learning

Jacob M Graving, Daniel Chae, Hemal Naik, Liang Li, Ben Koger, Blair R Costelloe, and Iain D Couzin

Quantitative behavioral measurements are important for answering questions across scientific disciplines—from neuroscience to ecology. State-of-the-art deep-learning methods offer major advances in data quality and detail by allowing researchers to automatically estimate locations of an animal's body parts directly from images or videos. However, currently available animal pose estimation methods have limitations in speed and robustness. Here, we introduce a new easy-to-use software toolkit, DeepPoseKit, that addresses these problems using an efficient multi-scale deep-learning model, called Stacked DenseNet, and a fast GPU-based peak-detection algorithm for estimating keypoint locations with subpixel precision. These advances improve processing speed >2x with no loss in accuracy compared to currently available methods. We demonstrate the versatility of our methods with multiple challenging animal pose estimation tasks in laboratory and field settings—including groups of interacting individuals. Our work reduces barriers to using advanced tools for measuring behavior and has broad applicability across the behavioral sciences.

# Bibliography

[1] J. Abdeljalil, M. Hamid, O. Abdel-mouttalib, et al. "The optomotor response: A robust first-line visual screening method for mice". en. In: *Vision Research* 45.11 (May 2005), pp. 1439–1446 (cit. on pp. 102, 105).

[2] K. Abson and I. Palmer. "Motion capture: capturing interaction between human and animal". In: *The Visual Computer* 31.3 (Mar. 2015), pp. 341–353 (cit. on p. 139).

[3] M. Adcock, M. Hutchins, and C. Gunn. *Augmented Reality Haptics: Using ARToolKit for Display of Haptic Applications*. 2003 (cit. on p. 19).

[4] V. AG. "Reference Point System - RPS". In: (1996) (cit. on pp. 59, 78).

[5] O. Akgul, H. Penekli, and Y. Genc. "Applying Deep Learning in Augmented Reality Tracking". In: *2016 12th International Conference on Signal-Image Technology and Internet-Based Systems (SITIS)*. Los Alamitos, CA, USA: IEEE Computer Society, Dec. 2016, pp. 47–54 (cit. on p. 44).

[6] M. Andriluka, L. Pishchulin, P. Gehler, and B. Schiele. "2D Human Pose Estimation: New Benchmark and State of the Art Analysis". In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. June 2014 (cit. on p. 144).

[7] R. C. Arkin, R. C. Arkin, et al. *Behavior-based robotics*. MIT press, 1998 (cit. on pp. 8, 100).

[8] B. Arnaldi, P. Guitton, and G. Moreau. "Virtual Reality and Augmented Reality". en. In: (), p. 375 (cit. on pp. 27, 29).

[9] A. Artaud and M. Richards. *The Theater and Its Double*. Evergreen book. Grove Press, 1958 (cit. on p. 27).

[10] S. Audet, M. Okutomi, and M. Tanaka. "Direct image alignment of projector-camera systems with planar surfaces". In: *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*. 00021. IEEE, 2010, pp. 303–310 (cit. on p. 87).

[11] R. T. Azuma. "A Survey of Augmented Reality". en. In: (), p. 48 (cit. on pp. 11, 13, 14, 16, 19).

[12] M. Bajura and U. Neumann. "Dynamic registration correction in video-based augmented reality systems". In: *Computer Graphics and Applications, IEEE* 15.5 (1995). 00213, pp. 52–60 (cit. on pp. 72, 84, 86–88).

[13] S. Balasubramanian, Y. M. Chukewad, J. M. James, G. L. Barrows, and S. B. Fuller. "An Insect-Sized Robot That Uses a Custom-Built Onboard Camera and a Neural Network to Classify and Respond to Visual Input". en. In: *2018 7th IEEE International Conference on Biomedical Robotics and Biomechatronics (Biorob)*. Enschede: IEEE, Aug. 2018, pp. 1297–1302 (cit. on p. 119).

[14] S. Baldauf, H. Kullmann, T. ThüNken, S. Winter, and T. Bakker. "Computer animation as a tool to study preferences in the cichlid *Pelvicachromis taeniatus*". en. In: *Journal of Fish Biology* 75.3 (Aug. 2009), pp. 738–746 (cit. on pp. 102, 107).

[15] D. Bandyopadhyay, R. Raskar, and H. Fuchs. "Dynamic shader lamps: Painting on movable objects". In: *Augmented Reality, 2001. Proceedings. IEEE and ACM International Symposium on*. 00175. IEEE, 2001, pp. 207–216 (cit. on pp. 21, 22, 83).

[16] P. Bateson and K. N. Laland. "Tinbergen's four questions: an appreciation and an update". en. In: *Trends in Ecology & Evolution* 28.12 (Dec. 2013), pp. 712–718 (cit. on p. 100).

[17] M. Bauer. "Tracking errors in augmented reality". PhD thesis. Technical University Munich, Germany, 2007 (cit. on pp. 64, 83).

[18] A. T. D. Bennett, I. C. Cuthill, J. C. Partridge, and E. J. Maier. "Ultraviolet vision and mate choice in zebra finches". In: *Nature* 380.6573 (Apr. 1996), pp. 433–435 (cit. on pp. 102, 104).

[19] C. Bichlmeier. "Immersive, Interactive and Contextual In-Situ Visualization for Medical Applications". Dissertation. München: Technische Universität München, 2010 (cit. on p. 16).

[20] B. Biggs, T. Roddick, A. Fitzgibbon, and R. Cipolla. "Creatures great and SMAL: Recovering the shape and motion of animals from video". en. In: *arXiv:1811.05804 [cs]* (Nov. 2018). arXiv: 1811.05804 (cit. on pp. 142, 143, 149).

[21] M. Billinghurst. "A Survey of Augmented Reality". en. In: (), p. 106 (cit. on pp. 14, 16, 23, 40).

[22] M. Billinghurst. "The Reality of Augmented Reality: Are we there yet?" en. In: (), p. 93 (cit. on p. 30).

[23] O. Bimber et al. "Projector-based augmentation". In: *Book Chapter in Emerging Technologies of Augmented Reality: Interfaces and Design* (2006), pp. 64–89 (cit. on pp. 13, 23).

[24] O. Bimber, F. Coriand, A. Kleppe, E. Bruns, S. Zollmann, and T. Langlotz. "Superimposing Pictorial Artwork with Projected Imagery". en. In: *IEEE MultiMedia* (2005), p. 11 (cit. on p. 21).

[25] O. Bimber, A. Emmerling, and T. Klemmer. "Embedded entertainment with smart projectors". In: *Computer* 38.1 (2005), pp. 48–55 (cit. on p. 22).

[26] O. Bimber, B. Fröhlich, D. Schmalstieg, and L. M. Encarnação. "The virtual showcase". In: *ACM SIGGRAPH 2006 Courses*. ACM, 2006, p. 9 (cit. on pp. 13, 19, 21–23).

[27] O. Bimber and R. Raskar. "Modern approaches to augmented reality". In: *ACM SIGGRAPH 2006 Courses*. ACM, 2006, p. 1 (cit. on pp. 13, 14, 21, 22).

[28] A. Bisazza, A. De santi, and G. Vallortigara. "Laterality and cooperation: mosquitofish move closer to a predator when the companion is on their left side". In: *Animal Behaviour* 57.5 (1999), pp. 1145–1149 (cit. on pp. 102, 106).

[29] G. Bleser and D. Stricker. "Advanced tracking through efficient image processing and visual-inertial sensor fusion". In: *2008 IEEE Virtual Reality Conference*. 2008, pp. 137–144 (cit. on p. 40).

[30] T. Blum, V. Kleeberger, C. Bichlmeier, and N. Navab. "mirracle: An augmented reality magic mirror system for anatomy education". en. In: *2012 IEEE Virtual Reality (VR)*. Costa Mesa, CA, USA: IEEE, Mar. 2012, pp. 115–116 (cit. on pp. 14, 19, 44, 48, 129).

[31] C. J. Bohil, B. Alicea, and F. A. Biocca. "Virtual reality in neuroscience research and therapy". en. In: *Nature Reviews Neuroscience* 12.12 (Dec. 2011), pp. 752–762 (cit. on pp. 102, 108).

[32] A. Borst. "Drosophila's View on Insect Vision". en. In: *Current Biology* 19.1 (Jan. 2009), R36–R47 (cit. on pp. 102, 105).

[33] H. Bülthoff. "Drosophila mutants disturbed in visual orientation". In: *Biological Cybernetics* 45.1 (1982), pp. 63–70 (cit. on pp. 108, 109).

[34] T. Butkowski, W. Yan, A. M. Gray, R. Cui, M. N. Verzijden, and G. G. Rosenthal. "Automated Interactive Video Playback for Studies of Animal Communication". en. In: *Journal of Visualized Experiments* 48 (Feb. 2011) (cit. on pp. 102, 108).

[35] X. Cao and R. Balakrishnan. "Interacting with dynamically defined information spaces using a handheld projector and a pen". In: *Proceedings of the 19th annual ACM symposium on User interface software and technology*. ACM, 2006, pp. 225–234 (cit. on p. 22).

[36] Z. Cao, T. Simon, S.-E. Wei, and Y. Sheikh. "Realtime Multi-Person 2D Pose Estimation using Part Affinity Fields". en. In: *arXiv:1611.08050 [cs]* (Apr. 2017). arXiv: 1611.08050 (cit. on pp. 47, 144).

[37] L. CARMICHAEL. "The Study of Instinct. N. Tinbergen. New York: Oxford Univ. Press, 1951." In: *Science* 115.2990 (1952), pp. 438–439 (cit. on p. 101).

[38] T. J. Cashman and A. W. Fitzgibbon. "What Shape Are Dolphins? Building 3D Morphable Models from 2D Images". en. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35.1 (Jan. 2013), pp. 232–244 (cit. on p. 142).

[39] J. Y. Chen and G. Fragomeni, eds. *Virtual, Augmented and Mixed Reality. Applications and Case Studies: 11th International Conference, VAMR 2019, Held as Part of the 21st HCI International Conference, HCII 2019, Orlando, FL, USA, July 26–31, 2019, Proceedings, Part II*. en. Vol. 11575. Lecture Notes in Computer Science. Cham: Springer International Publishing, 2019 (cit. on p. 7).

[40] Y. Chen and E. S. Liu. "A Path-Assisted Dead Reckoning Algorithm for Distributed Virtual Environments". In: *2015 IEEE/ACM 19th International Symposium on Distributed Simulation and Real Time Applications (DS-RT)*. 2015, pp. 108–111 (cit. on p. 38).

[41] L. Chouinard-Thuly, S. Gierszewski, G. G. Rosenthal, et al. "Technical and conceptual considerations for using animated stimuli in studies of animal behavior". en. In: *Current Zoology* 63.1 (Feb. 2017), pp. 5–19 (cit. on pp. 102, 114).

[42] D. L. Clark and G. W. Uetz. "Video image recognition by the jumping spider, Maevia inclemens (Araneae: Salticidae)". In: *Animal Behaviour* 40.5 (1990), pp. 884–890 (cit. on pp. 102, 106).

[43] T. S. Collett. "Vision: simple stereopsis". In: *Current Biology* 6.11 (1996), pp. 1392–1395 (cit. on p. 116).

[44] I. D. Couzin, J. Krause, N. R. Franks, and S. A. Levin. "Effective leadership and decision-making in animal groups on the move". In: *Nature* 433.7025 (Feb. 2005), pp. 513–516 (cit. on p. 126).

[45] C. Cruz-Neira, D. J. Sandin, T. A. DeFanti, R. V. Kenyon, and J. C. Hart. "The CAVE: audio visual experience automatic virtual environment". In: *Communications of the ACM* 35.6 (June 1992), pp. 64–72 (cit. on pp. 32, 101, 108).

[46] R. B. D'EATH. "Can video images imitate real stimuli in animal behaviour experiments?" In: *Biological Reviews* 73.3 (1998), pp. 267–292 (cit. on pp. 102, 104–106, 114–116).

[47] H. Dahmen, V. L. Wahl, S. E. Pfeffer, H. A. Mallot, and M. Wittlinger. "Naturalistic path integration of *Cataglyphis* desert ants on an air-cushioned lightweight spherical treadmill". en. In: *The Journal of Experimental Biology* 220.4 (Feb. 2017), pp. 634–644 (cit. on pp. 108, 112, 117).

[48] A. Datta, J.-S. Kim, and T. Kanade. "Accurate Camera Calibration using Iterative Refinement of Control Points". en. In: (), p. 8 (cit. on p. 51).

[49] H. G. Debarba, M. E. de Oliveira, A. Lädermann, S. Chagué, and C. Charbonnier. "Augmented Reality Visualization of Joint Movements for Physical Examination and Rehabilitation". In: *2018 IEEE Conference on Virtual Reality and 3D User Interfaces (VR)*. 2018, pp. 537–538 (cit. on p. 151).

[50] N. A. Del Grosso, J. J. Graboski, W. Chen, E. B. Hernández, and A. Sirota. "Virtual Reality system for freely-moving rodents." In: *bioRxiv* (2017), p. 161232 (cit. on pp. 108, 114).

[51] A. I. Dell, J. A. Bender, K. Branson, et al. "Automated image-based tracking and its application in ecology". en. In: *Trends in Ecology & Evolution* 29.7 (July 2014), pp. 417–428 (cit. on p. 106).

[52] A. I. Dell, J. A. Bender, K. Branson, et al. "Automated image-based tracking and its application in ecology". In: *Trends in Ecology & Evolution* 29.7 (2014), pp. 417–428 (cit. on pp. 126, 134).

[53] T. Denayer, T. Stöhr, and M. Van Roy. "Animal models in translational medicine: Validation and prediction". In: *New Horizons in Translational Medicine* 2.1 (2014), pp. 5–11 (cit. on p. 100).

[54] M. Dill, R. Wolf, and M. Heisenberg. "Visual pattern recognition in Drosophila involves retinotopic matching". en. In: *Nature* 365.6448 (Oct. 1993), pp. 751–753 (cit. on pp. 108, 109).

[55] S. DiVerdi and T. Hollerer. "Image-space Correction of AR Registration Errors Using Graphics Hardware". In: *IEEE Virtual Reality Conference (VR 2006)*. 2006, pp. 241–244 (cit. on p. 87).

[56] F. L. Dolins, K. Schweller, and S. Milne. "Technology advancing the study of animal cognition: using virtual reality to present virtually simulated environments to investigate nonhuman primate spatial cognition". en. In: *Current Zoology* 63.1 (Feb. 2017), pp. 97–108 (cit. on pp. 102, 108).

[57] D. A. Dombeck, A. N. Khabbaz, F. Collman, T. L. Adelman, and D. W. Tank. "Imaging Large-Scale Neural Activity with Cellular Resolution in Awake, Mobile Mice". en. In: *Neuron* 56.1 (Oct. 2007), pp. 43–57 (cit. on pp. 108, 112).

[58] D. A. Dombeck and M. B. Reiser. "Real neuroscience in virtual worlds". en. In: *Current Opinion in Neurobiology* 22.1 (Feb. 2012), pp. 3–10 (cit. on pp. 102, 108).

[59] D. A. Dombeck, C. D. Harvey, L. Tian, L. L. Looger, and D. W. Tank. "Functional imaging of hippocampal place cells at cellular resolution during virtual navigation". en. In: *Nature Neuroscience* 13.11 (Nov. 2010), pp. 1433–1440 (cit. on pp. 108, 112).

[60] C. Domnisoru, A. A. Kinkhabwala, and D. W. Tank. "Membrane potential dynamics of grid cells". en. In: *Nature* 495.7440 (Mar. 2013), pp. 199–204 (cit. on pp. 108, 112).

[61] E. Dotto. "Drawing Hands. The Themes of Representation in Steinberg and Escher's Images". en. In: *Proceedings* 1.9 (Nov. 2017), p. 1090 (cit. on pp. 12, 13).

[62] B. Drost and S. Ilic. "3D Object Detection and Localization Using Multimodal Point Pair Features". In: *2012 Second International Conference on 3D Imaging, Modeling, Processing, Visualization Transmission*. 2012, pp. 9–16 (cit. on pp. 48, 81).

[63] P. Edgcumbe, P. Pratt, G.-Z. Yang, C. Nguan, and R. Rohling. "Pico lantern: A pick-up projector for augmented reality in laparoscopic surgery". In: *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2014*. Springer, 2014, pp. 432–439 (cit. on p. 22).

[64] P. Eisert, K. Polthier, and J. Hornegger. "A Mathematical Model and Calibration Procedure for Galvanometric Laser Scanning Systems". In: () (cit. on p. 49).

[65] M. C. Escher and J. W. Vermeulen. *Escher on escher exploring the infinite*. 1989 (cit. on pp. 8, 13).

[66] C. S. Evans and P. Marler. "On the use of video images as social stimuli in birds: audience effects on alarm calling". In: *Animal Behaviour* 41.1 (1991), pp. 17–26 (cit. on pp. 102, 104, 106).

[67] *Extend3D GmbH : www.extend3d.com* (cit. on pp. 23, 55, 64, 65, 69, 157).

[68] A. Febretti, A. Nishimoto, T. Thigpen, et al. "CAVE2: a hybrid reality environment for immersive simulation and information analysis". en. In: ed. by M. Dolinsky and I. E. McDowall. Burlingame, California, USA, Mar. 2013, p. 864903 (cit. on pp. 32–34).

[69] A. C. Ferreira, L. R. Silva, F. Renna, et al. "Deep learning-based methods for individual recognition in small birds". In: *Methods in Ecology and Evolution* n/a.n/a (). eprint: `https://besjournals.onlinelibrary.wiley.com/doi/pdf/10.1111/2041-210X.13436` (cit. on p. 142).

[70] P. Fite-Georgel. "Is there a reality in industrial augmented reality?" In: *Mixed and Augmented Reality (ISMAR), 2011 10th IEEE International Symposium on*. IEEE, 2011, pp. 201–210 (cit. on pp. 4, 14, 55–57, 157).

[71] A. Flack, M. Nagy, W. Fiedler, I. D. Couzin, and M. Wikelski. "From local collective behavior to global migratory patterns in white storks". In: *Science* 360.6391 (2018), pp. 911–914. eprint: `https://science.sciencemag.org/content/360/6391/911.full.pdf` (cit. on p. 126).

[72] J. Foster, P. Nuyujukian, O. Freifeld, et al. "A freely-moving monkey treadmill model". In: *J. of Neural Engineering* 11.4 (2014), p. 046020 (cit. on p. 128).

[73] J. Fraser. "A NEW VISUAL ILLUSION OF DIRECTION". In: *British Journal of Psychology* 2 (1908), pp. 297–320 (cit. on p. 9).

[74] W. Friedrich. "ARVIKA Augmented Reality for Development, Production, and Service". en. In: (), p. 13 (cit. on p. 13).

[75] S. N. Fry, M. Bichsel, P. Müller, and D. Robert. "Tracking of flying insects using pan-tilt cameras". In: *Journal of Neuroscience Methods* 101.1 (2000), pp. 59–67 (cit. on pp. 105, 113).

[76] S. N. Fry, N. Rohrseitz, A. D. Straw, and M. H. Dickinson. "TrackFly: Virtual reality for a behavioral system analysis in free-flying fruit flies". en. In: *Journal of Neuroscience Methods* 171.1 (June 2008), pp. 110–117 (cit. on pp. 102, 105, 107, 108, 112, 113).

[77] Y. Furukawa and J. Ponce. "Accurate Camera Calibration from Multi-View Stereo and Bundle Adjustment". en. In: (), p. 8 (cit. on p. 51).

[78] A. E. F. D. Gama, T. M. Chaves, L. S. Figueiredo, et al. "MirrARbilitation: A clinically-related gesture recognition interactive tool for an AR rehabilitation system". In: *Computer Methods and Programs in Biomedicine*. Oct. 2016 (cit. on p. 17).

[79] G. Geiger. "Optomotor responses of the fly Musca domestica to transient stimuli of edges and stripes". en. In: *Kybernetik* 16.1 (1974), pp. 37–43 (cit. on p. 105).

[80] R. Gerlai, Y. Fernandes, and T. Pereira. "Zebrafish (Danio rerio) responds to the animated image of a predator: Towards the development of an automated aversive task". en. In: *Behavioural Brain Research* 201.2 (Aug. 2009), pp. 318–324 (cit. on pp. 102, 107).

[81] N. D. Glossop and Z. Wang. "Laser projection augmented reality system for computer-assisted surgery". en. In: *International Congress Series* 1256 (June 2003). 00045, pp. 65–71 (cit. on p. 13).

[82] J. M. Graving, D. Chae, H. Naik, et al. "DeepPoseKit: a software toolkit for fast and robust pose estimation using deep learning". In: *eLife* 8 (2019), e47994 (cit. on pp. 119, 120, 126, 128, 142, 144, 171).

[83] J. R. Gray, V. Pawlowski, and M. A. Willis. "A method for recording behavior and multineuronal CNS activity from tethered insects flying in virtual space". en. In: *Journal of Neuroscience Methods* 120.2 (Oct. 2002), pp. 211–223 (cit. on pp. 10, 101, 102, 105, 108, 111).

[84] S. Günel, H. Rhodin, D. Morales, J. Campagnolo, P. Ramdya, and P. Fua. "DeepFly3D, a deep learning-based approach for 3D limb and appendage tracking in tethered, adult *Drosophila*". In: *eLife* 8 (Oct. 2019). Ed. by T. O'Leary, R. L. Calabrese, and J. W. Shaevitz, e48571 (cit. on pp. 120, 128, 142, 143, 149).

[85] R. Hartley and A. Zisserman. *Multiple view geometry in computer vision*. Cambridge, UK; New York: Cambridge University Press, 2003 (cit. on pp. 15, 30, 37, 43, 45–47, 49, 51, 68).

[86] C. D. Harvey, F. Collman, D. A. Dombeck, and D. W. Tank. "Intracellular dynamics of hippocampal place cells during virtual navigation". en. In: *Nature* 461.7266 (Oct. 2009), pp. 941–946 (cit. on pp. 108, 112).

[87] D. Hatsushika, K. Nagata, and Y. Hashimoto. "SCUBA VR: Submersible-Type Virtual Underwater Experience System". In: *2019 IEEE Conference on Virtual Reality and 3D User Interfaces (VR)*. 2019, pp. 962–963 (cit. on pp. 28–30).

[88] M. L. Heilig. *Sensorama Simulator Patent No. US3050870A* (cit. on pp. 11, 25).

[89] S. Hess, E. Oberhummer, R. Burlaud, et al. "Animated images as a tool to study visual communication: a case study in a cooperatively breeding cichlid". en. In: *Behaviour* 151.12-13 (Oct. 2014), pp. 1921–1942 (cit. on pp. 102, 107).

[90] R. L. Holloway. "Registration errors in augmented reality systems". 00145. PhD thesis. Citeseer, 1995 (cit. on p. 86).

[91] C. Holscher. "Rats are able to navigate in virtual environments". en. In: *Journal of Experimental Biology* 208.3 (Feb. 2005), pp. 561–569 (cit. on p. 111).

[92] B. K. P. Horn. "Closed-form solution of absolute orientation using unit quaternions". en. In: *Journal of the Optical Society of America A* 4.4 (Apr. 1987), p. 629 (cit. on pp. 47, 77).

[93] R. Inger, J. Bennie, T. W. Davies, and K. J. Gaston. "Potential biological and ecological effects of flickering artificial light". In: *PloS one* 9.5 (2014), e98631 (cit. on p. 115).

[94] C. C. Ioannou, V. Guttal, and I. D. Couzin. "Predatory Fish Select for Coordinated Collective Motion in Virtual Prey". In: *Science* 337.6099 (2012), pp. 1212–1215. eprint: `https://science.sciencemag.org/content/337/6099/1212.full.pdf` (cit. on pp. 102, 107, 108, 120, 126).

[95] C. Ionescu, D. Papava, V. Olaru, and C. Sminchisescu. "Human3.6M: Large Scale Datasets and Predictive Methods for 3D Human Sensing in Natural Environments". en. In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 36.7 (July 2014), pp. 1325–1339 (cit. on pp. 128, 130, 142–144, 153).

[96] K. Iskakov, E. Burkov, V. Lempitsky, and Y. Malkov. "Learnable Triangulation of Human Pose". en. In: *arXiv:1905.05754 [cs]* (May 2019). arXiv: 1905.05754 (cit. on p. 149).

[97] N. T. Jafferis, E. F. Helbling, M. Karpelson, and R. J. Wood. "Untethered flight of an insect-sized flapping-wing microscale aerial vehicle". en. In: *Nature* 570.7762 (June 2019), pp. 491–495 (cit. on pp. 8, 100, 119).

[98] T. A. Jenssen. "Female response to filmed displays of Anolis nebulosus (Sauria, Iguanidae)". In: *Animal Behaviour* 18 (1970), pp. 640–647 (cit. on pp. 102, 106).

[99] T. Johnson and H. Fuchs. "Real-Time Projector Tracking on Complex Geometry Using Ordinary Imagery". en. In: (), p. 8 (cit. on p. 22).

[100] T. Johnson, G. Welch, H. Fuchs, E. La Force, and H. Towles. "A distributed cooperative framework for continuous multi-projector pose estimation". In: *Virtual Reality Conference, 2009. VR 2009. IEEE*. 00019. IEEE, 2009, pp. 35–42 (cit. on pp. 22, 23).

[101] P. Joshi. *Jigs and Fixtures*. Tata McGraw-Hill Education, 1998 (cit. on p. 59).

[102] S. Jung and T. Whangbo. "Study on inspecting VR motion sickness inducing factors". In: *2017 4th International Conference on Computer Applications and Information Processing Technology (CAIPT)*. 2017, pp. 1–5 (cit. on p. 25).

[103] A. Kanazawa, S. Kovalsky, R. Basri, and D. Jacobs. "Learning 3D Deformation of Animals from 2D Images". en. In: *Computer Graphics Forum* 35.2 (May 2016), pp. 365–374 (cit. on pp. 128, 143).

[104] A. Kanazawa, S. Tulsiani, A. A. Efros, and J. Malik. "Learning Category-Specific Mesh Reconstruction from Image Collections". en. In: (), p. 20 (cit. on pp. 128, 143, 145).

[105] P. Karashchuk, K. L. Rupp, E. S. Dickinson, et al. "Anipose: a toolkit for robust markerless 3D pose estimation". In: *bioRxiv* (2020). eprint: `https://www.biorxiv.org/content/early/2020/05/29/2020.05.26.117325.full.pdf` (cit. on p. 128).

[106] P. Keitler. "Management of tracking: and tracking accuracy in industrial augmented reality environments". PhD thesis. Technical University Munich, 2011 (cit. on p. 64).

[107] K. Kim, J. Hwang, H. Zo, and H. Lee. "Understanding users' continuance intention toward smartphone augmented reality applications". In: *Information Development* 32.2 (2016), pp. 161–174. eprint: `https://doi.org/10.1177/0266666914535119` (cit. on p. 21).

[108] B. A. Klein, J. Stein, and R. C. Taylor. "Robots in the service of animal behavior". en. In: *Communicative & Integrative Biology* 5.5 (Sept. 2012), pp. 466–472 (cit. on pp. 104, 120).

[109] G. Klein and D. Murray. "Parallel Tracking and Mapping for Small AR Workspaces". In: *2007 6th IEEE and ACM International Symposium on Mixed and Augmented Reality*. 2007, pp. 225–234 (cit. on p. 47).

[110] K. Klein, B. Sommer, H. Nim, et al. "Fly with the flock: immersive solutions for animal movement visualization and analytics". English. In: *Journal of the Royal Society Interface* 16.153 (Apr. 2019) (cit. on p. 153).

[111] J. C. Koblitz, O. Frokosh, M. Nagy, N. Carlson, H. Naik, and C. Iain. "Turning birds into bats—Multi-modal tracking to study collective behaviour". In: *The Journal of the Acoustical Society of America* 144.3 (2018), pp. 1886–1886. eprint: https://doi.org/10.1121/1.5068270 (cit. on p. 171).

[112] J. Krause, A. F. Winfield, and J.-L. Deneubourg. "Interactive robots in experimental biology". en. In: *Trends in Ecology & Evolution* 26.7 (July 2011), pp. 369–375 (cit. on pp. 100, 104, 119, 120).

[113] F. Kretschmer, M. Tariq, W. Chatila, B. Wu, and T. C. Badea. "Comparison of optomotor and optokinetic reflexes in mice". en. In: *Journal of Neurophysiology* 118.1 (July 2017), pp. 300–316 (cit. on pp. 102, 105).

[114] A. Kulcke, C. Gurschler, C. Gasser, and A. Niel. "Image Processing Based Calibration of High Precision Laser Projection Systems". In: *INTERNATIONAL ARCHIVES OF PHOTOGRAMMETRY REMOTE SENSING AND SPATIAL INFORMATION SCIENCES* 34.3/B (2002), pp. 134–137 (cit. on p. 49).

[115] S. N. Kundu, N. Muhammad, and F. Sattar. "Using the augmented reality sandbox for advanced learning in geoscience education". In: *2017 IEEE 6th International Conference on Teaching, Assessment, and Learning for Engineering (TALE)*. 2017, pp. 13–17 (cit. on p. 22).

[116] R. Künzler and T. C. Bakker. "Female preferences for single and combined traits in computer animated stickleback males". In: *Behavioral Ecology* 12.6 (2001), pp. 681–685 (cit. on pp. 102, 107).

[117] M. Lai, S. Skyrman, C. Shan, et al. "Fusion of augmented reality imaging with the endoscopic view for endonasal skull base surgery; a novel application for surgical navigation based on intraoperative cone beam computed tomography and optical tracking". en. In: *PLOS ONE* 15.1 (Jan. 2020). Ed. by I. H. El-Sayed, e0227312 (cit. on p. 20).

[118] N. LaLone, S. A. Alharthi, and Z. O. Toups. "A Vision of Augmented Reality for Urban Search and Rescue". In: *Proceedings of the Halfway to the Future Symposium 2019*. HTTF 2019. Nottingham, United Kingdom: Association for Computing Machinery, 2019 (cit. on p. 20).

[119] M. Land. "Eye movements in man and other animals". en. In: *Vision Research* 162 (Sept. 2019), pp. 1–7 (cit. on pp. 102, 105).

[120] T. Landgraf, D. Bierbach, H. Nguyen, N. Muggelberg, P. Romanczuk, and J. Krause. "RoboFish: increased acceptance of interactive robotic fish with realistic eyes and natural motion patterns by live Trinidadian guppies". en. In: *Bioinspiration & Biomimetics* 11.1 (Jan. 2016), p. 015001 (cit. on pp. 100, 104, 120, 124).

[121] J. Larsch and H. Baier. "Biological Motion as an Innate Perceptual Mechanism Driving Social Affiliation". In: *Current Biology* 28.22 (Nov. 2018), 3523–3532.e4 (cit. on p. 121).

[122] T. Le, M. Nguyen, and T. Nguyen. "Human posture recognition using human skeleton provided by Kinect". In: *2013 International Conference on Computing, Management and Telecommunications (ComManTel)*. 2013, pp. 340–345 (cit. on pp. 47, 129).

[123] K. A. Leighty and D. M. Fragaszy. "Primates in cyberspace: using interactive computer tasks to study perception and action in nonhuman animals". en. In: *Animal Cognition* 6.3 (Sept. 2003), pp. 137–139 (cit. on pp. 102, 107).

[124] V. Lepetit, F. Moreno-Noguer, and P. Fua. "EPnP: An Accurate O(n) Solution to the PnP Problem". In: *International Journal Of Computer Vision* 81 (2009), pp. 155–166 (cit. on p. 46).

[125] M. Loper, N. Mahmood, J. Romero, G. Pons-Moll, and M. J. Black. "SMPL: A Skinned Multi-Person Linear Model". In: *ACM Trans. Graphics (Proc. SIGGRAPH Asia)* 34.6 (Oct. 2015), 248:1–248:16 (cit. on p. 143).

[126] J. Lugrin, S. Oberdorfer, M. E. Latoschik, A. Wittmann, C. Seufert, and S. Grafe. "VR-Assisted vs Video-Assisted Teacher Training". In: *2018 IEEE Conference on Virtual Reality and 3D User Interfaces (VR)*. 2018, pp. 625–626 (cit. on p. 28).

[127] T. Luhmann, S. Robson, S. Kyle, and J. Boehm. *Close-Range Photogrammetry and 3D Imaging*. Walter de Gruyter GmbH, 2019 (cit. on pp. 37, 42–44, 48, 49, 66, 67, 71, 80).

[128] T. Luo, Z. Liu, Z. Pan, and M. Zhang. "A Virtual-real Occlusion Method Based on GPU Acceleration for MR". In: *2019 IEEE Conference on Virtual Reality and 3D User Interfaces (VR)*. 2019, pp. 1068–1069 (cit. on p. 30).

[129] M. MA, P. Fallavollita, S. Habert, S. Weider, and N. Navab. "Device and System Independent Personal Touchless User Interface for Operating Rooms". In: *International Conference on Information Processing in Computer-Assisted Interventions (IPCAI)*. June 2016 (cit. on pp. 16, 18, 44).

[130] M. S. Madhav, R. P. Jayakumar, F. Savelli, H. T. Blair, N. J. Cowan, and J. J. Knierim. "Place cells in virtual reality dome reveal interaction between conflicting self-motion and landmark cues". In: *Society for Neuroscience*. Chicago, IL, USA, Oct. 2015 (cit. on pp. 108, 112).

[131] S. Manjrekar, S. Sandilya, D. Bhosale, S. Kanchi, A. Pitkar, and M. Gondhalekar. "CAVE: An Emerging Immersive Technology – A Review". In: *2014 UKSim-AMSS 16th International Conference on Computer Modelling and Simulation*. 2014, pp. 131–136 (cit. on p. 34).

[132] S. Mann, J. C. Havens, J. Iorio, Y. Yuan, and T. Furness. "All Reality: Values, taxonomy, and continuum, for Virtual, Augmented, eXtended/MiXed (X), Mediated (X,Y), and Multimediated Reality/Intelligence". en. In: (2018), p. 10 (cit. on pp. 10–12, 25, 27).

[133] J. Martinez, R. Hossain, J. Romero, and J. J. Little. "A Simple Yet Effective Baseline for 3d Human Pose Estimation". en. In: *2017 IEEE International Conference on Computer Vision (ICCV)*. Venice: IEEE, Oct. 2017, pp. 2659–2668 (cit. on pp. 144, 149).

[134] A. Mathis, P. Mamidanna, T. Abe, et al. "Markerless tracking of user-defined features with deep learning". en. In: *arXiv:1804.03142 [cs, q-bio, stat]* (Apr. 2018). arXiv: 1804.03142 (cit. on p. 47).

[135] A. Mathis, P. Mamidanna, K. M. Cury, et al. "DeepLabCut: markerless pose estimation of user-defined body parts with deep learning". en. In: *Nature Neuroscience* 21.9 (Sept. 2018), pp. 1281–1289 (cit. on pp. 120, 128, 144).

[136] M. W. Mathis and A. Mathis. "Deep learning tools for the measurement of animal behavior in neuroscience". en. In: *arXiv:1909.13868 [cs, q-bio]* (Oct. 2019). arXiv: 1909.13868 (cit. on pp. 120, 142, 144).

[137] P. Mealy. *Virtual & augmented reality for dummies*. en. 1st edition. Indianapolis, IN: John Wiley and Sons, 2018 (cit. on p. 7).

[138] D. Mehta, O. Sotnychenko, F. Mueller, et al. "Single-Shot Multi-person 3D Pose Estimation from Monocular RGB". en. In: *2018 International Conference on 3D Vision (3DV)*. Verona: IEEE, Sept. 2018, pp. 120–130 (cit. on p. 47).

[139] C. Menk, E. Jundt, and R. Koch. "Evaluation of Geometric Registration Methods for Using Spatial Augmented Reality in the Automotive Industry". In: *15th International Workshop on Vision, Modeling, and Visualization, VMV 2010, Siegen, Germany, November 15-17, 2010*. Ed. by R. Koch, A. Kolb, and C. Rezk-Salama. Eurographics Association, 2010, pp. 243–250 (cit. on pp. 85, 87).

[140] P. Milgram and F. Kishino. "A taxonomy of mixed reality visual displays". In: *IEICE TRANSACTIONS on Information and Systems* 77.12 (1994), pp. 1321–1329 (cit. on pp. 11, 12).

[141] M. Milinski and T. C. M. Bakker. "Female sticklebacks use male coloration in mate choice and hence avoid parasitized males". In: *Nature* 344.6264 (Mar. 1990), pp. 330–333 (cit. on pp. 102, 104).

[142] P. Mistry, P. Maes, and L. Chang. "WUW - wear Ur world: a wearable gestural interface". en. In: *Proceedings of the 27th international conference extended abstracts on Human factors in computing systems - CHI EA '09*. Boston, MA, USA: ACM Press, 2009, p. 4111 (cit. on p. 22).

[143] L. Miyashita, T. Yamazaki, K. Uehara, Y. Watanabe, and M. Ishikawa. "Portable Lumipen: Dynamic SAR in Your Hand". In: *2018 IEEE International Conference on Multimedia and Expo (ICME)*. July 2018, pp. 1–6 (cit. on pp. 48, 120).

[144] B. Mölzer and M. Strobelt. "Dimensional Management in Vehicle Development". en. In: (), p. 6 (cit. on p. 59).

[145] K. Müller, I. Smielik, J.-M. Hütwohl, S. Gierszewski, K. Witte, and K.-D. Kuhnert. "The virtual lover: variable and easily guided 3D fish animations as an innovative tool in mate-choice experiments with sailfin mollies-I. Design and implementation". en. In: *Current Zoology* 63.1 (Feb. 2017), pp. 55–64 (cit. on p. 107).

[146] R. Mur-Artal and J. D. Tardós. "Fast relocalisation and loop closing in keyframe-based SLAM". In: *2014 IEEE International Conference on Robotics and Automation (ICRA)*. 2014, pp. 846–853 (cit. on p. 47).

[147] M. Nagy, G. Vásárhelyi, B. Pettit, I. Roberts-Mariani, T. Vicsek, and D. Biro. "Context-dependent hierarchies in pigeons". In: *Proceedings of the National Academy of Sciences* 110.32 (2013), pp. 13049–13054. eprint: `https://www.pnas.org/content/110/32/13049.full.pdf` (cit. on pp. 124, 126).

[148] H. Naik, M. Bahaa, F. Tombari, P. Keitler, and N. Navab. "Frustration Free Pose Computation For Spatial AR Devices in Industrial Scenario". en. In: *2016 IEEE International Symposium on Mixed and Augmented Reality (ISMAR-Adjunct)*. Merida, Yucatan, Mexico: IEEE, Sept. 2016, pp. 121–122 (cit. on pp. 72, 75, 80–82, 171).

[149] H. Naik, R. Bastien, N. Navab, and I. D. Couzin. "Animals in Virtual Environments". en. In: *IEEE Transactions on Visualization and Computer Graphics* 26.5 (May 2020), pp. 2073–2083 (cit. on pp. 3, 4, 8, 11, 25, 26, 30, 34, 99, 100, 124, 157, 171).

[150] H. Naik, Y. Oyamada, P. Keitler, and N. Navab. "[POSTER] Exploiting Photogrammetric Targets for Industrial AR". en. In: *2015 IEEE International Symposium on Mixed and Augmented Reality*. Fukuoka, Japan: IEEE, Sept. 2015, pp. 144–147 (cit. on pp. 65–67, 71).

[151] H. Naik, F. Tombari, C. Resch, P. Keitler, and N. Navab. "[POSTER] A Step Closer To Reality: Closed Loop Dynamic Registration Correction in SAR". en. In: *2015 IEEE International Symposium on Mixed and Augmented Reality*. Fukuoka, Japan: IEEE, Sept. 2015, pp. 112–115 (cit. on pp. 63, 72, 75, 92, 93, 171).

[152] G. Narita, Y. Watanabe, and M. Ishikawa. "Dynamic Projection Mapping onto Deforming Non-Rigid Surface Using Deformable Dot Cluster Marker". In: *IEEE Transactions on Visualization and Computer Graphics* 23.3 (Mar. 2017), pp. 1235–1248 (cit. on pp. 48, 120).

[153] N. Navab. "Developing killer apps for industrial augmented reality". en. In: *IEEE Computer Graphics and Applications* 24.3 (May 2004), pp. 16–20 (cit. on pp. 4, 17, 55–57).

[154] N. Navab, S. Zokai, Y. Genc, and E. M. Coelho. "An on-line evaluation system for optical see-through augmented reality". In: *IEEE Virtual Reality 2004*. 2004, pp. 245–246 (cit. on p. 151).

[155] Z. W. Neil GLOSSOP, C. WEDLAKE, J. MOORE, and T. PETERS. "Augmented reality laser projection device for surgery". In: *Medicine Meets Virtual Reality 12: Building a Better You: the Next Tools for Medical Education, Diagnosis, and Care* 98 (2004). 00004, p. 104 (cit. on p. 64).

[156] X. J. Nelson and N. Fijn. "The use of visual media as a tool for investigating animal behaviour". en. In: *Animal Behaviour* 85.3 (Mar. 2013), pp. 525–536 (cit. on pp. 102, 105).

[157] A. Newell, K. Yang, and J. Deng. "Stacked Hourglass Networks for Human Pose Estimation". en. In: *arXiv:1603.06937 [cs]* (Mar. 2016). arXiv: 1603.06937 (cit. on pp. 44, 47, 149).

[158] A. K. T. Ng, L. K. Y. Chan, and H. Y. K. Lau. "A Study of Cybersickness and Sensory Conflict Theory Using a Motion-Coupled Virtual Reality System". In: *2018 IEEE Conference on Virtual Reality and 3D User Interfaces (VR)*. 2018, pp. 643–644 (cit. on pp. 25, 30).

[159] R. Nith and J. Rekimoto. "Falconer: A Tethered Aerial Companion for Enhancing Personal Space". In: *2019 IEEE Conference on Virtual Reality and 3D User Interfaces (VR)*. 2019, pp. 1550–1553 (cit. on p. 20).

[160] V. Nityananda, G. Tarawneh, S. Henriksen, D. Umeton, A. Simmons, and J. C. Read. "A Novel Form of Stereo Vision in the Praying Mantis". en. In: *Current Biology* (Feb. 2018) (cit. on p. 115).

[161] V. Nityananda, G. Tarawneh, R. Rosner, J. Nicolas, S. Crichton, and J. Read. "Insect stereopsis demonstrated using a 3D insect cinema". en. In: *Scientific Reports* 6.1 (May 2016) (cit. on pp. 114, 116).

[162] D. Novotny, N. Ravi, B. Graham, N. Neverova, and A. Vedaldi. "C3DPO: Canonical 3D Pose Networks for Non-Rigid Structure From Motion". In: *Proceedings of the IEEE International Conference on Computer Vision*. 2019 (cit. on pp. 120, 143–145).

[163] D. Novotny, N. Ravi, B. G. N. Neverova, and A. Vedaldi. "C3DPO: Canonical 3D Pose Networks for Non-Rigid Structure From Motion". en. In: (), p. 13 (cit. on p. 142).

[164] F. Okura, M. Kanbara, and N. Yokoya. "Interactive exploration of augmented aerial scenes with free-viewpoint image generation from pre-rendered images". In: *2013 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*. 2013, pp. 279–280 (cit. on p. 30).

[165] T. J. Ord, R. A. Peters, C. S. Evans, and A. J. Taylor. "Digital video playback and visual communication in lizards". en. In: *Animal Behaviour* 63.5 (May 2002), pp. 879–890 (cit. on pp. 102, 106, 108).

[166] C. Papachristos and K. Alexis. "Augmented reality-enhanced structural inspection using aerial robots". In: *2016 IEEE International Symposium on Intelligent Control (ISIC)*. 2016, pp. 1–6 (cit. on p. 20).

[167] D. Pavllo, C. Feichtenhofer, D. Grangier, and M. Auli. "3D human pose estimation in video with temporal convolutions and semi-supervised training". en. In: *arXiv:1811.11742 [cs]* (Mar. 2019). arXiv: 1811.11742 (cit. on pp. 144, 149).

[168] T. Peckmezian and P. W. Taylor. "A virtual reality paradigm for the study of visually mediated behaviour and cognition in spiders". en. In: *Animal Behaviour* 107 (Sept. 2015), pp. 87–95 (cit. on p. 112).

[169] A. Pérez-Escudero, J. Vicente-Page, R. C. Hinz, S. Arganda, and G. G. de Polavieja. "idTracker: tracking individuals in a group by automatic identification of unmarked animals". en. In: *Nature Methods* 11.7 (July 2014), pp. 743–748 (cit. on pp. 117, 142).

[170] T. J. Pitcher and J. E. T. Lawrence. "A simple stereo television system with application to the measurement of three-dimensional coordinates of fish in schools". In: *Behavior Research Methods, Instruments, & Computers* 16.6 (1984), pp. 495–501 (cit. on p. 106).

[171] Plato. *Plato's The Repulic*. New York:Books Inc., 1943 (cit. on pp. 7, 32).

[172] P. Pons, J. Jaen, and A. Catala. "Detecting Animals' Body Postures Using Depth-Based Tracking Systems". en. In: (2016), p. 5 (cit. on p. 142).

[173] R. Raskar, G. Welch, and H. Fuchs. "Spatially Augmented Reality". en. In: (), p. 7 (cit. on pp. 13, 21, 23).

[174] R. Raskar, R. Ziegler, and T. Willwacher. "Cartoon Dioramas in Motion". en. In: (2002), p. 8 (cit. on pp. 21, 22).

[175] H. Regenbrecht, G. Baratoff, and W. Wilke. "Augmented reality projects in the automotive and aerospace industries". In: *IEEE Computer Graphics and Applications* 25.6 (2005), pp. 48–56 (cit. on pp. 13, 17).

[176] G. Reinhart, W. Vogl, and I. Kresse. "A projection-based user interface for industrial robots". In: *Virtual Environments, Human-Computer Interfaces and Measurement Systems, 2007. VECIMS 2007. IEEE Symposium on*. IEEE, 2007, pp. 67–71 (cit. on p. 64).

[177] C. Resch, P. Keitler, and G. Klinker. "Sticky projections: A new approach to interactive shader lamp tracking". en. In: *2014 IEEE International Symposium on Mixed and Augmented Reality (ISMAR)*. Munich, Germany: IEEE, Sept. 2014, pp. 151–156 (cit. on pp. 22, 44, 45, 49, 87, 89).

[178] C. Resch, H. Naik, P. Keitler, S. Benkhardt, and G. Klinker. "On-Site Semi-Automatic Calibration and Registration of a Projector-Camera System Using Arbitrary Objects with Known Geometry". en. In: *IEEE Transactions on Visualization and Computer Graphics* 21.11 (Nov. 2015), pp. 1211–1220 (cit. on pp. 50, 171).

[179] S. Rohwer. "Dyed birds achieve higher social status than controls in Harris' sparrows". In: *Animal Behaviour* 33.4 (1985), pp. 1325–1331 (cit. on pp. 102, 104).

[180] J. P. Rolland, R. L. Holloway, and H. Fuchs. "Comparison of optical and video see-through, head-mounted displays". en. In: ed. by H. Das. Boston, MA, Dec. 1995, pp. 293–307 (cit. on pp. 16, 18).

[181] R. B. Rusu. "Semantic 3D Object Maps for Everyday Manipulation in Human Living Environments". In: *KI - Künstliche Intelligenz* 24.4 (2010), pp. 345–348 (cit. on p. 81).

[182] P. Scarfe and A. Glennerster. "The Science Behind Virtual Reality Displays". In: *Annual Review of Vision Science* 5.1 (2019). PMID: 31283449, pp. 529–547. eprint: `https://doi.org/10.1146/annurev-vision-091718-014942` (cit. on p. 30).

[183] D. Schmalstieg and D. Wagner. "Experiences with Handheld Augmented Reality". In: *Proceedings of the 2007 6th IEEE and ACM International Symposium on Mixed and Augmented Reality*. ISMAR '07. USA: IEEE Computer Society, 2007, pp. 1–13 (cit. on p. 21).

[184] W. Schreiber, K. Zürl, and P. Zimmermann, eds. *Web-basierte Anwendungen Virtueller Techniken: Das ARVIDA-Projekt – Dienste-basierte Software-Architektur und Anwendungsszenarien für die Industrie*. de. Berlin, Heidelberg: Springer Berlin Heidelberg, 2017 (cit. on pp. 13, 55, 67, 75, 76, 96, 151).

[185] S. Schuster, R. Strauss, and K. G. Götz. "Virtual-reality techniques resolve the visual cues used by fruit flies to evaluate object distances". In: *Current Biology* 12.18 (2002), pp. 1591–1594 (cit. on pp. 101, 102, 105, 108, 110).

[186] B. Schwerdtfeger. "Pick-by-vision: bringing HMD-based augmented reality into the warehouse". PhD thesis. Technical University Munich, 2010 (cit. on pp. 24, 64).

[187] B. Schwerdtfeger, A. Hofhauser, and G. Klinker. "An augmented reality laser projector using marker-less tracking". In: *Demonstration at 15th ACM Symposium on Virtual Reality Software and Technology (VRST'08)*. 00005. 2008 (cit. on pp. 64, 88).

[188] B. Schwerdtfeger, D. Pustka, A. Hofhauser, and G. Klinker. "Using laser projectors for augmented reality". In: *Proceedings of the 2008 ACM symposium on Virtual reality software and technology*. 00026. ACM, 2008, pp. 134–137 (cit. on pp. 23, 64, 83, 85).

[189] J. Shotton, A. Fitzgibbon, M. Cook, et al. "Real-time human pose recognition in parts from single depth images". In: *CVPR 2011*. 2011, pp. 1297–1304 (cit. on pp. 45, 47, 48, 129).

[190] L. Sigal, A. O. Balan, and M. J. Black. "HumanEva: Synchronized Video and Motion Capture Dataset and Baseline Algorithm for Evaluation of Articulated Human Motion". en. In: *International Journal of Computer Vision* 87.1-2 (Mar. 2010), pp. 4–27 (cit. on pp. 142, 144).

[191] E. Sikström, A. de Götzen, and S. Serafin. "Wings and flying in immersive VR — Controller type, sound effects and experienced ownership and agency". In: *2015 IEEE Virtual Reality (VR)*. 2015, pp. 281–282 (cit. on pp. 29, 30).

[192] E. Sobel. "The locust's use of motion parallax to measure distance". en. In: *Journal of Comparative Physiology A* 167.5 (Nov. 1990) (cit. on pp. 102, 105).

[193] C. Stapleton and J. Davies. "Imagination: The third reality to the virtuality continuum". en. In: *2011 IEEE International Symposium on Mixed and Augmented Reality - Arts, Media, and Humanities*. Basel, Switzerland: IEEE, Oct. 2011, pp. 53–60 (cit. on pp. 11–13).

[194] J. R. Stowers, A. Fuhrmann, M. Hofbauer, et al. "Reverse engineering animal vision with virtual reality and genetics". In: *Computer* 47.7 (2014), pp. 38–45 (cit. on pp. 100–102, 105, 108, 111–115, 117).

[195] J. R. Stowers, M. Hofbauer, R. Bastien, et al. "Virtual reality for freely moving animals". In: *Nature Methods* 14.10 (Aug. 2017), pp. 995–1002 (cit. on pp. 4, 34, 99, 101, 108, 113, 114, 116, 117, 121, 124).

[196] A. Strandburg-Peshkin, D. R. Farine, I. D. Couzin, and M. C. Crofoot. "Shared decision-making drives collective movement in wild baboons". In: *Science* 348.6241 (2015), pp. 1358–1361. eprint: https://science.sciencemag.org/content/348/6241/1358.full.pdf (cit. on pp. 124, 126).

[197] R. Strauss, S. Schuster, and K. G. Götz. "Processing of artificial visual feedback in the walking fruit fly Drosophila melanogaster." In: *Journal of Experimental Biology* 200.9 (1997), pp. 1281–1296 (cit. on p. 105).

[198] D. J. Sumpter, J. Krause, R. James, I. D. Couzin, and A. J. Ward. "Consensus Decision Making by Fish". In: *Current Biology* 18.22 (2008), pp. 1773–1777 (cit. on pp. 124, 126).

[199] I. E. Sutherland. "A head-mounted three dimensional display". en. In: *Proceedings of the December 9-11, 1968, fall joint computer conference, part I on - AFIPS '68 (Fall, part I)*. San Francisco, California: ACM Press, 1968, p. 757 (cit. on pp. 11, 15, 25).

[200] J. Takalo, A. Piironen, A. Honkanen, et al. "A fast and flexible panoramic virtual reality system for behavioural and electrophysiological experiments". en. In: *Scientific Reports* 2.1 (Dec. 2012) (cit. on p. 112).

[201] M. J. Tarr and W. H. Warren. "Virtual reality in behavioral neuroscience and beyond". en. In: *Nature Neuroscience* 5.S11 (Nov. 2002), pp. 1089–1092 (cit. on p. 101).

[202] G. K. Taylor, M. Bacic, R. J. Bomphrey, et al. "New experimental approaches to the biology of flight control systems". en. In: *Journal of Experimental Biology* 211.2 (Jan. 2008), pp. 258–266 (cit. on pp. 8, 100, 102, 105).

[203] A. Tero, S. Takagi, T. Saigusa, et al. "Rules for Biologically Inspired Adaptive Network Design". en. In: *Science* 327.5964 (Jan. 2010), pp. 439–442 (cit. on p. 120).

[204] J. C. Theobald, D. L. Ringach, and M. A. Frye. "Dynamics of optomotor responses in Drosophila to perturbations in optic flow". en. In: *Journal of Experimental Biology* 213.8 (Apr. 2010), pp. 1366–1375 (cit. on pp. 102, 105).

[205] K. Thurley and A. Ayaz. "Virtual reality systems for rodents". en. In: *Current Zoology* 63.1 (Feb. 2017), pp. 109–119 (cit. on pp. 4, 10, 99, 102, 108–111, 116).

[206] N. Tinbergen. "On Aims and Methods of Ethology". en. In: *Zeitschrift fuer Tierpsychologie* (1963), p. 28 (cit. on pp. 100, 103).

[207] N. Tinbergen and A. C. Perdeck. "On the Stimulus Situation Releasing the Begging Response in the Newly Hatched Herring Gull Chick (Larus Argentatus Argentatus Pont.)" en. In: *Behaviour* 3.1 (1950), pp. 1–39 (cit. on pp. 101, 102, 104).

[208] B. Triggs, P. F. McLauchlan, R. I. Hartley, and A. W. Fitzgibbon. "Bundle Adjustment — A Modern Synthesis". en. In: *Vision Algorithms: Theory and Practice*. Ed. by G. Goos, J. Hartmanis, J. van Leeuwen, B. Triggs, A. Zisserman, and R. Szeliski. Vol. 1883. Series Title: Lecture Notes in Computer Science. Berlin, Heidelberg: Springer Berlin Heidelberg, 2000, pp. 298–372 (cit. on p. 51).

[209] R. Y. Tsai. "A Versatile Camera Calibration Techniaue for High-Accuracy 3D Machine Vision Metrology Using Off-the-shelf TV Cameras and Lenses". en. In: (), p. 22 (cit. on p. 50).

[210] W. Tsai, L. Su, T. Ko, C. Yang, and M. Hu. "Improve the Decision-making Skill of Basketball Players by an Action-aware VR Training System". In: *2019 IEEE Conference on Virtual Reality and 3D User Interfaces (VR)*. 2019, pp. 1193–1194 (cit. on p. 28).

[211] T. Veen, S. J. Ingley, R. Cui, et al. "anyFish: an open-source software to generate animated fish models for behavioural studies". In: *Evolutionary Ecology Research* 15.3 (2013), pp. 361–375 (cit. on p. 107).

[212] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. "The Caltech-UCSD Birds-200-2011 Dataset". en. In: (), p. 8 (cit. on pp. 143, 145).

[213] B. Webb. "What does robotics offer animal behaviour?" en. In: *Animal Behaviour* 60.5 (Nov. 2000), pp. 545–558 (cit. on pp. 100, 104, 120).

[214] D. Weng, D. Cheng, Y. Wang, and Y. Liu. "Display systems and registration methods for augmented reality applications". en. In: *Optik - International Journal for Light and Electron Optics* 123.9 (May 2012), pp. 769–774 (cit. on p. 19).

[215] J. Werfel, K. Petersen, and R. Nagpal. "Designing Collective Behavior in a Termite-Inspired Robot Construction Team". en. In: *Science* 343.6172 (Feb. 2014), pp. 754–758 (cit. on p. 120).

[216] K. L. Woo and G. Rieucau. "From dummies to animations: a review of computer-animated stimuli used in animal behavior studies". en. In: *Behavioral Ecology and Sociobiology* 65.9 (Sept. 2011), pp. 1671–1685 (cit. on pp. 102, 107, 108, 114, 116).

[217] M. Zaeh and W. Vogl. "Interactive laser-projection for programming industrial robots". In: *Proceedings of the 5th IEEE and ACM International Symposium on Mixed and Augmented Reality*. IEEE Computer Society, 2006, pp. 125–128 (cit. on p. 64).

[218] F. Zheng, D. Schmalstieg, and G. Welch. "Pixel-wise closed-loop registration in video-based augmented reality". In: *Mixed and Augmented Reality (ISMAR), 2014 IEEE International Symposium on*. 00001. IEEE, 2014, pp. 135–143 (cit. on p. 87).

[219] F. Zheng, R. Schubert, and G. Weich. "A general approach for closed-loop registration in AR". In: *Virtual Reality (VR), 2013 IEEE*. 00003. IEEE, 2013, pp. 47–50 (cit. on p. 86).

[220] J. Zhou, I. Lee, B. H. Thomas, A. Sansome, and R. Menassa. *Facilitating Collaboration with Laser Projector-Based Spatial Augmented Reality in Industrial Applications*. 00000. Springer, 2011 (cit. on p. 85).

[221] J. Zhou, I. Lee, B. Thomas, R. Menassa, A. Farrant, and A. Sansome. "In-Situ Support for Automotive Manufacturing Using Spatial Augmented Reality". In: *International Journal of Virtual Reality* 11.1 (2012). 00004 (cit. on pp. 49, 85).

[222] S. Zollmann and O. Bimber. "Imperceptible calibration for radiometric compensation". In: *Proceedings Eurographics 2007, Short Paper* (2007). 00024 (cit. on p. 23).

[223] Y. Zou, W. Zhang, and Z. Zhang. "Liftoff of an Electromagnetically Driven Insect-Inspired Flapping-Wing Robot". In: *IEEE Transactions on Robotics* 32.5 (Oct. 2016), pp. 1285–1289 (cit. on p. 119).

[224] S. Zuffi, A. Kanazawa, T. Berger-Wolf, and M. J. Black. "Three-D Safari: Learning to Estimate Zebra Pose, Shape, and Texture from Images "In the Wild"". en. In: (), p. 10 (cit. on pp. 143, 153).

[225] S. Zuffi, A. Kanazawa, and M. J. Black. "Lions and Tigers and Bears: Capturing Non-rigid, 3D, Articulated Shape from Images". en. In: *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. Salt Lake City, UT: IEEE, June 2018, pp. 3955–3963 (cit. on pp. 120, 128, 142).

[226] S. Zuffi, A. Kanazawa, D. W. Jacobs, and M. J. Black. "3D Menagerie: Modeling the 3D Shape and Pose of Animals". en. In: IEEE, July 2017, pp. 5524–5532 (cit. on pp. 47, 128, 143).

# List of Figures

# List of Tables