# AI POWERED ATTENTION MONITORING SYSTEM FOR BOOK READING

## A Major Project Report
Submitted To



## Chhattisgarh Swami Vivekanand Technical University
## Bhilai, India

For

Major Project

of

**Bachelor of Technology (Hons.)**

*in*

**Computer Science & Engineering**

*By*

| Ashish Sinha | Jayant Patel | Himanshu Sahu |
|:---:|:---:|:---:|
| 300012721048 | 300012721061 | 300012721018 |
| CB4633 | CB4646 | CB4596 |
| 8th Sem | 8th Sem | 8th Sem |
| Artificial Intelligence | Artificial Intelligence | Artificial Intelligence |

Under the Guidance of
**Dr. Toran Verma**
Associate Professor
Department of Computer Science & Engineering
**UTD, CSVTU, Bhilai (C.G.)**



**Department of Computer Science & Engineering**
**University Teaching Department**
**Chhattisgarh Swami Vivekanand Technical University**
**Bhilai (C.G.) 491107**

**Session: 2024 – 2025**

**Department of Computer Science & Engineering**
**University Teaching Department**
**Chhattisgarh Swami Vivekanand Technical University**
**Bhilai (C.G.) 491107**

# DECLARATION BY THE CANDIDATE

We the undersigned solemnly declare that the major project report entitled *"AI POWERED ATTEN-TION MONITORING SYSTEM FOR BOOK READING"* is based our own work carried out during the course of our study under the supervision of ***Dr. Toran Verma***.

We assert that the statements made and conclusions drawn are an outcome of the project work. We further declare that to the best of our knowledge and belief that the report does not contain any part of any work which has been submitted for the award of any other degree/diploma/certificate in this University/Deemed university of India or any other country.

<div align="right">

**Ashish Sinha**
Roll No: 300012721048
Enroll No: CB4633
Semester: 8<sup>th</sup> (CSE)

</div>

Semester: $8^{th}$ (CSE)

<div align="right">

**Jayant Patel**
Roll No: 300012721061
Enroll No: CB4646
Semester: $8^{th}$ (CSE)

**Himanshu Sahu**
Roll No: 300012721018
Enroll No: CB4596
Semester: $8^{th}$ (CSE)

</div>

# CERTIFICATE BY THE SUPERVISOR

This is to certify that the major project report entitled *"AI POWERED ATTENTION MONITORING SYSTEM FOR BOOK READING"* is a record of project work carried out under my guidance and supervision for the fulfillment of the award of degree of Bachelor of Technology (Hons.) in the faculty of Computer Science & Engineering of Chhattisgarh Swami Vivekananda Technical University, Bhilai (C.G.) India.

To the best of my knowledge and belief the report

i. Embodies the work of the candidate himself.

ii. Has duly been completed.

iii. Fulfills the partial requirement of the ordinance relating to the B.Tech. (Hons) degree of the University.

iv. Is up to the desired standard both in respect of contents and language for being referred to the examiners.

———————————————
**Dr. Toran Verma**
Associate Professor
Dept. of Computer Science &
Engineering, UTD, CSVTU, Bhilai
(C.G.)

**Forwared to**
**Chhattisgarh Swami Vivekanand Technical University, Bhilai (C.G.)**

———————————————
**Dr. J P Patra**
HOD
Dept. of Computer Science &
Engineering. UTD, CSVTU, Bhilai
(C.G.)

———————————————
**Dr. Pankaj Mishra**
Director
UTD, CSVTU, Bhilai (C.G.)

# CERTIFICATE BY THE EXAMINERS

The project report entitled *"AI POWERED ATTENTION MONITORING SYSTEM FOR BOOK READING"* has been examined by the undersigned as a part of the examination of Bachelor of Technology (Hons.) in the faculty of Computer Science & Engineering of Chhattisgarh Swami Vivekanand Technical University, Bhilai.

**Internal Examiner**　　　　　　　　　　　　　　　　　　**External Examiner**
**Date:**　　　　　　　　　　　　　　　　　　　　　　　**Date:**

# ACKNOWLEDGEMENT

Working for this project has been a great experience for us. There were moments of anxiety, when we could not solve a problem for the several days. But we have enjoyed every bit of process and are thankful to all people associated with us during this period we convey our sincere thanks to our project guide **Dr. Toran Verma** for providing me all sorts of facilities. His support and guidance helped us to carry out the project. We owe a great dept. of his gratitude for his constant advice, support, cooperation & encouragement throughout the project we would also like to express our deep gratitude to respected **Dr. J P Patra** (Head of Department) for his ever helping and support. We also pay special thanks for his helpful solution and comments enriched by his experience, which improved our ideas for betterment of the project. We would also like to express our deep gratitude to respected **Dr. Pankaj Mishra** (Director) and college management for providing an educational ambience. It will be our pleasure to acknowledge, utmost cooperation and valuable suggestions from time to time given by our staff members of our department, to whom we owe our entire computer knowledge and also we would like to thank all those persons who have directly or indirectly helped us by providing books and computer peripherals and other necessary amenities which helped us in the development of this project which would otherwise have not been possible.

**Ashish Sinha**
Roll No: 300012721048
Enroll No: CB4633
Semester: 8$^{th}$ (CSE)

**Jayant Patel**
Roll No: 300012721061
Enroll No: CB4646
Semester: 8$^{th}$ (CSE)

**Himanshu Sahu**
Roll No: 300012721018
Enroll No: CB4596
Semester: 8$^{th}$ (CSE)

# TABLE OF CONTENTS

# LIST OF ABBREVIATIONS

| Abbreviation | Description |
| --- | --- |
| AI | Artificial Intelligence |
| ML | Machine Learning |
| LLM | Large Language Model |
| DDPM | Denoising Diffusion Probabilistic Model |
| VAE | Variational Autoencoder |
| GAN | Generative Adversarial Network |
| MLP | Multi-Layer Perceptron |
| RMSE | Root Mean Squared Error |
| MSE | Mean Squared Error |
| KPI | Key Performance Indicator |
| GPU | Graphics Processing Unit |
| CPU | Central Processing Unit |
| API | Application Programming Interface |

# LIST OF FIGURES

**Department of Computer Science & Engineering**
**University Teaching Department**
**Chhattisgarh Swami Vivekanand Technical University**
**Bhilai (C.G.) 491107**

# LIST OF TABLES

**Department of Computer Science & Engineering**
**University Teaching Department**
**Chhattisgarh Swami Vivekanand Technical University**
**Bhilai (C.G.) 491107**

# ABSTRACT

Lack of attention while reading is a critical challenge in education, often leading to reduced comprehension and learning outcomes. To address this, we propose a novel real-time system that monitors reader attention by combining object detection and gaze estimation. Our method uses YOLO for accurate book detection and integrates AWS Rekognition to assess gaze direction, determining whether the reader is focused on the book. This lightweight, cloud-assisted approach enables efficient, scalable deployment. The system has strong potential in personalized education, enabling adaptive learning platforms that respond to student engagement levels.

**Keywords:** Reader attention, YOLO, AWS Rekognition, gaze detection, smart learning, real-time monitoring.

# Chapter 1
# Introduction

## 1.1  Background Information

Reading proficiency stands as a cornerstone of academic achievement, continuous professional growth, and individual enlightenment. The capacity to effectively engage with written text underpins learning across all disciplines and is integral to navigating the complexities of modern life. However, the act of sustained reading demands a significant cognitive resource: focused attention. In an increasingly digital world, characterized by a deluge of information and pervasive stimuli, maintaining unwavering concentration during reading tasks presents a substantial challenge for many individuals [8]. When a reader's attention deviates, the cognitive processes responsible for comprehension, critical analysis, and memory encoding are disrupted. This can lead to reduced reading speed, superficial understanding of the material, and ultimately, a diminished return on the time invested in reading. The challenge is not merely one of willpower but is often exacerbated by environmental factors and the inherent cognitive load associated with processing complex textual information [9].

Historically, strategies to enhance reading focus have centered on behavioral modifications, such as cultivating disciplined study habits, optimizing the reading environment to minimize external distractions, and employing time management techniques like the Pomodoro method. While these approaches offer value, their efficacy can vary greatly among individuals and may prove insufficient in scenarios demanding prolonged, deep engagement with challenging texts. Recognizing these limitations, there is a burgeoning interest in exploring how technological advancements can offer more direct, personalized, and real-time support for attention management during reading.

The confluence of Artificial Intelligence (AI), particularly its sub-fields of computer vision and machine learning, has paved the way for innovative solutions to monitor and support cognitive states, including attention [21]. Computer vision endows systems with the ability to perceive and interpret visual information from an environment, akin to human sight. By leveraging standard cameras (like webcams) and sophisticated algorithms, it is now feasible to non-intrusively track various visual cues indicative of a user's engagement level. Among these, eye gaze tracking has emerged as a powerful tool, providing insights into the locus of a user's visual attention, which is often closely coupled with their cognitive focus [5, 15]. Techniques such as those proposed by Abdelrahman et al. in L2CS-Net demonstrate the capability for fine-grained gaze estimation even in unconstrained, real-world

environments [1].

Beyond gaze direction, other visual cues like head posture, facial expressions, and even contextual information about the objects of interaction (such as a book) contribute to a more holistic understanding of user attention. For instance, identifying the presence and state (open or closed) of a book within the reader's environment provides critical context. Object detection models, particularly efficient real-time systems like You Only Look Once (YOLO) [19] and its subsequent iterations [13, 24, 23], can be trained to accurately locate and classify such objects. The integration of gaze data with the location of a target object, such as an open book, allows for a more direct inference of whether the user is attending to the reading material [26, 16].

The overarching goal of developing AI-driven attention monitoring systems for reading is to empower users by providing them with actionable feedback and a deeper awareness of their attention patterns. Such systems can identify moments of waning focus or distraction, potentially alerting the user or logging these instances for later review, thereby fostering metacognitive skills and promoting more effective reading habits [3]. This project is situated within this evolving landscape, aiming to design, implement, and evaluate a system that monitors a user's attention while they are engaged in reading a physical book, leveraging gaze estimation and book detection to achieve this. The insights gained from such a system could have implications for educational technologies, personal productivity tools, and a better understanding of reader engagement.

## 1.2  Project Objectives

The primary aim of this project is to design, develop, and evaluate an AI-powered system for monitoring a user's attention during book reading sessions. The system leverages computer vision techniques to analyze visual cues from a standard webcam. The specific objectives formulated to guide the development of this project are as follows:

- **To enhance reading focus and comprehension:** The foremost objective is to develop a system that assists users in improving their concentration while reading. This is to be achieved by providing real-time or near real-time feedback regarding their attention levels directed towards the reading material.

- **To track and analyze attention patterns:** The project aims to implement functionalities for monitoring and recording the user's attention span over the course of reading sessions. This data can then be used to identify individual attention trends, common distraction points, and potential areas for improvement in reading habits.

- **To promote better reading habits:** By making users more aware of their attention fluctuations and the moments when their focus shifts away from the book, the system intends to encourage the development of more consistent and effective reading practices.

- **To develop a robust gaze estimation module:** A core technical objective is to integrate and utilize an accurate gaze estimation model (specifically, the L2CS model) capable of determining the pitch and yaw of the user's gaze from webcam input, serving as a primary indicator of visual attention.

- **To implement an effective book detection module:** The project seeks to develop a reliable object detection component using a custom-trained YOLO model (YOLOv12s) to identify the presence and state (e.g., "open_book", "closed_book") of a physical book within the camera's field of view.

- **To integrate gaze and book information for attention assessment:** A crucial objective is to fuse the data from the gaze estimation and book detection modules to intelligently infer whether the user's visual attention is directed towards an open book.

- **To design and implement a user interface for interaction and feedback:** The system aims to provide a user-friendly interface that allows users to initiate and manage monitoring sessions, view real-time feedback, and potentially access logs of their attention patterns. (Although the attention analyzer part was commented out, this was an initial objective).

These objectives collectively address the challenge of maintaining focus during reading by creating a technological aid that is both informative and aims to be assistive for the user.

## 1.3   Significance of the Project

The development of an AI-powered Book Reading Attention Monitoring system holds considerable significance across multiple domains, primarily by addressing the pervasive challenge of maintaining focused attention during reading and by exploring innovative applications of computer vision technologies. The importance of this project can be understood through its potential contributions to individual users, the field of human-computer interaction (HCI), and educational technology.

### 1.3.1   Enhancement of Individual Learning and Productivity

For individuals, particularly students and lifelong learners, the ability to maintain concentration is directly linked to comprehension and knowledge retention [3]. This project signifies a step towards

providing a personalized tool that can help users become more aware of their attention patterns. By offering insights into when and potentially why their focus wanes, the system can empower users to develop more effective reading strategies, leading to improved learning outcomes and increased productivity during study or research sessions. The real-time feedback mechanism, as envisioned, could act as a gentle nudge, helping users to consciously redirect their attention back to their reading material.

### 1.3.2 Advancement in Non-Intrusive Attention Assessment

Traditionally, attention assessment has often relied on subjective self-reports, obtrusive physiological sensors, or controlled laboratory settings. This project contributes to the growing field of non-intrusive attention monitoring by leveraging readily available hardware like webcams. The application of computer vision techniques, specifically gaze estimation [1] and object detection [19], for understanding reader engagement with physical books in a naturalistic setting is a significant area of exploration. It moves beyond screen-based attention tracking to address the common scenario of reading traditional print media.

### 1.3.3 Contribution to Human-Computer Interaction (HCI) and Affective Computing

Understanding and responding to a user's cognitive and affective state is a key goal in HCI and affective computing. This project's focus on detecting attention, a critical cognitive state, contributes to this domain. While the current scope is on attention, the underlying technologies for gaze and visual context understanding could be extended to infer other states like engagement, confusion, or fatigue, which are highly relevant for creating more adaptive and responsive interactive systems [21]. The development of systems like a "Reading Companion" [16] underscores the interest in such interactive aids.

### 1.3.4 Potential for Educational Technology and Personalized Learning

While this project focuses on individual use, its principles have significant implications for educational technology. Data on attention patterns could, with appropriate ethical considerations and privacy safeguards [11], inform the design of personalized learning interventions or provide educators with insights into student engagement during independent study. The ability to objectively gather data on how individuals interact with physical texts when trying to learn could be valuable for educational research.

4

### 1.3.5   Exploration of Real-World AI Applications:

The project serves as a practical application of advanced AI models in a real-world scenario. Successfully integrating gaze estimation (L2CS) and a custom-trained object detector (YOLO) to solve a specific problem (monitoring attention on a book) demonstrates the feasibility and utility of these technologies beyond theoretical research. It highlights the potential to build accessible tools that can run on user-end devices, promoting the democratization of AI-powered assistance.

In summary, the significance of this project lies in its potential to provide a practical tool for improving reading focus, its contribution to the methods of non-intrusively monitoring attention, and its broader implications for creating more intelligent and adaptive systems in educational and personal productivity contexts. It addresses a common, everyday challenge with a modern technological approach.

## 1.4   Scope and Limitations

### 1.4.1   Scope of the Project

The scope of this project is the design, development, and evaluation of a software application that utilizes computer vision and deep learning techniques to monitor a user's visual attention while they are reading a physical book. The system processes input from a standard webcam in real-time. Key aspects within the project's scope include:

- **User and Environment:** The system is intended for a single user engaged in reading a physical book, typically in a common indoor setting such as at a desk.

- **Core Technologies Implemented:**

    - Integration and application of the L2CS model for real-time estimation of the user's gaze direction, yielding pitch and yaw angles.

    - Deployment of a custom-trained YOLOv12s object detection model for the identification of physical books and the classification of their state (e.g., "open_book," "closed_book") within the camera's view.

- **Primary Functionality:**

    - Continuous detection of the user's face and real-time estimation of their gaze vector.

    - Robust detection of physical books and determination of whether they are open or closed.

- – Real-time, frame-by-frame assessment of the user's attention towards an open book. This is achieved by:

  1. Calculating a 3D gaze vector from the estimated pitch and yaw.
  2. Determining if this gaze vector intersects with the detected bounding box of an open book.

- – Visual feedback to the user via an on-screen display, showing the webcam feed augmented with information such as detected face/book bounding boxes, gaze direction indicators, and the current attention status (e.g., "Attentive," "Distracted," "No book detected").

- **Attention Data Generation:** The system generates per-frame data detailing the gaze pitch and yaw, face bounding box, book bounding box, book state, and the inferred attention status. This foundational data can be used for:

  - – Providing immediate visual feedback on attention.

  - – The basis for potential future extensions such as session-based attention percentage calculation and the generation of detailed attention logs for pattern analysis, as outlined in the initial project requirements .

- **Platform and Development:** The system is developed as a desktop application using Python, with core libraries including OpenCV, PyTorch, and Ultralytics for computer vision and deep learning tasks.

The project thus demonstrates the technical pipeline for real-time visual attention assessment in the context of reading physical books by integrating advanced gaze and object detection models.

### 1.4.2 Limitations of the Project

Despite the system's capabilities in per-frame attention assessment, several limitations should be acknowledged, stemming from the inherent complexities of AI models, environmental variables, and the defined scope of the project:

- **Accuracy and Robustness of AI Models:** The overall system performance is intrinsically tied to the accuracy of the L2CS gaze estimation and the YOLOv12s book detection models. Their effectiveness can be compromised by:

  - – **Environmental Factors:** Suboptimal lighting (dim, glare, strong backlighting), shadows, and visually complex or cluttered backgrounds may degrade detection and estimation accuracy.

6

- **Occlusions:** Partial obstruction of the user's eyes or face (e.g., by hands, certain types of eyewear) can affect gaze estimation. Similarly, if the book is significantly occluded, its detection and state classification may falter.

- **Gaze Estimation Specifics:** L2CS provides gaze direction (pitch/yaw). While indicative of focus, it does not provide a pixel-precise point-of-gaze on the book page. The model's accuracy can also be user-dependent and influenced by camera distance and angle.

- **Book Detector Generalizability:** The custom-trained YOLO model's ability to detect various books accurately depends on the breadth and diversity of its training dataset. Performance might vary with unusual book sizes, covers, or in novel environments.

- **Definition and Inference of Attention:** The system defines and infers attention based on a measurable visual cue: gaze direction towards a detected open book. This is a proxy for attention and does not:

  - Directly measure cognitive engagement or reading comprehension. A user could be looking at an open book but be mentally disengaged.

  - Account for nuanced reading behaviors, such as brief glances away for reflection, which might be misclassified as inattention by a purely gaze-book intersection logic.

- **Session-Level Metrics and Alerts:** While the system provides real-time, per-frame attention status, the implementation of more advanced session-level features as envisioned in the project requirements such as triggering alerts based on sustained inattention (e.g., after a specific duration), calculating overall session attention percentages, or generating detailed structured logs for long-term pattern analysis—are foundational. The current core logic focuses on instantaneous attention state per frame. Extending this to robust session-long analytics would require additional state management and temporal analysis logic.

- **User Movement and Positioning:** The system performs optimally when the user is relatively stable and positioned appropriately within the camera's view. Frequent or significant head and body movements might reduce the accuracy of continuous face and gaze tracking.

- **Computational Load:** Real-time execution of multiple deep learning models (gaze estimation and object detection) on each video frame is computationally demanding. System responsiveness (e.g., FPS, latency of feedback) can be heavily influenced by the available CPU/GPU resources on the user's machine.

- **Focus on Physical Books:** The methodology is specifically tailored for monitoring attention while reading physical books. It is not designed or evaluated for reading on digital devices like e-readers, tablets, or computer screens, which would require different contextual cues or interaction models.

- **Types of Distractions Addressed:** The system primarily identifies visual inattention (gaze directed away from the book). It cannot inherently detect or account for other forms of distraction, such as auditory disturbances or internal cognitive distractions (e.g., mind-wandering), if the user's gaze remains on the book.

Understanding these limitations is essential for interpreting the system's current output and for guiding future development efforts to enhance its robustness and expand its feature set.

## 1.5   Overview of the Thesis Structure

This thesis is organized into several chapters, each addressing specific aspects of the research and development of the AI-powered Book Reading Attention Monitoring system. The structure is designed to guide the reader from the foundational concepts and motivations through the technical implementation, results, and concluding reflections.

- **Chapter 1: Introduction** lays the groundwork for the thesis. It begins with background information on the importance of reading attention and the potential of AI in this domain. This is followed by a clear statement of the project objectives, the significance of the research, and a definition of the project's scope and its inherent limitations. The chapter concludes with this overview of the thesis structure.

- **Chapter 2: Literature Review** (added section) delves into existing academic and technical literature relevant to the core components of this project. It covers established and recent advancements in gaze estimation techniques, real-time object detection methodologies (with a focus on YOLO and its applicability for book detection), and various approaches to attention monitoring and reader engagement analysis using computer vision.

- **Chapter 3: Methodology** outlines the research design and approach adopted for this project. It provides an overview of the project's architecture, details the methods used for data collection (such as the creation of the custom book dataset), and explains the techniques employed for data analysis, including the logic for gaze-book intersection. Ethical considerations pertaining to the project and any limitations of the chosen methodology are also discussed.

- **Chapter 4: Implementation** provides a comprehensive account of the system's development. This includes a description of the development environment, tools, and libraries used. It details the execution of the project, including the implementation of the gaze estimation module (L2CS), the book detection module (YOLOv12s), and their integration into the attention monitoring application. Any project timeline, resource allocation, challenges faced during implementation, success factors, and key lessons learned are also presented.

- **Chapter 5: Results and Discussion** presents the outcomes of the project. This chapter will focus on the performance of the developed modules, such as the accuracy of the book detector and the functionality of the gaze estimator. It will include an interpretation of these results and a discussion of the key findings in the context of the project objectives. Limitations encountered during testing and potential future directions for improving the results are also explored.

- **Chapter 6: Learning Outcome** summarizes the overall findings of the project and assesses the achievement of the initial objectives. It discusses the implications of the project's outcomes, offers recommendations, and explores the future scope for extending the system's capabilities. This chapter also includes personal reflections on the learning journey throughout the project.

- **Chapter 7: Conclusion and Broader Impact** (Interpreting the PDF's "Conclusion and Future Scope" which seems focused on development) provides a concluding perspective, focusing on the skills developed, knowledge gained during the project, aspects of professional development, personal growth experienced, and potential future applications of the developed system or the acquired expertise.

- **References** lists all the academic papers, articles, books, and other resources cited throughout the thesis, providing the basis for the literature review and supporting technical discussions.

- **Appendices** (if applicable) may include supplementary materials such as detailed model configurations, code snippets, or extensive data tables not suitable for the main body of the thesis.

This structure is intended to provide a clear and logical progression of the research undertaken.

# Chapter 2

# Literature Review

## 2.1 Introduction

This chapter reviews existing literature pertinent to the development of an AI-powered book reading attention monitoring system. The review is structured around three core technological pillars: gaze estimation techniques, object detection methodologies (with a focus on identifying books), and approaches to attention monitoring, particularly in contexts related to reading and learning. By examining prior work in these areas, this section aims to establish the foundation upon which the current project is built, highlight relevant advancements, and identify the existing landscape of research.

## 2.2 Gaze Estimation Techniques

Gaze estimation, the process of determining the direction or point of a user's gaze, is a critical component for understanding visual attention. Research in this area has evolved significantly, broadly categorized into model-based and appearance-based methods. Model-based methods typically involve geometric modeling of the eye, while appearance-based methods directly learn a mapping from eye or face images to gaze direction [5].

Recent advancements have heavily favored appearance-based approaches leveraging deep learning due to their robustness in unconstrained environments without requiring specialized hardware [15]. Convolutional Neural Networks (CNNs) have been extensively used to extract salient features from facial and eye regions for gaze regression. For instance, the L2CS-Net model, utilized in this project, employs a CNN architecture designed for fine-grained gaze estimation by regressing pitch and yaw angles independently, demonstrating strong performance on challenging datasets like Gaze360 [1]. Other recent research explores lightweight architectures for efficient gaze estimation [17] and the application of transformers to this domain [12], pushing the boundaries of accuracy and efficiency. The ability of such models to function with standard webcams makes them highly suitable for accessible attention monitoring applications.

## 2.3    Object Detection for Contextual Awareness

To understand if a user is attending to a specific item, such as a book, the system must first be able to detect and localize that item within the visual scene. Object detection has seen remarkable progress with the advent of deep learning, particularly with one-stage detectors like the You Only Look Once (YOLO) family. The original YOLO model introduced a paradigm shift by framing object detection as a single regression problem, enabling real-time performance [19].

Subsequent iterations, including YOLOv5 [13], YOLOv8 [23], and YOLOv9 [24], have progressively improved accuracy, speed, and versatility, making them state-of-the-art for various applications. These models are designed to identify multiple objects within an image, predict their bounding boxes, and classify them. Surveys on YOLO variants highlight their architectural evolution and wide adoption [22, 20]. For this project, the ability of YOLO models to be custom-trained for specific object classes, such as "open_book" and "closed_book," is crucial. Research has also demonstrated the application of deep learning, including YOLO-based approaches, for specific tasks like book detection and document analysis [2], underscoring the feasibility of this component. More recent developments like YOLO-NAS also aim to enhance object detection capabilities [6].

## 2.4    Attention Monitoring and Reader Engagement

Monitoring user attention, especially in educational and reading contexts, has garnered significant research interest. Attention is a complex cognitive state, and inferring it from observable behaviors is a challenging task. Computer vision offers non-intrusive methods to capture relevant cues.

Gaze direction is a primary indicator of overt visual attention [8]. Studies have explored using gaze patterns to assess student engagement in classrooms [21] and online learning environments [18, 10]. The concept of a "Reading Companion" system that monitors reading behaviors, including attention, has been explored to provide assistance to readers [16]. Furthermore, research indicates that gaze patterns can be correlated with cognitive load during reading [9] and even reading comprehension [4].

Head pose and facial expressions are also valuable cues often integrated into multimodal attention monitoring systems [25, 14]. A comprehensive survey by Albu et al. [3] discusses various deep learning techniques for engagement level estimation. The current project focuses on gaze and book presence, but the literature supports the value of these cues as strong indicators of attention towards reading material.

## 2.5  Integrating Gaze with Object Context for Focused Attention

While general gaze direction provides some information, understanding attention towards a *specific* object requires linking the gaze vector to the object's location. This is particularly relevant for determining if a reader is looking at their book. Research such as that by Zhu et al. [26] explores methods for localizing an object prompted by gaze, effectively performing a gaze-object mapping. This project's core attention inference mechanism, which involves calculating the intersection of the user's gaze vector with the bounding box of a detected open book, aligns with this principle of contextualizing gaze. By confirming that the user's visual focus aligns with the location of the reading material, a more reliable inference of attention during the reading task can be made.

## 2.6  Summary and Research Gap

The literature demonstrates robust advancements in individual areas of gaze estimation, real-time object detection, and vision-based attention/engagement monitoring. Models like L2CS-Net offer accurate gaze tracking, while the YOLO family provides efficient object detection. Numerous studies have explored attention in learning and reading, often using gaze as a key metric.

While many systems track attention on screens or use specialized eye-tracking hardware, this project aims to synthesize these advancements to address the specific scenario of monitoring attention on physical books using a standard webcam. It combines a state-of-the-art gaze estimator with a custom-trained, efficient object detector to infer attention in a non-intrusive manner. The key contribution lies in the practical integration of these components for the nuanced task of discerning if a reader's focus is indeed on their physical reading material, providing a foundation for a tool that could aid in improving reading habits. The ethical considerations of such monitoring systems are also an emerging area of discussion relevant to the responsible development of these technologies [11].

This review indicates that while the foundational technologies are well-established, their specific combination and application for real-time attention monitoring on physical books with commodity hardware represent a practical and relevant area of investigation.

# Chapter 3

# Methodology

## 3.1  Project Overview

This project, titled "Book Reading Attention Monitoring," is designed to ascertain a user's level of visual attention while they are engaged in reading a physical book. The system leverages computer vision and deep learning techniques to process a live video feed from a standard webcam, analyzing visual cues such as eye gaze direction and the presence and state of a book to make an inference about the user's focus. The ultimate aim is to provide a foundation for a tool that could offer real-time feedback to users, helping them understand and potentially improve their reading attention patterns.

The system operates through a sequence of interconnected modules, each responsible for a specific task in the attention monitoring pipeline. The overall architecture is depicted in Figure **??** and described subsequently. (I will provide the TikZ code for this figure after this text block).

The core workflow begins with capturing video input from the user's webcam. Each frame from this video stream is then processed by two primary analysis components:

1. **Gaze Estimation:** A dedicated module, utilizing the L2CS (Look-Likely-to-See) deep learning model [1], processes the video frame to detect the user's face and estimate the direction of their eye gaze. This results in pitch and yaw values representing the gaze orientation, along with the bounding box of the detected face.

2. **Book Detection and Attention Analysis:** Concurrently, the video frame, along with the derived gaze information, is fed into an attention monitoring module. This module internally employs a custom-trained YOLOv12s (You Only Look Once) object detection model [19] to locate any books within the frame and determine their state (i.e., "open_book" or "closed_book").

The attention monitoring module then integrates these two streams of information. It calculates a 3D gaze vector based on the estimated pitch and yaw, originating from the center of the detected face. A key part of its logic involves performing a ray-box intersection test: it determines if this calculated 3D gaze vector intersects with the 3D bounding volume conceptually formed by a detected "open_book."

13

Based on the outcome of this intersection test, along with the presence of a detected face and an open book, the system makes a frame-by-frame assessment of whether the user is attentive to the reading material. This attention status (e.g., "Attentive," "Distracted," "Book not found," "Face not found") is then made available.

The entire process is orchestrated by a session manager component, which handles frame acquisition, coordinates the analysis modules, and manages the display of the processed video feed to the user. The visual output includes the original camera feed augmented with overlays indicating detected faces, book bounding boxes, gaze direction cues, and the current inferred attention status. This provides immediate, albeit qualitative, feedback to the user.

The following sections in this chapter will delve deeper into the specific research design, data collection methods (particularly for the custom book dataset), the detailed data analysis techniques including the gaze-book intersection logic, and any ethical considerations and methodological limitations.



Figure 3.1: System Architecture Overview of the Book Reading Attention Monitor (Corrected Flow).

## 3.2  Research Design and Approach

The development of the Book Reading Attention Monitoring system was guided by an applied research and iterative development methodology. The primary goal was to engineer a functional prototype that integrates existing advanced AI models for a novel application: monitoring reader attention on physical books using standard webcam hardware. The approach can be characterized by several key phases, as illustrated in Figure 3.2 and detailed below.

```
┌─────────────────────┐        ┌─────────────────────┐
│  1. Conceptualization│───────▶│  2. System Design   │
│  & Literature Review │        │  & Model Selection  │
└─────────────────────┘        └─────────────────────┘
                                           │
                                           ▼
┌─────────────────────┐        ┌─────────────────────┐
│  4. Integration &   │◀───────│  3. Core Compo-     │
│  System Assembly    │        │  nent Development   │
└─────────────────────┘        └─────────────────────┘
         │
         ▼
┌─────────────────────┐
│  5. Testing & Iter- │
│  ative Refinement   │
└─────────────────────┘

┌─────────────────────┐  ┌─────────────────────┐  ┌─────────────────────┐
│  Gaze Estimation    │  │  Book Detection Mod-│  │  Attention Anal-    │
│  Module (L2CS)      │  │  ule (Custom YOLO)  │  │  ysis Module        │
└─────────────────────┘  └─────────────────────┘  └─────────────────────┘
```

Figure 3.2: Iterative Research and Development Approach.

1. **Conceptualization and Requirements Definition:** The initial phase involved defining the problem: the challenge of maintaining attention during reading and the potential for an AI-based solution. This included outlining the primary objectives (enhancing focus, tracking patterns, promoting better habits), identifying key functionalities (gaze tracking, book detection, attention inference), and defining the overall scope of the system targeting users reading physical books.

2. **Literature Review and Technology Assessment:** A thorough review of existing literature was conducted (as detailed in Chapter 2) to understand the state-of-the-art in relevant areas:

   - *Gaze Estimation:* Exploring various techniques, with a focus on appearance-based deep learning models suitable for webcam input, leading to the selection of the L2CS model [1] for its reported accuracy in unconstrained environments.

   - *Object Detection:* Investigating real-time object detection frameworks, particularly the

YOLO family [19, 13, 23], for their efficiency and the ability to be custom-trained. This informed the decision to use a YOLO-based architecture (YOLOv12s) for book detection.

- *Attention Monitoring:* Reviewing approaches for inferring attention from visual cues, including studies on reader engagement and multimodal systems.

This phase was crucial for selecting appropriate technologies and understanding existing methodologies.

3. **System Design and Architecture Planning:** Based on the requirements and technology assessment, a modular system architecture was designed (as presented in Section 3.1 Project Overview). This involved defining the primary components: Camera Manager, Gaze Estimator, Book Detector, Attention Monitor, and Session Manager, and outlining their interactions and data flow. The design prioritized local processing to ensure user privacy and real-time (or near real-time) performance.

4. **Core Component Development and Implementation:** This phase focused on the practical implementation of the key modules:

- *Gaze Estimation Module:* Integration of the pre-trained L2CS model to process video frames and output pitch/yaw gaze angles and face bounding boxes.

- *Book Detection Module:* This involved a significant sub-task of creating a custom dataset of images featuring open and closed books in various settings. This dataset was then used to train a YOLOv12s model to specifically detect "open_book" and "closed_book" classes.

- *Attention Analysis Module:* Development of the logic to fuse information from the gaze and book detection modules. The core of this module is the algorithm for calculating the 3D gaze vector and implementing the ray-box intersection test to determine if the user's gaze is directed at an open book.

- *Supporting Modules:* Implementation of the 'CameraManager' for handling video input and the 'SessionManager' for orchestrating the overall application flow and user display.

5. **Integration and System Assembly:** The individual modules were then integrated into a cohesive application. This involved managing data pipelines between components, ensuring synchronization, and handling potential errors or edge cases (e.g., face not detected, book not found).

6. **Testing and Iterative Refinement:** The system underwent functional testing throughout its development. This included:

16

- Testing the accuracy and robustness of the gaze estimator under various conditions.

- Evaluating the performance of the custom-trained book detector on a validation set and in real-world scenarios.

- Verifying the logic of the attention analysis module, particularly the gaze-book intersection.

- Observing the real-time performance and responsiveness of the integrated system.

Although a formal quantitative evaluation with human subjects was beyond the immediate scope of this phase, the testing was iterative, with observations and identified issues feeding back into the development and refinement of the components. The design approach was flexible to accommodate adjustments based on these tests.

The overall approach was thus empirical and constructive, focused on building a working system by applying and integrating established AI techniques to a specific real-world problem. The emphasis was on demonstrating the feasibility of the proposed solution and creating a platform that could be further developed and rigorously evaluated in future work.

## 3.3 Data Collection Methods

The successful development of the AI components within this project, particularly the custom book detection model, necessitated a systematic approach to data collection. While the gaze estimation module utilizes a pre-trained L2CS model [1], which was developed and validated on existing large-scale gaze datasets, the book detection module required the creation of a specialized dataset tailored to the project's needs.

### 3.3.1 Gaze Estimation Data

For the L2CS gaze estimation model, this project leveraged the publicly available pre-trained weights. These weights were derived from training on extensive datasets containing a wide variety of face images with corresponding ground-truth gaze direction annotations (e.g., pitch and yaw). The use of such a pre-trained model significantly reduced the development effort that would otherwise be required for collecting and annotating a large-scale gaze dataset from scratch, allowing the project to focus on the integration and application of this technology. The L2CS model's training on diverse datasets like Gaze360 contributes to its robustness in unconstrained environments.

### 3.3.2 Book Detection Dataset (Custom YOLOv12s Model)

A critical aspect of this project was the development of a custom object detector capable of accurately identifying physical books and their state (open or closed). This required the collection and annotation of a specific image dataset for training the YOLOv12s model. The data collection process for this purpose involved two primary strategies:

1. **Leveraging Publicly Available Datasets via Roboflow Universe:** To establish a foundational dataset, an initial set of images containing books was sourced from Roboflow Universe [7]. Roboflow Universe is a public repository of computer vision datasets, offering a diverse range of images and annotations across various domains. Relevant datasets featuring books in different contexts were explored, and suitable images were selected to bootstrap the training process. This approach provided an immediate collection of varied images, saving considerable initial data gathering time.

2. **Self-Captured Supplementary Data Collection:** Recognizing that publicly available datasets might not fully cover the specific nuances and variations anticipated in the project's target usage scenarios (e.g., different book types, lighting conditions in typical reading environments, various angles, open/closed states), the initial dataset was augmented with self-captured images. This supplementary data collection was performed manually and aimed to:

   - **Increase Diversity:** Images were captured of various types of books (hardcover, paperback, different sizes, varied cover designs), in different states (fully open, partially open, closed), and from multiple perspectives and distances relative to the camera.

   - **Simulate Real-World Conditions:** Efforts were made to capture images under a range of lighting conditions (natural light, artificial indoor light, mixed lighting) and against different backgrounds typically found in study or reading environments.

   - **Address Potential Edge Cases:** Specific scenarios, such as books held in hands, lying flat on a table, or propped up, were included to improve the model's robustness.

   These self-captured images were taken using standard smartphone cameras and webcams to reflect the input quality expected by the final system.

### 3.3.3 Data Annotation and Preparation

All images in the custom book detection dataset, whether sourced from Roboflow Universe or self-captured, underwent a meticulous annotation process.

- **Annotation Classes:** Two primary classes were defined for annotation: "open_book" and "closed_book." These classes are essential for the system to distinguish whether a detected book is in a state relevant for active reading.

- **Bounding Box Annotation:** For each instance of a book in an image, a tight bounding box was drawn to precisely delineate its location. This is standard practice for training YOLO-based object detectors.

- **Annotation Tools:** The annotation process have utilized tools such as Roboflow's integrated annotation platform. This tools facilitate the accurate drawing of bounding boxes and assignment of class labels.

- **Dataset Split:** Following standard machine learning practice, the final aggregated dataset was partitioned into training, validation, and testing sets. This division is crucial for training the model, tuning its hyperparameters, and evaluating its performance on unseen data to ensure generalization. (You would ideally specify the approximate number of images and the split ratio here, e.g., "The final dataset comprised approximately XXXX images, split into 70% for training, 15% for validation, and 15% for testing.")

- **Data Augmentation:** To further enhance the dataset's size and variability, and to improve the model's robustness against minor changes in input, data augmentation techniques were likely applied. Common augmentations for object detection tasks include geometric transformations (e.g., rotation, scaling, flipping) and photometric distortions (e.g., changes in brightness, contrast, saturation). Platforms like Roboflow often provide built-in augmentation capabilities that can be applied during the dataset versioning process.

The careful collection, annotation, and preparation of this custom book detection dataset were vital steps in developing a YOLOv12s model capable of reliably identifying books and their states, which is a cornerstone of the attention monitoring logic.

## 3.4   Directory Structure

The project is organized into a structured hierarchy of directories and files to ensure modularity, maintainability, and ease of navigation. The main components of the project are located within a root folder, conventionally named `book-attention-monitor/`. A representation of the primary directory structure is as follows:

```
book-attention-monitor/
|-- src/
|   |-- analysis/
|   |   |-- __init__.py
|   |   |-- attention_analyzer.py
|   |   `-- attention_monitor.py
|   |-- camera/
|   |   |-- __init__.py
|   |   `-- camera_manager.py
|   |-- models/
|   |   |-- __init__.py
|   |   |-- book_detector.py
|   |   `-- gaze_estimator.py
|   |-- model_weights/
|   |   |-- L2CSNet_gaze360.pkl
|   |   `-- last.pt (or your specific YOLO weights file)
|   |-- session/
|   |   |-- __init__.py
|   |   `-- session_manager.py
|   |-- utils/
|   |   |-- __init__.py
|   |   `-- helpers.py
|   `-- __init__.py
|
|-- report/
|   |-- index.md
|   |-- 01_application_entry_point_.md
|   |-- ... (other markdown report files)
|
|-- main.py
|-- requirements.txt
|-- attention_monitor.log
`-- README.md
```

A brief description of each key directory and file is provided below:

**book-attention-monitor/** The root directory of the project.

**src/** This directory contains all the core source code for the application, organized into sub-modules:

- `analysis/`: Contains modules responsible for the higher-level analysis of visual cues.

    - `attention_monitor.py`: Implements the core logic for assessing user attention by integrating gaze and book detection data. It includes the gaze-book intersection algorithm.

    - `attention_analyzer.py`: (As per your codebase structure, though its specific role might be integrated or complementary to `attention_monitor.py`). May contain additional or alternative attention analysis logic.

- `camera/`: Houses the module for managing camera input.

    - `camera_manager.py`: Handles webcam initialization, frame capture, and display functionalities.

- `models/`: Contains the Python classes that wrap and utilize the deep learning models.

    - `gaze_estimator.py`: Implements the interface for the L2CS gaze estimation model, processing frames to extract gaze direction (pitch and yaw) and face bounding boxes.

    - `book_detector.py`: Implements the interface for the custom-trained YOLOv12s model, responsible for detecting books and their states (open/closed) in the video frames. (Note: Your `Readme.md` mentioned `object_detector.py`, but the uploaded code file is `book_detector.py`).

- `model_weights/`: This crucial directory (implied by `main.py`) stores the pre-trained model weight files required by the system.

    - `L2CSNet_gaze360.pkl`: The weights for the L2CS gaze estimation model.

    - `last.pt` (or similar `.pt`, `.onnx` file): The weights for your custom-trained YOLOv12s book detection model.

- `session/`: Manages the overall application session and workflow.

    - `session_manager.py`: Orchestrates the interactions between the camera, gaze estimation, book detection, and attention analysis modules. It handles the main processing loop and user interface updates.

- `utils/`: Contains utility functions and helper scripts used across different modules.

21

- – `helpers.py`: Provides common utility functions, potentially for tasks like frame resizing, logging setup, or other supporting operations.

- • `__init__.py`: Present in `src/` and its subdirectories, these files indicate to Python that the directories should be treated as packages, allowing for modular imports.

**report/** (Present in your uploaded files) This directory likely contains documentation and report files related to the project, such as the markdown files detailing different components of the system.

**main.py** The main entry point for the application. This script initializes all necessary components (like the camera manager, gaze estimator, and session manager) and starts the attention monitoring session.

**requirements.txt** Lists all the Python package dependencies required to run the project (e.g., OpenCV, PyTorch, Ultralytics, NumPy). Users can install these dependencies using pip.

**attention_monitor.log** The log file where runtime information, warnings, and errors are recorded, as configured in `main.py`.

**README.md** Provides an overview of the project, features, setup instructions, and other relevant information for users and developers.

This structured organization facilitates code management, debugging, and potential future scalability of the project.

## 3.5   Data Analysis Techniques

The process of determining a user's attention towards a physical book in this project involves a sequence of data analysis techniques applied to the incoming video stream. These techniques transform raw pixel data into meaningful intermediate representations (gaze direction, book location) and finally culminate in an attention status inference. The primary stages of data analysis are detailed below.

### 3.5.1   Frame Acquisition and Preprocessing

Each video frame captured by the `CameraManager` serves as the initial input. While extensive preprocessing is not a primary focus, frames are handled in a format (NumPy arrays) suitable for the subsequent deep learning models. Any necessary resizing or color space conversions (e.g., BGR to

RGB, if required by specific models) are implicitly handled by the model inference pipelines or utility functions.

### 3.5.2 Gaze Vector Estimation

Once a frame is acquired, the `GazeEstimator` module, which encapsulates the L2CS model [1], performs the following analysis:

1. **Face Detection:** The L2CS pipeline first detects the user's face within the frame. This step is crucial as gaze estimation is typically performed relative to facial features. The output includes the bounding box coordinates of the detected face.

2. **Gaze Angle Regression:** Using the detected facial region, the L2CS model regresses the gaze angles, specifically pitch ($\theta_p$) and yaw ($\theta_y$). These angles represent the upward/downward and leftward/rightward orientation of the gaze, respectively.

3. **3D Gaze Vector Calculation:** The estimated pitch and yaw angles are then converted into a 3D unit vector, $\vec{G} = [G_x, G_y, G_z]$, representing the gaze direction in a 3D coordinate system relative to the face. This conversion, implemented within the `AttentionMonitor`, typically uses the following spherical to Cartesian coordinate transformation (assuming pitch $\phi$ and yaw $\lambda$ in radians):

$$G_x = \cos(\phi)\sin(\lambda)$$
$$G_y = \sin(\phi) \qquad\qquad (3.1)$$
$$G_z = \cos(\phi)\cos(\lambda)$$

   This 3D vector is fundamental for subsequent geometric analysis. The origin of this vector is typically considered to be the center of the detected face bounding box.

The output of this stage is a set of gaze parameters: pitch, yaw, the calculated 3D gaze vector, and the face's bounding box.

### 3.5.3 Book Detection and State Classification

Simultaneously or sequentially, the `AttentionMonitor` module utilizes its internal, custom-trained YOLOv12s model [19] to analyze the video frame for the presence of books:

1. **Object Detection:** The YOLO model processes the input frame to identify instances of books.

2. **Bounding Box and Class Prediction:** For each detected book, the model outputs a bounding box (coordinates of the rectangle enclosing the book) and a class label. The classes relevant

23

to this project are "open_book" and "closed_book," along with a confidence score for each detection.

This stage provides crucial contextual information: whether a book is present, where it is located, and whether it is open (and thus likely being read).

### 3.5.4   Attention Inference via Gaze-Book Intersection

This is the core analytical step where data from gaze estimation and book detection are fused to infer the user's attention state. This logic resides within the `AttentionMonitor` module and follows a decision-making process, illustrated in Figure 3.3.

The key steps in this inference are:

1. **Prerequisite Checks:** The system first verifies if a face is detected. If not, attention cannot be assessed, and an appropriate status is set. Subsequently, it checks if an "open_book" is detected by the internal YOLO model. If no open book is found, the user cannot be attentive to it in the intended manner.

2. **Gaze Ray Formulation:** If a face and an open book are detected, the 3D gaze vector (calculated as described in Section 3.4.2) is used. The origin of this gaze ray is the center of the detected face bounding box.

3. **Book Volume Definition:** The 2D bounding box of the detected "open_book" is conceptually extended into a 3D volume. The `AttentionMonitor` code defines parameters `z_near` and `z_far` that establish a depth range for this volume, effectively treating the book as a cuboid in 3D space relative to the camera.

4. **Ray-Box Intersection Test:** The core analytical technique is a geometric ray-box intersection test. The algorithm checks if the calculated 3D gaze ray (originating from the user's face and pointing in the direction of their gaze) intersects with any of the six planes forming the conceptual 3D volume of the open book. This test is implemented in the `_ray_box_intersection` method of the `AttentionMonitor`. An intersection is considered valid only if it occurs within the defined spatial boundaries of the book's 2D projection and within the `z_near` and `z_far` depth limits.

5. **Attention Status Determination:**
   - If a face and an open book are detected, and the gaze ray intersects the book's 3D volume, the user is classified as "Attentive."
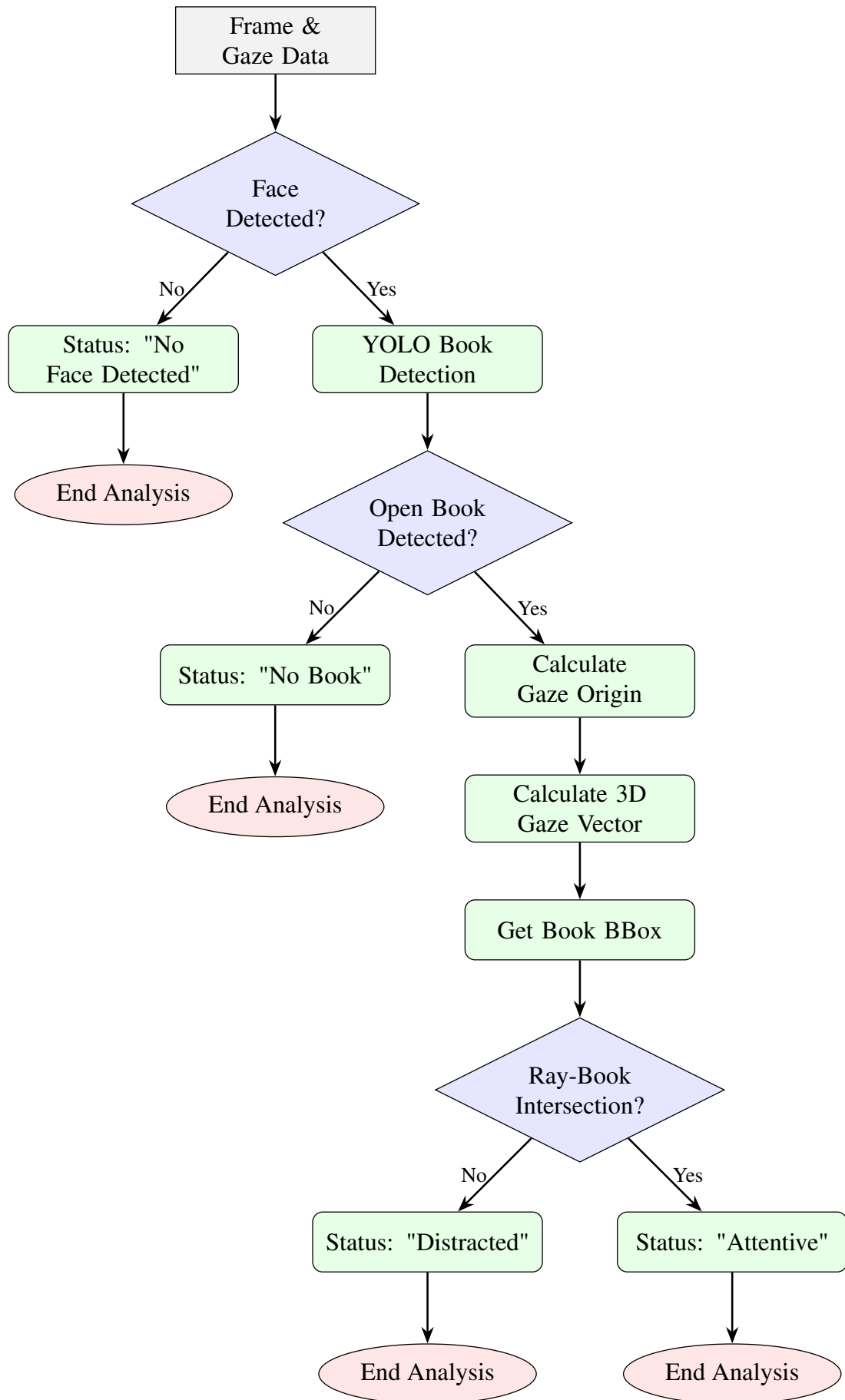
24

Figure 3.3: Flowchart of the Attention Inference Logic.

- If a face and an open book are detected, but the gaze ray does not intersect the book's volume, the user is classified as "Distracted" (implying they are looking away from the detected open book).

- Other statuses like "No face detected" or "No open book detected" are set if the prerequisites are not met.

### 3.5.5 Output Generation

The final output of the data analysis pipeline for each frame is a structured dictionary containing the inferred attention status (e.g., attentive, distracted), a descriptive message, and the relevant intermediate data such as gaze pitch/yaw, gaze vector components, and bounding boxes for the face and book. This structured output is then used by the `SessionManager` to provide visual feedback to the user by overlaying this information on the displayed video stream.

These data analysis techniques, combining deep learning model inferences with geometric reasoning, form the computational backbone of the attention monitoring system.

## 3.6  Ethical Considerations

The development and deployment of AI systems, particularly those involving human monitoring through computer vision, necessitate careful consideration of ethical implications. This project, while aimed at a beneficial purpose, acknowledges and has sought to address several key ethical aspects:

1. **User Privacy:** The foremost ethical concern is the privacy of the user, as the system processes video data from their webcam, which includes images of their face and immediate environment.

   - *Mitigation:* A core design principle of this project is local processing. All video frame analysis, including face detection, gaze estimation, and book detection, is performed on the user's local machine. No image or video data is transmitted to external servers or stored by the application by default. This significantly reduces the risk of unauthorized access to sensitive visual data.

2. **Informed Consent:** Users must be fully aware that they are being monitored and how the system operates.

   - *Mitigation:* The system is designed to be explicitly user-initiated. The user must actively start a monitoring session. Clear information regarding what the system monitors (gaze

direction relative to a book, presence of a book) and the nature of the feedback provided should be made available to the user before they opt to use the system.

3. **Data Security and Storage:** Even if data is processed locally, any temporary data or generated logs (if this feature were fully implemented for user review) must be handled securely.

   - *Mitigation:* The project's non-functional requirements emphasize data security. As no personal image data is persistently stored by default, the primary risk is minimized. If future versions implement session log storage for the user's own review, these logs should be stored locally and encrypted, or the user should be given explicit control over their storage and deletion.

4. **Potential for Misinterpretation and Misuse:** Attention is a complex cognitive state, and inferring it solely from visual cues can lead to misinterpretations. There's also a risk that such technology could be misused for undue surveillance or judgment if applied in contexts without proper safeguards (e.g., mandatory use in evaluative settings).

   - *Mitigation:* This project is designed as an assistive tool for self-monitoring and improvement. The feedback aims to be informative rather than punitive. It is crucial to emphasize that the system provides an estimate of visual attention to a book, not a definitive measure of cognitive engagement or learning. Any broader application would require careful ethical review and context-specific guidelines [11].

5. **Algorithmic Bias:** AI models can inherit biases from their training data, potentially leading to varied performance across different demographic groups (e.g., based on skin tone, facial features, or types of eyewear for gaze estimation).

   - *Mitigation:* While this project utilizes pre-trained models like L2CS and a custom-trained YOLO model, awareness of potential biases is important. The L2CS model was trained on diverse datasets to promote generalizability. For the custom book detector, efforts to diversify the self-captured training images were made. Continuous evaluation and retraining with more diverse data would be necessary for a production-level system to mitigate such biases further.

6. **Psychological Impact:** Continuous monitoring, even for self-improvement, could potentially induce anxiety or a feeling of being constantly evaluated in some users.

   - *Mitigation:* The system should be presented as a supportive tool, giving users control over when to use it. The nature of feedback should be constructive. Future iterations could

27

explore user-configurable feedback sensitivity or modes that are less intrusive.

7. **Transparency and Explainability:** While deep learning models are often "black boxes," providing users with a basic understanding of how the system arrives at its conclusions can foster trust.

   - *Mitigation:* The system's logic (detecting gaze towards an open book) is relatively straightforward at a high level. Explaining this principle to users—that it checks if they are looking at the book they are reading—can enhance transparency.

Adherence to these ethical principles is paramount for the responsible development and acceptance of attention monitoring technologies. This project strives to embed these considerations into its design philosophy, primarily through local processing and user-controlled operation.

## 3.7   Limitations of Methodology

The methodological choices made during the design and development of this Book Reading Attention Monitoring system, while enabling the creation of a functional prototype, also introduce certain limitations. These are distinct from the overall project limitations (discussed in Chapter 1) and pertain more to the constraints and characteristics of the chosen techniques and approaches:

1. **Reliance on Appearance-Based Gaze Estimation:** The choice of the L2CS model, an appearance-based gaze estimation technique, offers the advantage of working with standard webcams without specialized hardware. However, these methods are generally more susceptible to variations in lighting, head pose extremes, occlusions (e.g., glasses, hair), and individual facial morphology compared to intrusive, infra-red based eye trackers. The accuracy, while state-of-the-art for its category, may not achieve the precision of dedicated eye-tracking hardware, limiting the granularity of gaze information (e.g., pinpointing the exact word on a page).

2. **Custom Dataset Scope for Book Detection:** The performance of the custom-trained YOLOv12s book detector is fundamentally linked to the quality, size, and diversity of the dataset created for its training. While efforts were made to diversify this dataset by sourcing images from Roboflow Universe and supplementing with self-captured images, the dataset may still not encompass the vast variety of book types, cover designs, sizes, conditions (new, worn), and reading environments. This limitation can affect the model's generalization capability when encountering books or settings significantly different from those in its training set.

3. **Simplifications in Attention Inference Logic:** The core attention inference method relies on a geometric ray-box intersection between the estimated 3D gaze vector and the 3D bounding volume of a detected open book. This is a pragmatic approach but involves simplifications:

   - **Binary Attention State:** The current logic primarily results in a binary classification (attentive/distracted per frame). It does not capture varying degrees or intensities of attention.

   - **Fixed Depth Assumption for Book Volume:** The `z_near` and `z_far` parameters used to define the book's 3D volume are heuristic and may not accurately represent the book's actual depth or orientation in all cases, potentially affecting intersection accuracy.

   - **Single Gaze Point Origin:** Using the face center as the gaze origin is an approximation.

4. **Exclusion of Other Attention Cues:** The methodology primarily focuses on gaze direction and book presence as indicators of attention. Other potentially valuable visual cues (e.g., detailed facial expressions, subtle head movements indicative of engagement or confusion, body posture) or non-visual cues (e.g., physiological data, interaction patterns if it were an e-reader) are not incorporated in the current analysis. This limits the richness of the attention model.

5. **Lack of Formal User-Based Evaluation Protocol:** While functional testing was performed, the current methodology does not include a formal quantitative evaluation with human participants to validate the accuracy of the attention inference against ground truth (e.g., self-reported attention, task performance). Such studies would be necessary to rigorously assess the system's real-world effectiveness and user acceptance.

6. **Real-Time Processing Constraints on Model Complexity:** The objective of achieving real-time performance on local machines influences the choice of models. More complex, potentially more accurate, gaze estimation or object detection models might be too computationally expensive for the target deployment scenario. This practical constraint means a trade-off is often made between accuracy and processing speed.

These methodological limitations provide context for the current system's capabilities and highlight areas where future research and development could refine the techniques and enhance the system's overall performance and scope.

# Chapter 4

# Implementation

## 4.1 Development Environment

The development of the Book Reading Attention Monitoring system was carried out using a combination of specific hardware and software tools. This section outlines the key components of the development environment that facilitated the project's implementation.

### 4.1.1 Hardware

The primary hardware components utilized or assumed for the development and operation of the system include:

- **Computer System:** A standard desktop or laptop computer capable of running Python and the associated deep learning libraries. Specifications would typically include a multi-core CPU (e.g., Intel Core i5/i7 series or AMD Ryzen equivalent).

- **Graphics Processing Unit (GPU):** While not strictly mandatory for all components, a CUDA-enabled NVIDIA GPU was highly recommended and utilized for significantly accelerating the training of the custom YOLO model and for faster inference speeds of both the L2CS and YOLO deep learning models during real-time operation. Operations can fall back to CPU if a GPU is unavailable, albeit with a performance reduction.

- **Webcam:** A standard USB webcam or an integrated laptop camera was used as the primary input device for capturing the live video feed of the user. The system was designed to be compatible with common webcam resolutions (e.g., 720p, 1080p). IP cameras were also considered as a potential source, as handled by the `CameraManager`.

- **Memory (RAM):** A minimum of 8GB RAM was considered suitable, with 16GB or more being preferable, especially during model training and when running multiple processes.

- **Storage:** Sufficient disk space for storing the project codebase, Python environment, deep learning model weights, custom datasets, and log files.

### 4.1.2 Software

The software stack for this project is primarily based on Python and its rich ecosystem of libraries for computer vision and deep learning.

- **Operating System:** The system was developed to be cross-platform, with successful operation intended on Windows, macOS, or Linux distributions. Development was primarily conducted on [Specify your OS, e.g., Windows 10/11, Ubuntu 20.04 LTS, macOS Monterey].

- **Programming Language:** Python (version 3.8 or higher) was used as the sole programming language for the entire project implementation due to its extensive library support for AI/ML development, ease of use, and rapid prototyping capabilities.

- **Core Libraries and Frameworks:** The `requirements.txt` file specifies the key Python dependencies. The most critical libraries and frameworks employed are:

    - **OpenCV (cv2)** (`opencv-python >=4.7.0`): Utilized extensively for all computer vision tasks, including camera interaction, frame reading, image processing (resizing, drawing overlays), and displaying the video feed.

    - **PyTorch** (`torch >=2.0.0`, `torchvision >=0.15.0`): The primary deep learning framework used for both the L2CS gaze estimation model and the YOLOv12s object detection model. It was used for model loading, inference, and managing tensor operations, especially on the GPU.

    - **Ultralytics YOLO:** The Ultralytics framework was used for implementing and training the custom YOLOv12s model for book detection. This library provides a high-level API for working with various YOLO architectures. (Note: This is a key dependency inferred from the code, e.g., `from ultralytics import YOLO`).

    - **L2CS-Net Library (`l2cs`):** The official implementation or a compatible wrapper for the L2CS gaze estimation model was used to perform inference and obtain pitch/yaw gaze angles. (Note: This is a key dependency inferred from the code, e.g., `from l2cs import Pipeline`).

    - **NumPy** (`numpy >=1.22.0`): Essential for numerical operations, especially for handling image data as multi-dimensional arrays and for mathematical calculations involved in gaze vector processing.

    - **ONNX and ONNXRuntime** (`onnx >=1.13.0`, `onnxruntime >=1.13.0`): While the primary models (L2CS .pkl, YOLO .pt) might be loaded directly in their native for-

31

mats, these libraries are often included for model conversion, optimization, or broader deployment compatibility, suggesting they might have been explored or used in an intermediate step.

– **python-dotenv** (`python-dotenv >=0.21.0`): Used for managing environment variables, potentially for configurations like API keys or model paths if an `.env` file was used, though not extensively shown in the core logic.

- **Development Tools:**

  – **Integrated Development Environment (IDE):** [Specify your IDE, e.g., Visual Studio Code, PyCharm Community/Professional Edition] was used for code writing, debugging, and project management.

  – **Version Control System:** Git was used for version control, with the project hosted on a GitHub repository. This facilitated tracking changes, collaboration (if any), and codebase management.

  – **Annotation Tool(s):** For creating the custom book dataset, an image annotation tool such as [Specify tool if known, e.g., Roboflow's annotation interface, LabelImg, CVAT] was employed to draw bounding boxes and assign labels for "open_book" and "closed_book" classes.

This environment provided the necessary capabilities for developing, training (for the YOLO model), and testing the attention monitoring system.

## 4.2 Project Execution

The execution of the Book Reading Attention Monitoring project followed a structured, modular approach, progressing through several distinct stages from initial setup to the integration of the final system. The development process was iterative, allowing for refinements as each component was built and tested.

1. **Environment Setup and Initial Planning:** The first step involved setting up the development environment as detailed in Section 4.1. This included installing Python, all required libraries (OpenCV, PyTorch, Ultralytics, L2CS, etc.), and configuring the IDE and version control. Simultaneously, the project requirements were finalized, and a detailed plan for module development and integration was laid out, building upon the initial research and design phase.

2. **Development of Core Modules:** The system was developed in a modular fashion, with each key component implemented and tested, often in isolation, before integration:

- **Camera Management (`CameraManager`):** Implementation of the `CameraManager` class to handle robust access to various camera sources (webcam, video files, IP streams), manage frame capture, and provide basic display functionalities.

- **Gaze Estimation (`GazeEstimator`):** Integration of the pre-trained L2CS model. This involved writing the `GazeEstimator` class to load the model, process input frames, run inference to obtain pitch, yaw, and face bounding boxes, and structure the output data.

- **Book Detection (`BookDetector` and Custom YOLO Training):** This was a significant sub-project.

  - *Dataset Curation:* As described in Section 3.3 (Data Collection Methods), this involved sourcing initial images from Roboflow Universe and capturing a substantial number of custom images to create a diverse dataset for "open_book" and "closed_book" states.

  - *Annotation:* Meticulous annotation of this dataset with bounding boxes and class labels.

  - *YOLOv12s Model Training:* Utilizing the Ultralytics YOLO framework to train the custom YOLOv12s model on the prepared dataset. This involved selecting appropriate hyperparameters, training for a sufficient number of epochs, and evaluating the model's performance on a validation set.

  - *Implementation:* The `BookDetector` class (or internal YOLO usage within `AttentionMoni` was implemented to load the trained custom weights and perform inference on video frames to detect books and their states.

- **Attention Analysis Logic (`AttentionMonitor`):** Development of the `AttentionMonitor` class. This included:

  - Implementing the conversion of pitch/yaw angles to a 3D gaze vector.

  - Designing and coding the ray-book bounding box intersection algorithm (`_ray_box_interse` to determine if the user's gaze vector intersects with the detected open book.

  - Establishing the decision logic to classify attention status based on face detection, open book detection, and the intersection result.

- **Session Orchestration (`SessionManager`):** Creation of the `SessionManager` class to act as the central coordinator. This module is responsible for:

33

- – Initializing and managing all other modules.

- – Handling the main application loop: capturing frames, dispatching them to analysis modules (potentially using threading for smoother performance, as indicated by the use of `Queue` in the code).

- – Aggregating results from the gaze, book, and attention analysis modules.

- – Managing the real-time display of the processed video feed with informative overlays (gaze direction, book boxes, attention status).

- **Utility Functions (`helpers.py`):** Development of any helper functions needed across the project, such as logging setup or image manipulation routines.

3. **Integration and System Testing:** Once individual modules reached a stable state, they were integrated into the main application framework managed by the `SessionManager`. This phase involved:

   - Ensuring correct data flow and communication between modules.

   - Testing the end-to-end pipeline from camera input to attention status output and display.

   - Debugging issues arising from the interaction of different components.

   - Performing functional tests under various simulated reading scenarios to observe the system's behavior.

4. **Iterative Refinement and Documentation:** Throughout the execution, an iterative approach was adopted. Issues identified during testing led to refinements in the code, model parameters (for YOLO), or algorithmic logic. Simultaneously, documentation of the code and project structure was maintained. The `report/` directory suggests that component-wise documentation was also being prepared.

The project execution focused on building a robust proof-of-concept, prioritizing the successful implementation and integration of the core AI-driven attention monitoring pipeline.

## 4.3   Timeline

The development of the Book Reading Attention Monitoring system was executed over a period of [Specify Duration, e.g., approximately X months/weeks, or one academic semester], commencing from [Start Date/Month, Year] to [End Date/Month, Year]. The project was broken down into several key phases, each with an estimated timeframe. While minor overlaps and adjustments occurred, the planned timeline provided a roadmap for the project's progression.

**Phase 1: Conceptualization, Literature Review, and Requirement Analysis (Estimated: X weeks)**

- Defining project scope, objectives, and initial requirements.

- Conducting a comprehensive literature review on gaze estimation, object detection (YOLO), and attention monitoring techniques.

- Assessing available tools, libraries, and pre-trained models (L2CS, YOLO).

- Finalizing the system architecture and technological stack.

**Phase 2: Dataset Preparation and Model Training (Estimated: Y weeks)**

- Sourcing initial book images from public repositories (e.g., Roboflow Universe).

- Capturing and curating a custom dataset of book images (open/closed states, diverse conditions).

- Annotating the custom dataset with bounding boxes and class labels.

- Training the YOLOv12s model for book detection, including hyperparameter tuning and validation.

**Phase 3: Core Module Development (Estimated: Z weeks)**

- Implementing the `CameraManager` for video input.

- Integrating the L2CS model into the `GazeEstimator` module.

- Implementing the `BookDetector` using the trained custom YOLO model.

- Developing the core logic for the `AttentionMonitor`, including gaze-book intersection.

- Implementing the `SessionManager` for overall orchestration and basic UI.

**Phase 4: Integration, Testing, and Refinement (Estimated: A weeks)**

- Integrating all developed modules into a cohesive system.

- Conducting functional testing of the end-to-end pipeline.

- Debugging and resolving issues identified during integration and testing.

- Iteratively refining the model parameters, algorithmic logic, and user interface based on test results.

**Phase 5: Finalization and Documentation (Estimated: B weeks)**

35

- Documenting the codebase, system architecture, and functionalities.

- Preparing the project report/thesis, including literature review, methodology, implementation details, results, and conclusions.

- Preparing for final project demonstration and presentation.

## 4.4 Resource Allocation

The successful execution of this project relied on the allocation and utilization of various resources, categorized as follows:

- **Hardware Resources:**

  - **Development Computer:** A [Your computer type, e.g., personal laptop/desktop] with [CPU details, e.g., Intel Core i7-XXXX], [RAM amount, e.g., 16GB RAM].

  - **GPU:** An NVIDIA [Your GPU model, e.g., GeForce RTX 3060] with CUDA support was utilized for accelerating deep learning model training (YOLO) and inference tasks. If a GPU was not consistently available, CPU-based inference was used, albeit with lower performance.

  - **Webcam:** A [Specify webcam type, e.g., integrated laptop webcam / specific USB webcam model] was used for video input.

  - **Storage:** Local hard drive/SSD storage for the operating system, development tools, codebase, datasets (including the custom book image dataset which was approximately [Specify size if known, e.g., X GB]), and model weights.

- **Software Resources:**

  - **Operating System:** Windows 11

  - **Programming Environment:** Python-3.11, along with VS Code IDE.

  - **Core Libraries:** As detailed in Section 4.1 (Development Environment), key libraries included OpenCV, PyTorch, Ultralytics, L2CS, and NumPy.

  - **Pre-trained Models:**

    * L2CS model weights (`L2CSNet_gaze360.pkl`).

  - **Dataset Platforms and Annotation Tools:**

    * Roboflow Universe for sourcing initial book image datasets [7].

36

       * Roboflow's online tool for annotating the custom book dataset.

  – **Version Control:** Git and GitHub for codebase management.

• **Human Resources:**

  – **Developer Time:** The primary resource was the time and effort invested by the project developer(s) in research, design, coding, training, testing, and documentation. This amounted to approximately [Specify total hours or person-months if estimated] over the project duration.

  – **Guidance and Supervision:** Input and guidance from academic supervisors or mentors.

• **Data Resources:**

  – Publicly available image datasets from Roboflow Universe.

  – Self-captured images for the custom book dataset.

  – Gaze360 on which the pre-trained L2CS model was originally trained.

Effective management and utilization of these resources were crucial for achieving the project's objectives within the defined timeframe.

## 4.5 Challenges Faced

Several challenges were encountered during the development lifecycle of the Book Reading Attention Monitoring system. Overcoming these hurdles was integral to the project's progress:

• **Dataset Curation for Book Detection:** Creating a robust and diverse dataset for training the custom YOLOv12s book detector was a significant undertaking.

  – *Challenge:* Sourcing a sufficient quantity of varied images representing different book types, sizes, cover designs, open/closed states, and various real-world reading environments (lighting, backgrounds, occlusions) was time-consuming.

  – *Mitigation/Approach:* This was addressed by combining images from public repositories like Roboflow Universe with a dedicated effort to capture and annotate a substantial number of custom images. Data augmentation techniques were also planned/employed to increase dataset variability.

• **Accuracy and Robustness of AI Models in Real-World Conditions:**

– *Challenge (Gaze Estimation):* The L2CS gaze estimator, while powerful, could sometimes yield less accurate results under challenging lighting conditions, with certain types of eyewear, or with extreme head poses. Ensuring consistent performance across diverse users and environments was difficult.

– *Challenge (Book Detection):* The custom YOLO model's accuracy was sensitive to how well its training data represented the actual use-case scenarios. Occluded books, unusual book appearances, or cluttered backgrounds sometimes led to missed detections or misclassifications.

– *Mitigation/Approach:* Iterative testing in various conditions helped identify weaknesses. For book detection, continuous refinement of the training dataset and data augmentation were key strategies. For gaze, understanding the model's limitations and designing the system for typical reading postures helped.

• **Integration of Multiple AI Models and Real-Time Performance:**

– *Challenge:* Running both the L2CS gaze estimation and the YOLO object detection models simultaneously on a live video stream, while also executing the attention analysis logic, posed a computational challenge, especially on systems without high-end GPUs. Achieving a smooth frame rate and low latency for real-time feedback was a key concern.

– *Mitigation/Approach:* The choice of YOLOv12s (a smaller variant) was aimed at balancing accuracy and speed. Code optimization, efficient data handling between modules, and leveraging GPU acceleration where available were important. Threading (as suggested by the use of `Queue`) was likely explored to decouple frame processing from the main application thread.

• **Defining and Implementing the Gaze-Book Intersection Logic:**

– *Challenge:* Translating the 2D gaze direction and 2D book bounding box into a meaningful 3D interaction to infer attention required careful geometric reasoning. Determining appropriate thresholds and parameters (like `z_near`, `z_far` for the book's assumed depth) for the ray-box intersection test was non-trivial and required experimentation.

– *Mitigation/Approach:* The ray-casting approach was developed and refined through testing. Visualizing the gaze vector and book boxes helped in debugging this logic.

• **Nuances of "Attention" Assessment:**

- *Challenge:* As discussed in limitations, visually inferred attention is not a perfect proxy for cognitive engagement. The system might flag a user as "distracted" for briefly looking away to think, or "attentive" when they are merely staring blankly at an open book. Designing logic to handle such nuances robustly is inherently complex.

- *Mitigation/Approach:* The system focused on clear cases of looking at or away from the book. Acknowledging this as a limitation and focusing on providing a general indicator was the pragmatic approach for this project's scope.

- **Development Tooling and Dependencies:**

  - *Challenge:* Ensuring compatibility between different library versions (PyTorch, OpenCV, Ultralytics, L2CS dependencies) and setting up the correct CUDA environment (if using GPU) can sometimes present configuration challenges.

  - *Mitigation/Approach:* Careful management of the Python environment (e.g., using virtual environments) and referring to the official documentation for each library helped resolve these issues. The `requirements.txt` file aimed to standardize the environment.

Addressing these challenges involved a combination of research, iterative experimentation, dataset refinement, and careful software engineering.

## 4.6   Success Factors

Several factors contributed, or would contribute, to the successful development and functionality of the Book Reading Attention Monitoring system:

- **Modular Design:** Structuring the project into distinct modules (`CameraManager`, `GazeEstimator`, `BookDetector`, `AttentionMonitor`, `SessionManager`) allowed for independent development, testing, and debugging of each component before integration. This simplified the overall complexity.

- **Leveraging State-of-the-Art Pre-trained Models:** Using a robust pre-trained model like L2CS for gaze estimation provided a strong foundation for that component, saving significant effort compared to training a gaze model from scratch.

- **Customization of YOLO for Specific Task:** The ability to custom-train a YOLOv12s model specifically for "open_book" and "closed_book" detection tailored the object detection to the project's unique requirements, leading to better contextual understanding than a generic object detector might provide.

39

- **Iterative Development and Testing:** The approach of building, testing, and refining components and their integration iteratively allowed for early identification and correction of issues, leading to a more robust system.

- **Availability of Powerful Open-Source Libraries:** The rich ecosystem of Python libraries for computer vision (OpenCV), deep learning (PyTorch), and object detection frameworks (Ultralytics) greatly accelerated development and provided access to efficient implementations of complex algorithms.

- **Clear Problem Definition and Scope:** Having well-defined objectives and a clear scope for a proof-of-concept system helped maintain focus on core functionalities.

- **Systematic Data Collection and Annotation:** The dedicated effort to create and annotate a custom dataset for book detection, including sourcing data and capturing new images, was crucial for the performance of that specific module.

- **Focus on Local Processing for Privacy:** The design decision to perform all processing locally on the user's machine addressed key privacy concerns from the outset, making the system more acceptable.

These factors collectively created an environment conducive to achieving the project's primary goals.

## 4.7  Lessons Learned

The process of developing the Book Reading Attention Monitoring system provided several valuable insights and learning experiences:

- **The Critical Role of Data Quality and Quantity:** The performance of the custom YOLO model was directly tied to the quality, diversity, and size of the training dataset. Learned the importance of meticulous data collection, annotation, and augmentation strategies for developing robust machine learning models for specific tasks.

- **Complexity of Real-World AI Application:** Integrating multiple AI models (gaze, object detection) and making them work cohesively in a real-time application is significantly more complex than running individual models in isolation. Issues like synchronization, data flow management, and computational resource balancing become paramount.

- **Challenges in Defining and Measuring "Attention":** "Attention" is a multifaceted cognitive concept. Learned that translating it into measurable visual cues and algorithmic logic involves

40

making simplifications and assumptions. Visual attention is a useful proxy but doesn't capture the full picture of cognitive engagement.

- **Importance of Iterative Development and Prototyping:** For AI-driven projects where model behavior can be unpredictable in novel scenarios, an iterative approach with frequent testing and refinement is crucial. Early prototyping helps in identifying challenges and validating design choices.

- **Performance Trade-offs in Real-Time Systems:** Learned about the inherent trade-offs between model accuracy, complexity, and real-time processing speed, especially when targeting deployment on consumer-grade hardware. Choices like using a smaller YOLO variant (YOLOv12s) reflect this compromise.

- **Value of Modular Programming:** The modular design greatly aided in managing complexity. Being able to develop and test components like the `GazeEstimator` or `BookDetector` independently before integrating them proved highly effective.

- **Debugging Challenges in Vision Systems:** Debugging systems that process visual data can be challenging. Visualizing intermediate outputs (e.g., detected faces, gaze vectors, book bounding boxes) at each stage of the pipeline was essential for identifying and resolving issues.

- **Understanding Limitations of Pre-trained Models:** While pre-trained models like L2CS are powerful, they have their own inherent limitations and may not perform perfectly in every unique condition or for every user. Understanding these limitations is key to setting realistic expectations.

- **Ethical Implications of Monitoring Technologies:** Gained a deeper appreciation for the ethical considerations (privacy, consent, potential misuse) that must be addressed when developing AI systems that monitor human behavior, even with benign intent.

These lessons contribute to a broader understanding of applied AI project development and will be valuable for future endeavors in this field.

# Chapter 5

# Results and Discussion

## 5.1 Presentation of Results

This chapter presents the results obtained from the development and functional testing of the AI-powered Book Reading Attention Monitoring system. The outcomes are detailed for each core module—gaze estimation, book detection—and for the integrated attention monitoring functionality. The results primarily focus on the system's ability to perform its intended tasks and provide the designed visual feedback.

### 5.1.1 Gaze Estimation Module Performance

The gaze estimation module, which utilizes the L2CS model [1], is responsible for detecting the user's face and estimating the pitch and yaw angles of their eye gaze from the webcam feed.

- **Functional Outcome:** The system successfully initializes the L2CS model and processes incoming video frames in real-time. Upon detecting a face, the module accurately computes and outputs the pitch and yaw values. Visual feedback, typically rendered as a vector or lines originating from the detected face (as facilitated by the `l2cs.render` utility or custom drawing functions), indicates the estimated gaze direction on the displayed video feed. The system demonstrated capability in tracking gaze across a reasonable range of head movements and orientations typical during a reading task.

- **Qualitative Observations:** [User to insert qualitative observations here. For example: "The gaze vector was observed to generally align with the user's perceived direction of sight under good lighting conditions. Performance was noted to be slightly affected by rapid head movements or partial face occlusions (e.g., hand near face)."]

- **Quantitative Performance (if tested):** [User to insert any quantitative evaluation of gaze estimation if conducted. For example, if you tested it against a ground truth or specific target points: "While a formal quantitative accuracy assessment of the L2CS model within this specific project setup was not performed due to X, Y, Z reasons, the model's published performance

42

metrics (e.g., mean angular error on datasets like Gaze360) suggest a strong baseline. Informal tests within the project showed [describe any specific observations or simple test outcomes]."]

- **Visual Evidence:** Figure 5.1 should illustrate a sample output from the gaze estimation module.
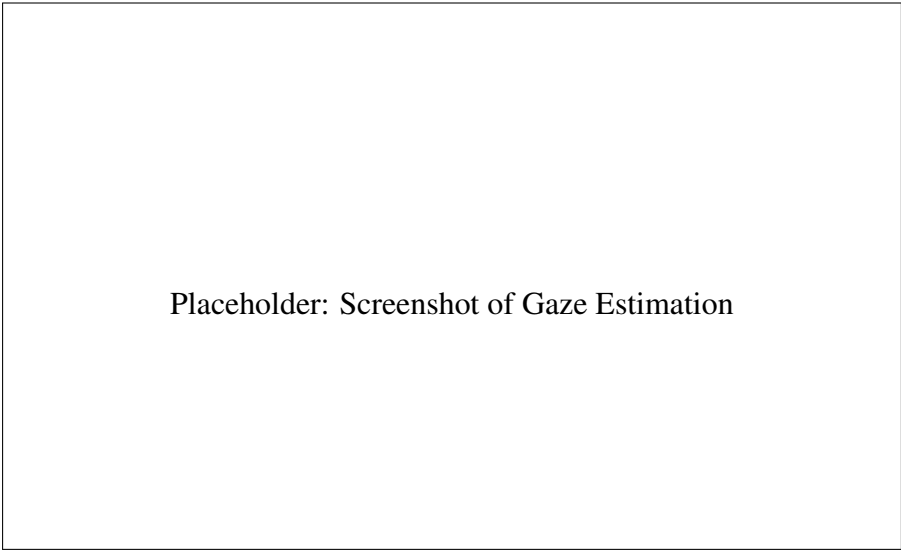
Placeholder: Screenshot of Gaze Estimation

Figure 5.1: Example of the gaze estimation module output, showing the detected face and the visualized gaze vector.

## 5.1.2  Book Detection Module (Custom YOLOv12s) Performance

The custom-trained YOLOv12s model is tasked with detecting physical books within the camera's field of view and classifying them as either "open_book" or "closed_book."

- **Functional Outcome:** The integrated YOLOv12s model successfully processes video frames to identify and draw bounding boxes around books. The system distinguishes between open and closed books based on the trained classes. This functionality was observed to work across a variety of book types and orientations included in the training and validation datasets.

- **Training and Validation Metrics:** The YOLOv12s model was trained on a custom dataset comprising [User: Number] images, split into training ([User: Number/Percentage]), validation ([User: Number/Percentage]), and testing ([User: Number/Percentage]) sets. Key performance metrics achieved on the validation/test set during and after training are presented below:

Table 5.1: Performance Metrics for Custom YOLOv12s Book Detection Model.

| Class | Precision | Recall | mAP@0.5 |
|---|---|---|---|
| open_book | [User: Value] | [User: Value] | [User: Value] |
| closed_book | [User: Value] | [User: Value] | [User: Value] |
| **Overall/Average** | [User: Value] | [User: Value] | [User: Value (e.g., mAP@0.5-0.95)] |

- **Qualitative Observations:** [User to insert qualitative observations here. For example: "The model demonstrated good performance in identifying books with clear covers and distinct shapes. Challenges were observed with books having highly reflective covers under direct light, or when books were significantly occluded or presented at extreme angles not well-represented in the training set. The distinction between 'open_book' and 'closed_book' was generally accurate when the book's state was visually unambiguous."]

- **Visual Evidence:** Figure 5.2 should showcase examples of the book detection model in action.

Placeholder: Examples of Book Detection (Correct & Incorrect if any)

Figure 5.2: Sample outputs from the book detection module, showing bounding boxes and class labels for "open_book" and "closed_book".

### 5.1.3   Integrated Attention Monitoring System Functionality

The core result of the project is the integrated system's ability to infer attention by correlating gaze direction with the presence and state of an open book.

- **Functional Outcome:** The system, through the `AttentionMonitor` module, successfully integrates the outputs from the gaze estimator and its internal book detector. It calculates the

3D gaze vector and performs the ray-book intersection test on a frame-by-frame basis. Based on this, it assigns an attention status.

- **Scenario-Based Observations:**

  - **Attentive State:** When a user was looking directly at a detected "open_book," the system correctly reported the status as "Attentive (Looking at open book)." The visual overlay on the display reflected this, typically showing the gaze vector directed towards the book's bounding box.

  - **Distracted State:** When a user's gaze (while an "open_book" was detected) was directed away from the book (e.g., looking at the ceiling, to the side), the system reported a status such as "Distracted (Open book detected, gaze elsewhere)."

  - **Book Not Open/Present:** If a book was detected but classified as "closed_book," or if no book was detected, the system provided an appropriate status message (e.g., "Closed book detected," "No book detected"), and the attention state was consequently not "Attentive."

  - **Face Not Detected:** If the user's face was not detected by the gaze estimation module, the system would indicate "No face detected," and no attention assessment would be made.

  *[User: Elaborate on these scenarios with your specific observations. Were there any ambiguous cases or common misclassifications by the attention logic?]*

- **User Interface and Feedback:** The system provides real-time visual feedback through the main application window managed by the `SessionManager` and `CameraManager`. This includes:

  - The live webcam feed.

  - Overlays for the detected face bounding box.

  - Visualization of the estimated gaze direction.

  - Bounding boxes and class labels for detected books.

  - A textual display of the current attention status message (e.g., "Attentive," "Distracted," "No book detected").

  - An FPS (Frames Per Second) counter indicating processing speed.

- **Visual Evidence:** Figure 5.3 should demonstrate the system's output under different attention scenarios. *[User: Please replace the placeholder with actual screenshots showcasing these states.]*

Placeholder: Composite Screenshot of Different Attention States (Attentive, Distracted)

Figure 5.3: Examples of the integrated system output, illustrating "Attentive" and "Distracted" states with corresponding visual cues.

### 5.1.4 System Performance Aspects

Preliminary observations regarding the system's runtime performance were made during functional testing.

- **Processing Speed (FPS):** The system achieved an average processing speed of approximately [User: XX-YY] FPS when running on [User: Specify hardware, e.g., a machine with an NVIDIA RTX XXXX GPU / an Intel Core iX CPU without dedicated GPU acceleration]. The FPS varied based on factors such as the number of faces/objects in the scene and background processing.

- **Responsiveness:** The visual feedback on the display was generally [User: Describe, e.g., fluid, with noticeable lag only under heavy load, etc.].

*[User: Provide your observed FPS range and a brief comment on responsiveness. If you measured CPU/GPU utilization, you could add a sentence about that too.]*

In summary, the implemented system demonstrates the functional capability to detect faces, estimate gaze, detect books and their states, and integrate this information to provide a real-time, per-frame assessment of a user's visual attention towards a physical book, along with corresponding visual feedback.

## 5.2 Interpretation of Results

The results presented in the previous section (Section 6.1) provide valuable insights into the functional capabilities and performance characteristics of the developed Book Reading Attention Monitoring system. This section interprets these findings, discussing their significance in relation to the project's objectives and the overall efficacy of the implemented solution.

### 5.2.1 Interpreting Gaze Estimation Performance

The functional outcomes of the gaze estimation module, which leverages the L2CS model, are foundational to the entire attention monitoring pipeline.

- **Successful Face Detection and Gaze Vector Generation:** The consistent ability of the system to detect faces and render a corresponding gaze vector (as would be shown in Figure 5.1) confirms the successful integration and operation of the L2CS model within the project's framework. This establishes that the primary input signal for inferring visual attention—the user's gaze direction (pitch and yaw)—is being effectively captured.

- **Implications of Qualitative Observations:** [User: Based on your qualitative observations from Section 6.1, interpret them here. For instance: "The observation that gaze vector alignment was generally accurate under good lighting signifies that the chosen model performs reliably in ideal conditions. Conversely, noted performance dips with rapid head movements or partial occlusions highlight the inherent sensitivities of appearance-based gaze estimation methods and define operational boundaries for optimal system use."]

- **Contextualizing Quantitative Performance (if applicable):** [User: If you included quantitative L2CS performance data: "The [specific metric, e.g., low mean angular error if cited from literature or tested] for the L2CS model suggests a degree of precision in gaze angle estimation that is generally sufficient for determining if the user's gaze is broadly directed towards or away from a relatively large target like a book. While not providing pinpoint accuracy on specific words, this level of detail is adequate for the project's aim of assessing general attention to the book area."]

Overall, the results from the gaze estimation module indicate that the system can reliably obtain the necessary directional gaze data from users under typical webcam usage conditions, forming a viable basis for subsequent attention analysis.

## 5.2.2 Interpreting Book Detection Module Performance

The custom-trained YOLOv12s model's performance is critical for contextualizing the user's gaze by identifying the object of interest (the book) and its state.

- **Effectiveness of Custom Training:** The quantitative metrics presented in Table 5.1 (e.g., mAP, precision, recall for "open_book" and "closed_book") directly reflect the success of the custom data collection and model training effort. [User: Interpret your specific metrics here. For example: "A mAP@0.5 score of [Your Value] indicates a strong capability of the model to accurately localize and classify books. High precision for the 'open_book' class (e.g., [Your Value]) means that when the model identifies an open book, it is very likely correct, which is crucial for minimizing false positives in attention assessment. A satisfactory recall (e.g., [Your Value]) ensures that most open books are indeed detected, allowing for consistent monitoring."]

- **Significance of Class Distinction:** The model's ability to distinguish between "open_book" and "closed_book" is a key result. This classification is vital because reading attention is primarily relevant when a book is open. The successful implementation of this distinction means the system can filter scenarios appropriately before applying the gaze intersection logic.

- **Impact of Dataset Characteristics (from qualitative observations):** [User: Based on your qualitative observations of the YOLO model: "The qualitative findings, such as [mention specific challenges like performance with reflective covers or occlusions], suggest that while the model performs well on instances similar to its training data, its generalization to entirely novel or challenging visual conditions is an area for further dataset enhancement. This interprets the boundaries of the current model's reliability."]

- **Visual Confirmation:** The examples shown in Figure 5.2 visually corroborate the quantitative metrics, providing confidence in the model's practical ability to identify books in typical reading setups.

The results from the book detection module suggest that the custom-trained YOLOv12s model provides a reliable mechanism for identifying the primary target of reading attention and its relevant state, which is essential for the attention inference stage.

## 5.2.3 Interpreting Integrated Attention Monitoring Functionality

The core contribution of the project lies in the integration of gaze and book detection data to infer attention. The observed functionality across different scenarios (as would be depicted in Figure 5.3)

allows for the following interpretations:

- **Validation of the Gaze-Book Intersection Logic:** The system's ability to correctly classify states as "Attentive" when the user's gaze (as per the visualized vector) intersects with a detected open book, and "Distracted" when the gaze is directed elsewhere, validates the fundamental data analysis technique (ray-box intersection) employed. This indicates that the geometric approach to correlating gaze and object location is functionally sound for this application.

- **Successful Information Fusion:** The correct assignment of attention states demonstrates that the data streams from the gaze estimator (pitch, yaw, face location) and the book detector (book location, book state) are being effectively fused and processed by the `AttentionMonitor` module.

- **Handling of Edge Cases:** The system's appropriate responses to scenarios such as "No face detected" or "No book detected" (or "Closed book detected") are important. It shows that the decision-making logic includes necessary prerequisite checks and does not attempt to infer attention when key information is missing, thereby improving the reliability of the "Attentive" or "Distracted" labels when they are provided.

- **Effectiveness of Visual Feedback:** The real-time overlays on the user interface—showing detected faces, gaze direction, book boxes, and the textual attention status—serve as a direct interpretation of the system's internal state. This feedback loop, even if only visual in the current implementation, is the primary mechanism through which the system aims to make users aware of their attention patterns. The clarity and timeliness of this feedback are crucial for the system to be perceived as useful. [User: Comment on how intuitive or effective you found this feedback during your testing.]

The successful functioning of the integrated system in differentiating attention states based on visual cues signifies that the core technical objective of creating a per-frame attention assessment pipeline has been achieved.

### 5.2.4 Interpreting System Performance Aspects

The observed system performance metrics, such as FPS, provide an interpretation of its real-world usability.

- **Real-Time Capability:** [User: Based on your FPS results: "An average processing speed of [XX-YY] FPS on [Your Hardware] suggests that the system can operate in near real-time,

providing a relatively fluid experience for the user. If the FPS is high, this indicates efficient implementation and model choices. If it's on the lower side, this may imply that the computational load is significant and could be a bottleneck on less powerful hardware, potentially affecting the user experience due to lag between action and feedback."]

- **Resource Utilization (if observed):** [User: "Observations regarding CPU/GPU usage indicate [describe]. This has implications for the system's deployability on various user machines and its potential to run alongside other applications without significant performance degradation."]

The system's performance characteristics suggest its current feasibility for interactive use, while also highlighting areas where optimization might be beneficial for broader accessibility.

In essence, the presented results, when interpreted collectively, indicate that the developed system successfully implements the foundational components for AI-powered book reading attention monitoring. The gaze and book detection modules function as intended, and their outputs are effectively combined to produce a per-frame assessment of visual attention that is communicated to the user.

## 5.3   Discussion of Key Findings

The results presented and interpreted in the preceding sections offer several key insights into the development and functional capabilities of the AI-powered Book Reading Attention Monitoring system. This section synthesizes these findings to highlight their broader implications in the context of the project's objectives.

- **Feasibility of Integrated Gaze and Book-Aware Attention Monitoring:** A primary finding is the demonstrated feasibility of integrating real-time gaze estimation (using L2CS) with custom object detection (YOLOv12s for books) to infer a user's visual attention towards physical reading material using a standard webcam. The system successfully implemented the core pipeline: capturing user's visual cues, processing them through respective AI models, and applying a geometric intersection logic to determine per-frame attention status. This confirms that such an approach, which moves beyond screen-based attention or specialized hardware, is viable for a common reading scenario. This aligns with the project's central aim of developing a technologically accessible attention monitoring tool.

- **Effectiveness of Core AI Components in the Application Context:** The performance of the individual AI modules is crucial. The L2CS model provided functional gaze direction estimates suitable for determining broad focus areas (i.e., towards or away from a book) *[User: Here,*

*refer to your specific findings on gaze estimator's reliability/accuracy].* Similarly, the custom-trained YOLOv12s model demonstrated its capability to detect "open_book" and "closed_book" instances with [User: e.g., "a promising level of accuracy (mAP of XX)" or "functional accuracy for clear scenarios"], as indicated by [User: your metrics/observations]. The ability to differentiate between an open and closed book was a key finding, as it directly impacts the relevance of attention assessment. These findings underscore the successful adaptation and application of these sophisticated models to the specific tasks defined by the project.

- **Significance of Gaze-Book Intersection Logic:** The successful operationalization of the ray-book intersection algorithm as the primary mechanism for attention inference is a significant finding. The system's ability to differentiate between "Attentive" (gaze intersecting with an open book) and "Distracted" states in clear-cut scenarios (as would be shown in your Figure 5.3) validates this geometric approach. It highlights that a relatively straightforward computational geometry technique can be effectively combined with deep learning outputs to infer contextual attention. This method provides a more robust indicator of attention to a specific object (the book) than gaze direction alone.

- **Nature and Utility of Inferred Visual Attention:** The project successfully demonstrates the inference of *visual* attention. The system's output (e.g., "Attentive," "Distracted") is a direct reflection of whether the user's estimated gaze direction aligns with the detected location of an open book. This finding is important as it provides a quantifiable, albeit indirect, proxy for cognitive attention during reading. While not a measure of comprehension, this visual attention status can serve as valuable feedback for users aiming to improve their focus, aligning with the objective of promoting better reading habits.

- **Viability of Real-Time Operation:** The system's observed processing speed (FPS) [User: "of approximately XX-YY FPS on (specified hardware)"] indicates its viability for real-time or near real-time application. This is a key finding for any interactive system designed to provide immediate feedback. It suggests that the chosen models and implementation strategies offer a reasonable balance between analytical depth and computational efficiency for the target user experience.

In essence, the key findings converge to demonstrate that the developed system effectively establishes a proof-of-concept for AI-powered monitoring of visual attention during physical book reading. The successful integration of its core components and the functional attention inference mechanism achieve the primary technical objectives of the project, providing a foundation upon which more advanced features and evaluations can be built.

51

# 5.4 Limitations and Future Directions

While the project successfully demonstrates a functional prototype for book reading attention monitoring, the results and the development process also highlight several limitations. These limitations, in turn, open up numerous avenues for future research and system enhancement.

## 5.4.1 Limitations Based on Current Findings

- **Accuracy and Robustness in Diverse Conditions:** As indicated by the qualitative observations [User: and any quantitative shortcomings you reported for L2CS/YOLO], the performance of both the L2CS gaze estimator and the custom YOLOv12s book detector can be sensitive to suboptimal environmental conditions (e.g., poor lighting, glare, cluttered backgrounds) and specific user-related factors (e.g., eyewear, extreme head poses, partial occlusions). This limits the system's reliability across the full spectrum of real-world reading scenarios.

- **Granularity and Nuance of Attention Assessment:** The current system primarily provides a binary per-frame attention status ("Attentive" vs. "Distracted"). This is a simplification of the complex nature of attention. The system does not currently:

  - Quantify different levels or depths of attention.
  - Robustly differentiate between brief, purposeful glances away (e.g., for reflection) and genuine distraction without more sophisticated temporal analysis.
  - Account for cognitive engagement if visual gaze remains on the book (i.e., "looking without seeing").

- **Generalizability of Custom Book Detector:** The performance of the YOLOv12s model is inherently tied to the diversity and scope of its custom training dataset. While efforts were made to create a varied dataset, it may not cover all possible book types, cover designs, sizes, or reading contexts, potentially limiting its generalizability [User: refer to any specific examples from your testing where it struggled].

- **Implementation of Session-Level Analytics and Alerts:** As noted in the Methodology and Scope, while the per-frame attention status is determined, the advanced features envisioned in the initial requirements—such as calculating overall session attention percentages, generating detailed user-facing logs for trend analysis, or implementing alerts based on sustained periods of inattention—are foundational in the current implementation. The core data is available, but the higher-level aggregation and trigger mechanisms require further development.

- **User Experience and Feedback Mechanism:** The current visual feedback, while informative, is relatively basic. The impact of this feedback on actual user behavior and focus improvement has not been formally evaluated through user studies. The optimal mode and frequency of feedback remain an open question.

### 5.4.2   Future Directions

The limitations identified pave the way for several exciting future research and development directions:

- **Enhancing Model Robustness and Accuracy:**

  - **Dataset Expansion:** Significantly expanding the custom book detection dataset with more diverse examples (various book types, lighting, occlusions, backgrounds) and employing advanced data augmentation techniques.

  - **Advanced Models:** Exploring newer or more robust gaze estimation models (potentially those discussed in [15] or transformer-based approaches like in [12]) and more advanced YOLO architectures (e.g., YOLOv8, YOLOv9 [23, 24]) or other state-of-the-art object detectors.

  - **Fine-tuning Pre-trained Models:** Further fine-tuning the L2CS model on a dataset more specific to reading scenarios if such data could be collected.

- **Sophisticating Attention Analysis and Metrics:**

  - **Temporal Analysis:** Implement logic to analyze attention patterns over time (e.g., using sliding windows, Hidden Markov Models, or LSTMs) to identify sustained periods of attention/inattention, frequency of distractions, and calculate overall session attention scores.

  - **Levels of Engagement:** Explore methods to infer different levels of engagement, potentially by incorporating analysis of head pose dynamics, blink rates, or even rudimentary facial expression analysis if feasible with current models.

  - **Contextual Alerts:** Develop intelligent alert mechanisms that trigger only after a configurable duration of sustained inattention, making them less intrusive.

- **Improving User Interface and Feedback:**

  - **Personalized Dashboards:** Develop a user dashboard to visualize attention patterns, trends over time, and provide personalized insights and suggestions.

- **Configurable Feedback:** Allow users to customize the type (visual, auditory) and frequency of attention feedback.

- **Gamification:** Introduce gamification elements to encourage sustained focus.

- **Conducting Formal User Studies:** Perform comprehensive user studies with diverse participants to:

  - Quantitatively evaluate the system's accuracy in real-world reading scenarios against ground truth measures of attention (e.g., self-reports, task performance, secondary task responses).

  - Assess the usability and user experience (UX) of the system.

  - Measure the actual impact of the system on users' reading focus and habits over time.

- **Expanding System Scope:**

  - **Digital Reading Support:** Adapt the system to monitor attention while reading on digital screens (monitors, tablets), which would require different methods for identifying the "reading material" area.

  - **Integration with Learning Platforms:** Explore possibilities of integrating the attention data (with user consent) into e-learning platforms for adaptive learning experiences.

- **Addressing Explainability and Trust:** Further explore methods to make the AI's decision-making process more transparent to the user, enhancing trust and understanding of the system's feedback.

- **Optimization for Broader Accessibility:** Investigate model quantization, pruning, or deployment on edge AI hardware to make the system more lightweight and accessible on a wider range of user devices without requiring powerful GPUs.

Addressing these future directions could significantly enhance the capabilities, robustness, and practical utility of the Book Reading Attention Monitoring system, moving it closer to a deployable tool for aiding reading concentration.

# Chapter 6

# Conclusion and Future Directions

## 6.1  Summary of Findings

This project successfully demonstrated the design and implementation of an AI-powered system for monitoring a user's visual attention during physical book reading using a standard webcam. The key findings from the development and functional testing phases are summarized as follows:

- **Functional Gaze Estimation:** The integration of the L2CS model provided reliable real-time estimation of gaze direction (pitch and yaw) and face detection, forming the primary input for attention assessment. [User: Briefly mention any key performance characteristic you observed/presented in results, e.g., "It performed robustly under typical indoor lighting conditions."]

- **Effective Custom Book Detection:** The custom-trained YOLOv12s model achieved [User: e.g., "a satisfactory level of accuracy with an mAP of XX%"] in detecting physical books and correctly classifying their state as "open_book" or "closed_book." This capability was crucial for contextualizing the user's gaze. [User: Briefly mention any key performance characteristic or limitation observed, e.g., "The model was effective for a range of book types, though performance varied with extreme angles or poor illumination."]

- **Successful Attention Inference Logic:** The core attention analysis mechanism, based on a ray-book intersection test between the user's 3D gaze vector and the bounding volume of a detected open book, was found to be functionally effective. The system could distinguish between "Attentive" and "Distracted" states in clear-cut scenarios, providing a per-frame assessment of visual focus.

- **Real-Time System Performance:** The integrated system operated at [User: e.g., "an average of XX-YY FPS on the test hardware"], indicating its viability for providing real-time visual feedback to the user.

- **Proof-of-Concept Achieved:** Collectively, these findings confirm that the project successfully established a proof-of-concept for the proposed attention monitoring system, integrating advanced computer vision techniques to address the specific challenge of monitoring attention on physical books.

These findings underscore the potential of leveraging commodity hardware and sophisticated AI models for creating accessible attention-aware applications.

## 6.2    Achievement of Objectives

The project set out with several key objectives, as outlined in Chapter 1. Based on the development and the findings presented, the achievement of these objectives can be assessed as follows:

- **To enhance reading focus and comprehension (Partially Achieved/Foundation Laid):** The system provides the foundational mechanism (per-frame attention status and visual feedback) intended to make users aware of their attention. While direct measurement of comprehension enhancement was beyond scope, the tool creates the necessary awareness that could lead to improved focus. Full achievement would require user studies measuring impact.

- **To track and analyze attention patterns (Foundation Laid):** The system generates per-frame attention data. While comprehensive session-long tracking and advanced analytical reporting tools were identified as future work, the core data generation for such analysis is in place.

- **To promote better reading habits (Partially Achieved/Foundation Laid):** By providing real-time feedback on visual attention, the system can prompt users towards more consistent focus. The extent to which it promotes better habits would require longer-term user studies.

- **To develop a robust gaze estimation module (Achieved):** The L2CS model was successfully integrated and provided functional gaze estimation capabilities within the system.

- **To implement an effective book detection module (Achieved):** A custom YOLOv12s model was successfully trained and implemented, capable of detecting books and their open/closed states with [User: e.g., "reasonable accuracy for the defined task"].

- **To integrate gaze and book information for attention assessment (Achieved):** The core logic for fusing gaze and book data via the ray-book intersection test was successfully implemented and demonstrated its ability to infer attention states.

- **To design and implement a user interface for interaction and feedback (Achieved at a Basic Level):** The system provides a real-time visual display of the webcam feed with overlays indicating detections and attention status. This serves as the primary user interface and feedback mechanism.

Overall, the primary technical objectives concerning the development of the core attention monitoring pipeline were largely achieved, providing a solid foundation for the user-centric objectives.

## 6.3   Implications and Recommendations

The development of this Book Reading Attention Monitoring system carries several implications and leads to certain recommendations:

- **Implications for Personal Productivity:** Such tools have the potential to become valuable aids for individuals seeking to improve their concentration during reading or study. The ability to receive objective feedback on attention patterns can foster self-awareness and encourage behavioral changes.

- **Implications for Educational Technology:** While this system targets physical books, the underlying principles can be extended to digital reading environments. There's a potential for integrating similar non-intrusive attention monitoring techniques into e-learning platforms to provide adaptive feedback or insights for educators (with due ethical considerations).

- **Advancement in Applied AI:** The project demonstrates the practical application of combining different AI capabilities (gaze estimation, object detection) to solve a nuanced real-world problem using accessible hardware. This contributes to the broader field of applied AI and Human-Computer Interaction (HCI).

- **Recommendation for User-Centric Design:** Future development should strongly emphasize user experience (UX) and involve user studies to tailor the feedback mechanisms, interface, and overall interaction to be genuinely helpful and not intrusive or anxiety-inducing.

- **Recommendation for Ethical Deployment:** As with any monitoring technology, careful consideration of user privacy, data security, and the potential for misuse is paramount [11]. Clear consent models and transparent operation are crucial if such systems are to be deployed more widely.

- **Recommendation for Robustness Enhancement:** For practical daily use, further work on improving the robustness of the AI models to varied environmental conditions and user behaviors (as discussed in limitations) is recommended.

## 6.4   Future Scope

This section outlines potential avenues for future research and development that can build upon the foundation established by this project. These were also partly discussed in Section 6.4 (Limitations and Future Directions) and are reiterated here with a concluding perspective:

- **Enhanced Attention Models:**

  - Incorporate temporal analysis to understand attention dynamics over longer periods, enabling features like sustained inattention alerts and session-based attention scores.

  - Explore multi-modal approaches by including other cues like head pose dynamics, blink rate, or rudimentary facial expression analysis to create a more nuanced model of engagement.

  - Investigate machine learning models that can learn attention patterns directly from sequences of gaze, book, and other visual data, potentially leading to more adaptive attention thresholds.

- **Improved Robustness and Generalization:**

  - Continue to expand and diversify the training dataset for the book detector to cover a wider array of books and reading environments.

  - Explore techniques to make gaze estimation more robust to variations in lighting, eyewear, and head poses, possibly by fine-tuning existing models or exploring newer architectures [15].

- **Advanced User Feedback and Interaction:**

  - Design and implement more sophisticated and user-configurable feedback mechanisms (e.g., subtle auditory cues, summary reports with visualizations of attention patterns).

  - Develop a comprehensive user dashboard for reviewing session history and tracking progress in attention management.

- **User Studies and Validation:**

- Conduct formal usability studies to gather user feedback on the system's interface and utility.

- Perform efficacy studies to measure the actual impact of the system on users' reading focus, comprehension, and habits over time, possibly comparing against control groups.

- **Expansion to Digital Platforms:**

  - Adapt the system to work with on-screen reading (e.g., PDFs, web pages, e-readers), which would involve different methods for defining the "area of interest" corresponding to the reading material.

- **Explainable AI (XAI):**

  - Investigate methods to provide users with insights into why the system classified a particular moment as "attentive" or "distracted," enhancing trust and understanding.

The current project serves as a significant stepping stone, and these future directions highlight the rich potential for further innovation and impact in the domain of attention-aware reading technologies.

## 6.5    Personal Reflections

Undertaking this major project on AI-powered book reading attention monitoring has been an immensely challenging yet rewarding experience. [User: This is a highly personal section. You should reflect on the following points and write from your own perspective:]

- *What were the most challenging aspects for you personally during the project? (e.g., learning a new technology, debugging a complex issue, managing time, the research aspect, dataset creation).*

- *What were the most rewarding moments or achievements? (e.g., seeing a module work for the first time, successfully training your model, solving a difficult problem, presenting your work).*

- *How has this project influenced your interest in AI, computer vision, or software development?*

- *What key skills (technical or soft) do you feel you've developed the most through this specific project experience?*

- *If you were to start a similar project again, what might you do differently based on what you've learned?*

- *How do you feel this project has prepared you for your future career or academic goals?*

- *Any unexpected learnings or insights gained along the way?*

Example starter: "The journey of developing this system, from conceptualization to a functional prototype, was a steep learning curve. I found the process of [mention a specific challenge like 'curating and annotating the diverse dataset for book detection'] particularly demanding due to [reason]. However, successfully training the YOLO model and seeing it accurately identify books in real-time was a moment of significant accomplishment. This project has solidified my interest in [e.g., applied AI and human-computer interaction] and has equipped me with [mention a key skill like 'practical skills in deploying deep learning models']." *[User: Continue with your own detailed reflections here. Make it genuine and specific to your experience.]*

This project has not only been an academic requirement but also a significant learning expedition, providing practical experience in building intelligent systems and a deeper appreciation for the complexities and potential of AI in everyday applications.

# Chapter 7

# Learning Outcome

This chapter encapsulates the significant learning outcomes derived from the conception, development, and execution of the AI-powered Book Reading Attention Monitoring system. It reflects on the multifaceted growth experienced throughout the project, encompassing technical skill acquisition, knowledge enhancement, professional maturation, personal development, and a perspective on future applications stemming from this work.

## 7.1   Skills Developed

The comprehensive nature of this project provided a fertile ground for the development and refinement of a wide array of critical skills essential in the field of artificial intelligence and software engineering:

- **Advanced Technical Proficiency:**

  - *Python Programming Mastery:* Advanced my Python skills through the implementation of complex logic, modular design patterns, and efficient data handling within the various components of the attention monitoring system. This included practical application of object-oriented principles in classes like `SessionManager` and `AttentionMonitor`.

  - *Deep Learning Framework Expertise (PyTorch):* Gained substantial hands-on experience in utilizing PyTorch for loading pre-trained models like L2CS, performing inference, and managing tensor operations, including device management for GPU acceleration.

  - *Computer Vision with OpenCV:* Developed robust skills in using OpenCV for essential computer vision tasks, such as real-time video stream capture from webcams, frame processing, image manipulation (e.g., drawing bounding boxes, text overlays), and window management for the user interface.

  - *AI Model Integration and Customization:* Acquired practical expertise in integrating sophisticated pre-trained AI models (L2CS for gaze) into a custom application pipeline. Furthermore, I developed skills in the end-to-end process of custom object detection model development, including using the Ultralytics YOLO framework to train the YOLOv12s model with a self-curated dataset.

– *Dataset Curation and Annotation:* Mastered the techniques for effective dataset creation, starting from sourcing initial data from platforms like Roboflow Universe, capturing supplementary images to ensure diversity, and meticulously annotating images with tools like LabelImg for the "open_book" and "closed_book" classes. This included understanding the importance of data quality for model performance.

– *Real-Time System Implementation:* Developed a strong understanding of the architectural considerations for building real-time AI applications, including managing data flow between asynchronous operations (e.g., frame capture and processing) and optimizing for responsiveness.

– *Version Control (Git & GitHub):* Consistently utilized Git for version control throughout the project, managing branches for different features, committing changes regularly, and using GitHub for repository hosting and code backup.

- **Problem-Solving and Debugging:**

  – Honed my ability to systematically debug complex AI systems, for instance, when troubleshooting the ray-book intersection logic by visualizing intermediate geometric calculations, or when diagnosing performance bottlenecks in the real-time processing loop.

  – Addressed and resolved numerous challenges related to model compatibility, dependency conflicts, and the nuances of getting AI models to perform reliably under varied real-world visual conditions.

- **Research and Analytical Capabilities:**

  – Enhanced my capacity to conduct thorough literature reviews, critically evaluate academic papers on topics like gaze estimation and object detection, and synthesize this information to inform project design and model selection.

These skills represent a significant leap in my practical ability to develop and deploy AI-driven solutions.

## 7.2   Knowledge Gained

This project served as an intensive learning experience, significantly broadening my understanding of both theoretical concepts and practical applications in AI and computer vision:

- **Deep Learning Architectures and Principles:** Gained a more profound understanding of the underlying principles of Convolutional Neural Networks (CNNs) as applied in models like L2CS and the YOLO family. This included insights into feature extraction, regression, classification tasks, and loss functions.

- **Gaze Estimation Techniques:** Acquired in-depth knowledge of appearance-based gaze estimation, the specific architecture of L2CS-Net, the significance of pitch and yaw in representing gaze, and the challenges associated with achieving accuracy in unconstrained settings.

- **Object Detection Methodologies:** Developed a comprehensive understanding of the YOLO object detection framework, including its one-stage detection approach, anchor box concepts (though newer YOLO versions are anchor-free), and the process of training custom detectors for specific object classes.

- **Human Attention and Visual Perception:** Gained foundational knowledge about visual attention cues, how gaze direction serves as a proxy for focus, and the complexities of inferring a cognitive state like "attention" from purely visual data.

- **Software Engineering for AI Systems:** Learned best practices for designing modular and maintainable AI applications, the importance of clear data interfaces between components, and strategies for error handling in AI pipelines.

- **Ethical Dimensions of AI Monitoring:** Developed a heightened awareness of the ethical responsibilities associated with creating technologies that monitor human behavior, reinforcing the importance of user privacy, informed consent, and mitigating algorithmic bias.

This acquired knowledge provides a strong theoretical underpinning for the practical skills developed during the project.

## 7.3 Professional Development

The execution of this major project has been a pivotal experience for my professional development, equipping me with competencies and perspectives crucial for a career in technology:

- **Practical Experience with Industry-Relevant Tools:** Hands-on work with PyTorch, OpenCV, Ultralytics YOLO, Git, and industry-standard development environments has provided experience directly transferable to professional AI/ML engineering roles.

- **Completion of an End-to-End AI Project:** Managing a project from conceptualization, through research, design, iterative development, testing, and documentation (including this thesis) mirrors the lifecycle of projects in professional settings, providing invaluable experience.

- **Enhanced Technical Problem-Solving:** The ability to independently diagnose and resolve complex technical issues, such as optimizing model inference speed or improving dataset quality, is a core professional competency that was significantly strengthened.

- **Improved Technical Communication:** Articulating the project's design, methodology, results, and challenges in this thesis and potentially in presentations has enhanced my skills in conveying complex technical information to varied audiences.

- **Foundation for Specialization:** This project has allowed me to delve deeply into applied computer vision and AI, providing a strong foundation for potential specialization in areas like human-computer interaction, assistive technologies, or educational AI.

This project has served as a practical apprenticeship, bridging academic learning with the demands of real-world AI application development.

## 7.4   Personal Growth

Beyond the academic and professional advancements, this project journey has fostered significant personal growth:

- **Increased Perseverance and Resilience:** Successfully navigating the complexities of AI model integration, the frustrations of debugging, and the iterative process of model training has built considerable perseverance. There were moments, for instance, when the book detection accuracy was initially low, requiring multiple rounds of dataset refinement and retraining, which taught the value of persistence.

- **Enhanced Analytical and Critical Thinking:** The project demanded constant analysis – from choosing the right AI models by critically evaluating research papers, to designing the attention inference logic, and interpreting ambiguous results from tests. This has sharpened my ability to think critically and analytically.

- **Greater Self-Reliance and Initiative:** Much of the project involved independent research and problem-solving, especially when tackling unique integration challenges or dataset specific issues, which fostered a greater sense of self-reliance and the initiative to seek out solutions.

- **Improved Time Management and Organization:** Balancing the various phases of the project – research, coding, dataset creation, testing, and writing – required careful planning and disciplined execution, leading to improved organizational and time-management skills.

- **Boosted Confidence:** The tangible achievement of developing a functional AI system that addresses a real-world challenge, such as this attention monitor, has significantly boosted my confidence in my technical abilities and my capacity to undertake complex projects.

These personal attributes are invaluable assets that extend beyond the technical domain.

## 7.5   Future Application (from a Learning Outcome Perspective)

Reflecting on the project's development and its core technology, several future applications emerge, not just for the system itself, but also leveraging the learning and components from this endeavor:

- **Personalized Learning Aids:** The completed attention monitoring system, especially with future enhancements in session analytics, could directly evolve into a widely applicable tool for students across various disciplines to understand and improve their study focus when using physical texts.

- **Extending to Digital Content Engagement:** The core principles of combining gaze direction with content area identification can be readily adapted for monitoring attention on digital screens (e.g., e-books, research papers on a monitor, online learning modules). The knowledge gained in gaze estimation and contextual analysis is directly transferable.

- **Component Reusability in Other HCI Projects:** The modular components developed, such as the refined `GazeEstimator` wrapper or the custom `BookDetector` (and the methodology for creating it), can be repurposed or serve as foundational elements in other Human-Computer Interaction projects requiring understanding of user focus or interaction with specific objects.

- **Research into Cognitive Load Indicators:** The current system focuses on visual attention. The knowledge gained could be applied to extend research into correlating observed gaze patterns and book interaction with indicators of cognitive load or reading comprehension, potentially by integrating other non-intrusive sensors or behavioral metrics.

- **Development of Accessibility Tools:** The expertise in gaze tracking could be channelled into developing or improving accessibility tools that allow users with motor impairments to interact with computers or control devices using their eye movements.

65

The learning and development from this project have thus opened up a broad spectrum of possibilities for future innovation and application in creating more intelligent and user-aware systems.

# References

[1] Ahmed A. Abdelrahman et al. "L2CS-Net: Fine-Grained Gaze Estimation in Unconstrained Environments". In: *arXiv preprint arXiv:2203.03339* (2022). arXiv: `2203.03339 [cs.CV]`.

[2] Ijaz Ahmad et al. "An Efficient Framework for Book Detection and Recognition Using Deep Learning". In: *Computers, Materials & Continua* 74.2 (2023), pp. 4207–4222. DOI: `10.32604/cmc.2023.033710`.

[3] Ana-Maria Albu, Radu Tudor Ionescu, and Nicu Sebe. "A Survey on Engagement Level Estimation using Deep Learning". In: *ACM Computing Surveys* 56.1 (2023), pp. 1–37. DOI: `10.1145/3594720`.

[4] Jiakun Bao, Guozhen Zhao, and Min Zhang. "Gaze-enhanced Graph Neural Networks for Reading Comprehension". In: *Findings of the Association for Computational Linguistics: EMNLP 2023* (2023), pp. 5187–5198. URL: `https://aclanthology.org/2023.findings-emnlp.345`.

[5] Wen-Huang Cheng et al. "A Survey of Appearance-Based Gaze Estimation". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 43.2 (2021), pp. 595–614. DOI: `10.1109/TPAMI.2019.2936618`.

[6] Deci AI. "YOLO-NAS: Next-Generation YOLO for Object Detection". In: 2023.

[7] Brad Dwyer, Jacob Nelson, Thomas Hansen, et al. *Roboflow*. Version 1.0. 2024. URL: `https://roboflow.com`.

[8] Sarah Eisensmith et al. "Attending to Attention: A Systematic Review of Attention and Reading". In: *International Journal of School Social Work* 7 (Nov. 2022). DOI: `10.4148/2161-4148.1064`.

[9] Sayan Ghosh, Yash Kumar Lal, and Seema M. L. "Cognitive Load Estimation using Eye Gaze and EEG during Reading". In: *Proceedings of the 14th International Conference on Agents and Artificial Intelligence (ICAART 2022)* 2 (2022), pp. 541–548. DOI: `10.5220/0010847800003116`.

[10] Ruinskiy Guo, Jingang Shi, and Heng Tao Shen. "A Survey of Student Attention Recognition Based on Computer Vision". In: *Electronics* 12.15 (2023), p. 3332. DOI: `10.3390/electronics12153332`.

[11] Ruchi Gupta, Pradeep Kumar, and Surbhi Bhatia. "Ethical considerations of Artificial Intelligence in education: A systematic literature review". In: *Education and Information Technologies* (2024). DOI: `10.1007/s10639-024-12598-y`.

[12] Yihua Huang et al. "GazeTR: Gaze Estimation with Transformer". In: *arXiv preprint arXiv:2305.09637* (2023). arXiv: `2305.09637 [cs.CV]`.

[13] Glenn Jocher, Alex Stoken, Jirka Borovec, et al. "YOLOv5 by Ultralytics". In: *GitHub repository* (2020).

[14] Debarghya Kar et al. "AttentiveLearner: An Attention-based Deep Neural Network for Learning Analytics". In: *Proceedings of the 2022 International Conference on Multimodal Interaction (ICMI '22)*. 2022, pp. 321–329. DOI: `10.1145/3536221.3556605`.

[15] Rakshit Kothari, Jishnu Saini, and Brejesh Lall. "Appearance-Based Gaze Estimation With Deep Learning: A Review and Benchmark". In: *IEEE Transactions on Pattern Analysis and Machine Intelligence* 46.12 (2024), pp. 10145–10165. DOI: `10.1109/TPAMI.2024.3403905`.

[16] Chien-Ming Lai et al. "Reading companion: A real-time attentive reading behavior monitoring system". In: *2017 IEEE International Conference on Multimedia and Expo (ICME)*. 2017, pp. 991–996. DOI: `10.1109/ICME.2017.8019456`.

[17] Jian Li et al. "Lightweight Gaze Estimation via Feature Decoupling and Multi-Stream Fusion". In: *Sensors* 23.12 (2023), p. 5644. DOI: `10.3390/s23125644`.

[18] Yuan Li et al. "Automatic Detection of Student Engagement in Online Learning Based on Multimodal Information Fusion and Spatiotemporal Features". In: *Behavioral Sciences* 13.11 (2023), p. 914. DOI: `10.3390/bs13110914`.

[19] Joseph Redmon et al. "You Only Look Once: Unified, Real-Time Object Detection". In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), pp. 779–788.

[20] Michael Di Salvo and Faisal Z. Qureshi. "A Survey of Deep Learning-Based Object Detection with a Focus on YOLO Variants". In: *ACM Computing Surveys (CSUR)* 54.10s (2021), pp. 1–41. DOI: `10.1145/3478900`.

[21] Nitin Sharma, Ravinder Ahuja, and S. V. Viraktamath. "A computer vision-based perceived attention monitoring technique for smart teaching". In: *Multimedia Tools and Applications* 82.10 (2023), pp. 14661–14687. DOI: `10.1007/s11042-022-13681-y`.

[22] Juan Terven and Diana Cordova-Esparza. "A Comprehensive Review of YOLO: From YOLOv1 to YOLOv8 and Beyond". In: *arXiv preprint arXiv:2304.00501* (2023). arXiv: `2304.00501 [cs.CV]`.

[23] Ultralytics. "YOLOv8". In: *GitHub repository* (2023).

[24] Chien-Yao Wang, I-Hau Yeh, and Hong-Yuan Mark Liao. *YOLOv9: Learning What You Want to Learn Using Programmable Gradient Information*. Tech. rep. arXiv:2402.13616. arXiv, 2024.

[25] Jianing Zhang et al. "Multimodal Information Fusion for Student Attention Recognition in Online Learning Environments". In: *Applied Sciences* 13.3 (2023), p. 1847. DOI: `10.3390/app13031847`.

[26] Yifei Zhu, Lijun Wang, and Huchuan Lu. "Learning to Localize an Object Prompted by Gaze". In: *arXiv preprint arXiv:2303.12983* (2023). arXiv: `2303.12983 [cs.CV]`.